

The Bayesian Confidence Interval for Coefficient of Variation of Zero-inflated Poisson Distribution with Application to Daily COVID-19 Deaths in Thailand

Sunisa Junnumtuam¹, Sa-Aat Niwitpong^{1*}, Suparat Niwitpong¹

¹ Department of Applied Statistics, Faculty of Applied Science, King Mongkut's University of Technology North Bangkok, Bangkok 10800, Thailand

Abstract

Coronavirus disease 2019 (COVID-19) has spread rapidly throughout the world and has caused millions of deaths. However, the number of daily COVID-19 deaths in Thailand has been low with most daily records showing zero deaths, thereby making them fit a Zero-Inflated Poisson (ZIP) distribution. Herein, confidence intervals for the Coefficient Of Variation (CV) of a ZIP distribution are derived using four methods: the standard bootstrap (SB), percentile bootstrap (PB), Markov Chain Monte Carlo (MCMC), and the Bayesian-based highest posterior density (HPD), for which using the variance of the CV is unnecessary. We applied the methods to both simulated data and data on the number of daily COVID-19 deaths in Thailand. Both sets of results show that the SB, MCMC, and HPD methods performed better than PB for most cases in terms of coverage probability and average length. Overall, the HPD method is recommended for constructing the confidence interval for the CV of a ZIP distribution.

Keywords:

Bootstrap;
Markov Chain Monte Carlo;
Highest Posterior Density.

Article History:

Received:	15	February	2021
Revised:	28	April	2021
Accepted:	12	May	2021
Published:	29	May	2021

1- Introduction

The outbreak of coronavirus disease 2019 (COVID-19) caused by severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) was first reported in Wuhan (Hubei, China) in December 2019 [1]. The disease has been prevalent in Thailand from 31 December 2019 until now, albeit at a low level of infection. COVID-19 can affect people in different ways; most people develop mild to moderate symptoms and recover without hospitalization whereas death occurs in acute cases. Even though the number of deaths around the world has crept into the millions, the number of daily COVID-19 deaths in Thailand has remained low throughout this period, with most daily records showing zero deaths.

Data on the number of daily COVID-19 deaths in Thailand from 3/12/2019 to 30/06/2020 comprising 176 days of daily observations in total were used in this study; there were 0 deaths on 147 days, 1 death on 13 days, 2 deaths on 6 days, 3 deaths on 7 days, and 4 deaths on 3 days. The frequency distribution of these presented in Figure 1 shows that they are clearly overdispersed with the variance exceeding the mean.

Since one of the properties of a Poisson distribution is equidispersion [2] (i.e., the mean is equal to the variance), it is not suitable for analyzing them. Some models that depart from standard count models, such as zero-inflated (ZI) and hurdle models, have been proposed to solve the overdispersion problem as both of them provide an alternative way to

* **CONTACT:** Sa-aat.n@sci.kmutnb.ac.th

DOI: <http://dx.doi.org/10.28991/esj-2021-SPER-05>

© 2021 by the authors. Licensee ESJ, Italy. This is an open access article under the terms and conditions of the Creative Commons Attribution (CC-BY) license (<https://creativecommons.org/licenses/by/4.0/>).

analyze count data with excess zeros. Lambert [3] proposed the ZI Poisson (ZIP) model where the probability of zeros is logit and the base count density is Poisson. The probability mass function (pmf) for this is given by:

$$f(x; \lambda, \omega) = \begin{cases} \omega + (1 - \omega)e^{-\lambda}; & x = 0 \\ \frac{(1 - \omega)e^{-\lambda}\lambda^x}{x!}; & x = 1, 2, 3, \dots \end{cases} \quad (1)$$

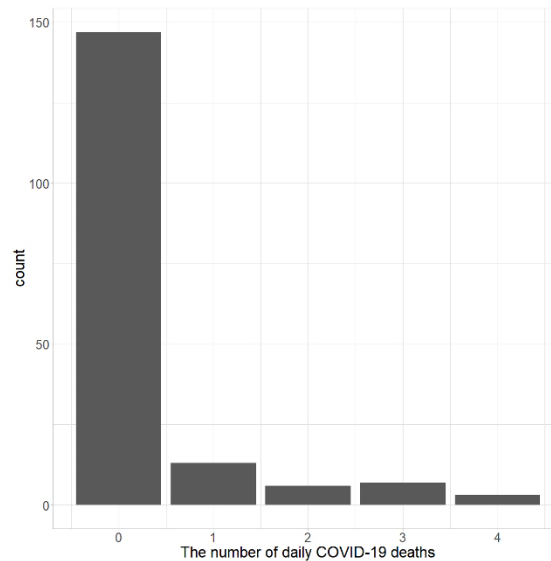


Figure 1. The frequency of the number of daily COVID-19 deaths in Thailand.

With the mean equal to $(1 - \omega)\lambda$ and the variance equal to $(1 - \omega)\lambda + (\omega/1 - \omega)((1 - \omega)\lambda)^2$; $\lambda > 0$, where ω indicates the proportion of zeros. There has been an abundance of research involving ZIP distributions. For instance, in a biomedical application, Bohning *et al.* [4] showed that the ZIP model can be used to analyze the decayed, missing and filled teeth index. When analyzing natural resources, Lee and Kim [5] recommended a ZIP model from a practical viewpoint for the number of torrential rainfall occurrences at the Daegu and Busan rain gauges in South Korea and compared it with the Poisson distribution, the generalized Poisson distribution (GPD), and the ZI generalized Poisson (ZIGP), and the Bayesian ZIGP model. In the field of insurance, Boucher *et al.* [6] adopted a ZIP distribution model for analyzing insurance data, while Kusuma and Purwono [7] showed that a ZIP regression model is more suitable than an ordinary Poisson regression model for modeling the frequency data of claims from the health insurance company PT.XYZ.

Previously, the Bayesian method has been extensively used to study the parameters of the ZIP distribution. For example, Rodrigues [8] studied the ZIP distribution by using the Bayesian method with noninformative priors to estimate the number of roots produced by 270 shoots of an apple cultivar *Trajan*; the ZIP distribution provided excellent fitting of the zero points. Xu *et al.* [9] investigated non-informative priors for a ZIP distribution with two ZIP model parameters and presented point estimations and 95% confidence intervals for the two parameters. Unhapipat *et al.* [10] applied Bayesian predictive inference with three types of prior distribution; the generalized noninformative prior, Jeffrey's noninformative prior, and beta-gamma prior for the ZIP distribution and illustrated their efficacies with real-life data on public health, natural catastrophes, and vehicle accidents. Furthermore, parameters of the ZIP distribution have been studied by applying other methods. For instance, Wagh and Kamalja [11] proposed a new probability estimator for the inflation parameter of ZIP distribution based on a moment estimator of the mean parameter. Srisuradetchai and Junnumtuam [12] studied a Wald-based confidence interval for the parameters in the Bernoulli component of ZIP and zero-altered Poisson (ZAP) models. Waguespack [13] used the likelihood and bootstrap approaches to construct confidence intervals for the mean of a ZIP distribution. Junnumtuam *et al.* [14] proposed confidence intervals for the mean of a ZIP distribution using the Markov chain Monte Carlo (MCMC) and highest posterior density (HPD) methods and applied them to analyze the number of new daily COVID-19 cases in Laos. Zou *et al.* [15] used the generalized fiducial inference to construct confidence intervals for the mean of ZIP and Poisson hurdle models.

Since the coefficient of variation (CV) is a relative measure of dispersion calculated as the standard deviation divided by the mean, it is unit invariant, and so is a useful statistic for comparing datasets with different units. Confidence intervals for the CV have been applied to various distributions. For example, Vangel [16] developed confidence intervals for the CVs of normal and modified McKay (Vangel) distributions. Panichkitkosolkul [17] introduced confidence intervals for the CV of a Poisson distribution using four different methods: Wald, Score, Wald Interval with Continuity Correction, and Variance Stabilizing. Meanwhile, [18] proposed a new asymptotic confidence interval constructed by

using a confidence interval for the Poisson mean for the CV of a Poisson distribution. Niwitpong [19] proposed a new confidence interval for the CV of a lognormal distribution. Yosboonruang *et al.* [20] studied confidence intervals for the CV of a delta-lognormal distribution constructed by using the generalized fiducial interval (GFI) and the method of variance estimates recovery (MOVER). The CV (θ) for a ZIP distribution can be expressed as:

$$\theta = \frac{\sigma}{\mu} = \frac{\sqrt{(1-\omega)\lambda + \frac{\omega}{1-\omega}((1-\omega)\lambda)^2}}{(1-\omega)\lambda} \quad (2)$$

With the sample estimate of θ being given by:

$$\hat{\theta} = \frac{\hat{\sigma}}{\hat{\mu}} = \frac{\sqrt{(1-\hat{\omega})\hat{\lambda} + \frac{\hat{\omega}}{1-\hat{\omega}}((1-\hat{\omega})\hat{\lambda})^2}}{(1-\hat{\omega})\hat{\lambda}} = \sqrt{\frac{1+\hat{\omega}\hat{\lambda}}{(1-\hat{\omega})\hat{\lambda}}} \quad (3)$$

Where $\hat{\omega}$ and $\hat{\lambda}$ are the estimators for ω and λ , respectively. We can see that the CV of a ZIP distribution is complex when formed with two parameters, thereby making it difficult to find its variance. To overcome this, we present four methods: the standard bootstrap (SB), percentile bootstrap (PB), MCMC, and HPD, for which finding the variance of the CV to construct the confidence interval is unnecessary. The efficiencies of the confidence intervals were compared via their coverage probabilities (CPs) and average lengths (ALs).

2- The Bootstrap-based Confidence Intervals

Let $x_i = (x_1, x_2, \dots, x_n)$ be a random sample from ZIP distribution, and let $\hat{\theta}$ represent the estimator of CV. The bootstrap procedure requires the following steps.

Algorithm 1

- Sample x_i^* with replacement from $\{x_1, \dots, x_n\}$ for $i \in \{1, \dots, n\}$.
- Calculate $\hat{\theta}^*$.
- Repeat 1 and 2 a total of B times to obtain the bootstrapped distribution of $\hat{\theta}$.

Efron and Tibshirani [21] indicated that a minimum of 1000 bootstrap samples are usually sufficient to compute reasonably accurate confidence interval estimates.

2-1- The SB Confidence Interval

From $B = 1000$ bootstrap estimates of $\hat{\theta}^*$, the sample average and standard deviation are calculated as:

$$\bar{\theta}^* = \frac{1}{1000} \sum_{i=1}^{1000} \hat{\theta}_b^*, \quad (4)$$

$$S_{\hat{\theta}^*}^* = \sqrt{\left(\frac{1}{999}\right) \sum_{i=1}^{1000} (\hat{\theta}_{(i)}^* - \bar{\theta}^*)^2}. \quad (5)$$

Algorithm 2

- Sample x_i^* with replacement from $\{x_1, \dots, x_n\}$ for $i \in \{1, \dots, n\}$.
- Calculate $\hat{\theta}^*$ using Equation 3.
- Repeat 1 and 2 a total of B times to obtain the bootstrap distribution of $\hat{\theta}$.
- Calculate the sample average and standard deviation using Equations 4 and 5, respectively.

Afterward, the $100(1 - \alpha)\%$ SB confidence interval for θ is calculated as follows:

$$CI_{SB} = \bar{\theta}^* \pm Z_{1-\frac{\alpha}{2}} S_{\hat{\theta}^*}^* \quad (6)$$

Where $Z_{1-\frac{\alpha}{2}}$ is obtained by using the $\left(1 - \frac{\alpha}{2}\right)^{\text{th}}$ quantile of the standard normal distribution.

2-2- The PB Confidence Interval

Algorithm 3

- Sample x_i^* with replacement from $\{x_1, \dots, x_n\}$ for $i \in \{1, \dots, n\}$.

- Calculate $\hat{\theta}^*$ using Equation 3.
- Repeat 1 and 2 a total of B times to obtain the bootstrap distribution of $\hat{\theta}$.
- Order $\hat{\theta}^*(i)$ and take the $100(\alpha/2)\%$ and the $100(1 - \alpha/2)\%$ points as the end points.

Afterward, the $100(1 - \alpha)\%$ PB confidence interval of θ is calculated as:

$$CI_{PB} = \left(\hat{\theta}_{B(\frac{\alpha}{2})}^*, \hat{\theta}_{B(1-\frac{\alpha}{2})}^* \right) \quad (7)$$

3- Bayesian Analysis of the ZIP Distribution

Suppose that $X = (X_1, \dots, X_n)$ is a vector of independent random variables generated from a ZIP distribution.

Let $A = x_i: x_i = 0, i = 1, \dots, n$ and m be the number of A , then the likelihood function [8] is given by:

$$L[\lambda, \omega] = [\omega + (1 - \omega)p(0|\lambda)]^m (1 - \omega)^{n-m} \prod_{x_i \notin A} p(x_i|\lambda) \quad (8)$$

Since the elements of set A can be generated from two different parts: (1) the real zeros part and (2) the Poisson distribution, then the unobserved latent allocation variable can be defined as:

$$I_i = \begin{cases} 1 & ; p(\lambda, \omega) \\ 0 & ; 1 - p(\lambda, \omega) \end{cases} \quad (9)$$

Where $i = 1, \dots, m$ and

$$p(\lambda, \omega) = \frac{\omega}{\omega + (1 - \omega)p(0|\lambda)}. \quad (10)$$

Thus, the likelihood function based on augmented data $D = \{X, I\}$, where $I = (I_1, \dots, I_m)$ [22], is in the form

$$L[\lambda, \omega|D] = L[\omega, \lambda] \prod_{i=1}^m p(\lambda, \omega)^{I_i} (1 - p(\lambda, \omega))^{1-I_i} = \omega^S (1 - \omega)^{n-S} p(0|\lambda)^{m-S} \prod_{x_i \notin A} p(x_i|\lambda), \quad (11)$$

Where $S = \sum_{i=1}^m I_i \sim \text{Bin}[m, p(\lambda, \omega)]$. Thus, the likelihood function based on the augmented data [8] is given by:

$$L[\lambda, \omega] \propto \omega^S (1 - \omega)^{n-S} \lambda^{\sum_{x_i \notin A} x_i} e^{-(n-S)\lambda} \quad (12)$$

and

$$p(\lambda, \omega) = \frac{\omega}{\omega + (1 - \omega)e^{-\lambda}}. \quad (13)$$

The likelihood function suggests the following independent priors:

$$\pi(\lambda) \sim \text{Gamma}[a, b] \quad (14)$$

$$\pi(\omega) \sim \text{Beta}[c, d]. \quad (15)$$

Hence, the joint posterior distribution for (λ, ω) given D becomes

$$\pi(\lambda, \omega|D) \propto \omega^{S+c-1} (1 - \omega)^{n-S+d-1} \lambda^{\sum_{x_i \notin A} x_i + a-1} e^{-(n-S+b)\lambda}. \quad (16)$$

Since ω and λ are independent given D , thus the marginal posterior distribution of ω is a Beta distribution; i.e.,

$$\pi(\omega|D) = \text{Beta}(S + c, n - S + d), \quad (17)$$

And the marginal posterior distribution of λ is:

$$\pi(\lambda|D) \propto \lambda^{\sum_{x_i \notin A} x_i + a-1} e^{-(n-S+b)\lambda}. \quad (18)$$

3-1- Posterior Simulation using the MCMC Algorithm: Gibbs Sampling

Algorithm 4

Given a, c , and $d = 0.5$ and $b = 0$ for the non-informative prior when $X \sim \text{ZIP}(\lambda^{(0)}, \omega^{(0)})$ and $t = 1, \dots, 10$, then;

- Calculate $p(\lambda^{(0)}, \omega^{(0)}) = \frac{\omega^{(0)}}{\omega^{(0)} + (1 - \omega^{(0)})e^{-\lambda^{(0)}}}$.
- Generate $S^{(t)}$ from $\text{Bin}(m, p(\lambda^{(t-1)}, \omega^{(t-1)}))$.
- Generate $\omega^{(t)}$ from $\text{Beta}(S^{(t)} + c, n - S^{(t)} + d)$.
- Generate $\lambda^{(t)}$ from $\text{Gamma}(\sum_{i=1}^n x_i + a, n - S^{(t)} + b)$.
- Repeat steps 2–4 t times to update the sample.
- Collect $\omega^{(t)}$ and $\lambda^{(t)}$ for 5,000 samples.
- Burn in 1,000 samples and calculate the estimator of θ using Equation 3.

Subsequently, the $100(1 - \alpha)\%$ approximately Bayesian confidence interval for θ is calculated as

$$CI_{MCMC} = (L.CI, U.CI), \quad (19)$$

Where $L.CI = \text{quantile}(\hat{\theta}, \alpha/2)$ and $U.CI = \text{quantile}(\hat{\theta}, 1 - \alpha/2)$.

3-2- The Bayesian-based HPD Interval

Chen and Shao [23] explained that credible intervals are easy to obtain either analytically or by using the MCMC method. The Bayesian credible interval or the HPD is the shortest interval containing $100(1 - \alpha)\%$ of the posterior probability such that the density within the interval has a higher probability than outside of it. This means that the HPD is more desirable when the marginal distribution is not symmetric. The two main properties of the HPD interval are as follows [24]:

- (a) The density for each point inside the interval is greater than that for each point outside the interval.
- (b) The HPD interval has the shortest length for a given probability (say $1 - \alpha$).

In this study, we used the MCMC method to estimate HPD intervals for the CV of a ZIP distribution. This approach only requires an MCMC sample generated from the marginal posterior distributions of the two parameters: λ and ω . In the simulation and computation, the HPD intervals were computed by using the *HDInterval* package version 0.2.0 from RStudio (<https://rstudio.com>).

Algorithm 5

Given that a, c , and $d = 0.5$ and $b = 0$ for the non-informative prior when $X \sim \text{ZIP}(\lambda^{(0)}, \omega^{(0)})$ and $t = 1, \dots, 10$, then;

- Calculate $p(\lambda^{(0)}, \omega^{(0)}) = \frac{\omega^{(0)}}{\omega^{(0)} + (1 - \omega^{(0)})e^{-\lambda^{(0)}}}$.
- Generate $S^{(t)}$ from $\text{Bin}(m, p(\lambda^{(t-1)}, \omega^{(t-1)}))$.
- Generate $\omega^{(t)}$ from $\text{Beta}(S^{(t)} + c, n - S^{(t)} + d)$.
- Generate $\lambda^{(t)}$ from $\text{Gamma}(\sum_{i=1}^n x_i + a, n - S^{(t)} + b)$.
- Repeat steps 2–4 t times to update the sample.
- Collect $\omega^{(t)}$ and $\lambda^{(t)}$ for 5,000 samples.
- Burn in 1,000 samples.
- Compute the HPD intervals $100(1 - \alpha)\%$ for θ .

4- Simulation Results

The simulation data were generated using the *gamlss.dist* package version 5.1-6 from RStudio. In the simulation study, sample size n was set as 30, 50, 100, or 200; $\lambda = 1, 5, 10, 15, 20$, or 25; and $\omega = 0.1(0.1)0.9$. The number of replications was set as 5,000 for SB and PB, and 1,000 for MCMC and HPD. The nominal confidence level was 0.95. In this study, the criterions to compare the efficiencies of the confidence intervals (CIs) are coverage probabilities (CPs) and average lengths (ALs). First, the confidence intervals were considered by the CPs. Since the nominal confidence level was 0.95, then the CIs which provided the CPs 0.95 or better are selected. After that, the ALs of these CIs are considered to find the shortest length to be the best CI. The simulation results in Table 1 are reported as the CPs and ALs of the confidence intervals for $n = 30$ and 100 (the results for 50 and 200 are not listed here).

For $n = 30$ and $\lambda = 1$, the SB, PB, and MCMC methods performed well for most cases and the HPD method performed well for all cases. When λ was increased to 5, 10, 15, 20, or 25, the performances of all of the methods dropped; this is especially true for PB, which obtained CPs lower than 0.95 in most cases. This is clearly evident in Figure 3, which presents the CPs for the CV of a ZIP distribution by using all four methods separated into 6 graphs according to the level of λ . Meanwhile, the ALs of the methods are shown in Figure 4.

SB provided CPs of approximately 0.95 for $n = 100$ and $\lambda = 1$, while the PB confidence interval produced CPs higher than 0.95 for only a few cases. However, MCMC and HPD performed well for most cases, with HPD attaining the shortest ALs. When $\lambda = 5$, the performances of PB, MCMC, and HPD dramatically plunged, especially that of the PB method with no cases providing CPs of more than 0.95. However, when λ was increased to 10, 15, 20, or 25, the performances of SB, MCMC, and HPD were better, with the HPD method still providing the shortest ALs. Comparisons of the CPs and ALs obtained by the four methods can be seen in Figures 5 and 6, respectively.

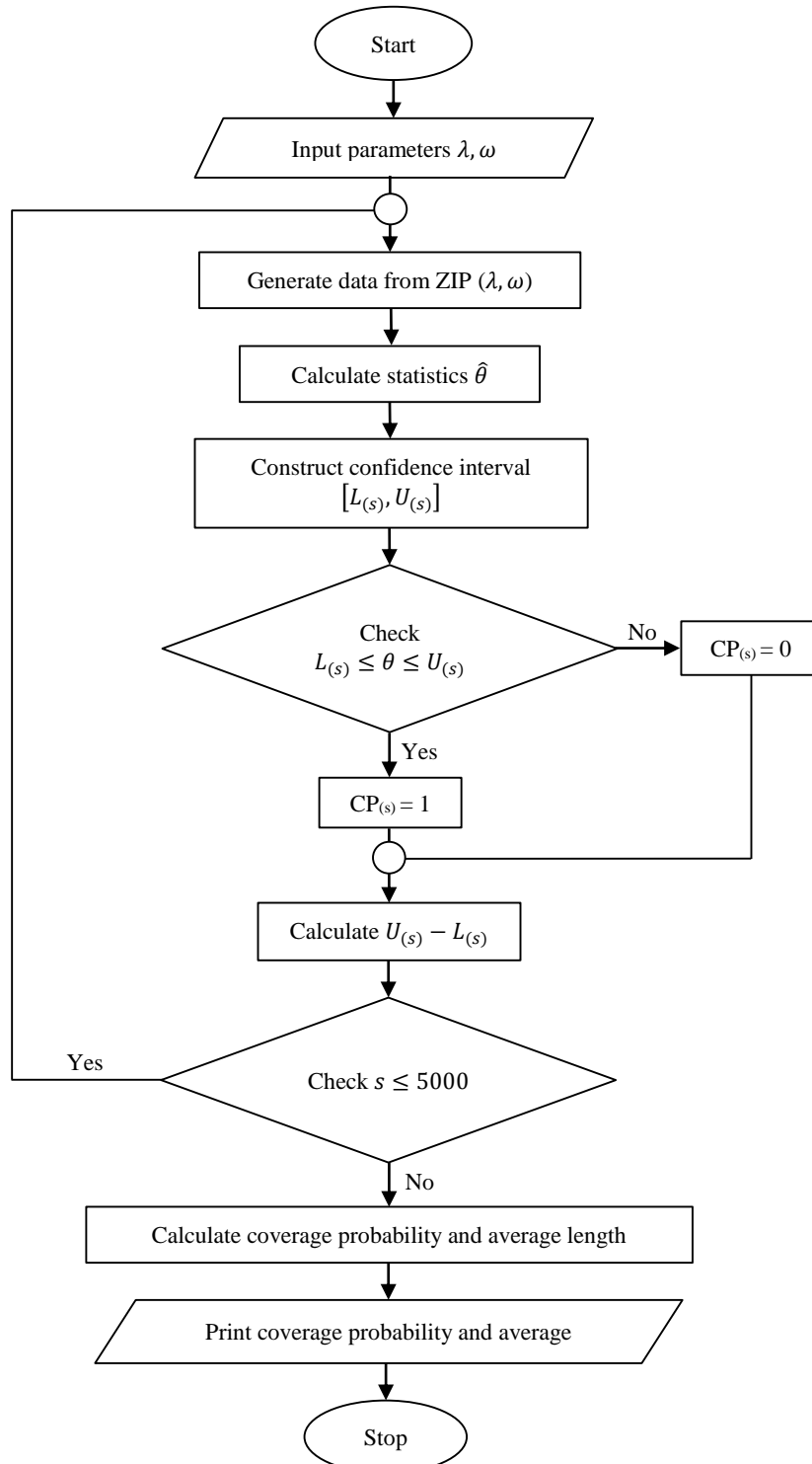


Figure 2. A flowchart of the simulation study.

Table 1. Performance metrics of the four methods for the 95% confidence intervals of the CV of a ZIP distribution.

n	λ	ω	k	Coverage probabilities (Average lengths)			
				SB	PB	MCMC	HPD
30	1	0.1	1.1055	0.9454	0.9466	0.9600	0.9730
				0.7195	0.7151	0.5856	0.5661
		0.2	1.2247	0.9518	0.9430	0.9540	0.9580
				0.8378	0.8314	0.7159	0.6910
		0.3	1.3628	0.9548	0.9514	0.9570	0.9620
				0.9975	0.9856	0.8627	0.8314
		0.4	1.5275	0.9562	0.9518	0.9570	0.9700
				1.2216	1.2025	1.0724	1.0296
		0.5	1.7321	0.9584	0.9524	0.9610	0.9710
				1.5621	1.5312	1.3592	1.2959
30	5	0.6	2.0000	0.9550	0.9442	0.9420	0.9550
				2.0527	2.0143	1.8475	1.7377
		0.7	2.3805	0.9508	0.9492	0.9550	0.9610
				2.7070	2.6138	2.7211	2.4906
		0.8	3.0000	0.9670	0.9520	0.9710	0.9670
				3.4482	3.0945	4.9293	4.2834
		0.9	4.3589	0.9606	0.9758	0.9560	0.9530
				3.9419	3.0752	10.1019	8.2707
30	10	0.1	0.5774	0.9210	0.9244	0.9370	0.9380
				0.3518	0.3503	0.2870	0.2746
		0.2	0.7071	0.9374	0.9360	0.9310	0.9380
				0.4472	0.4457	0.3941	0.3827
		0.3	0.8452	0.9490	0.9440	0.9490	0.9460
				0.5488	0.5467	0.4952	0.4829
		0.4	1.0000	0.9546	0.9460	0.9520	0.9530
				0.6753	0.6716	0.6205	0.6050
		0.5	1.1832	0.9646	0.9470	0.9440	0.9490
				0.8575	0.8497	0.7915	0.7699
30	5	0.6	1.4142	0.9660	0.9508	0.9520	0.9470
				1.1543	1.1366	1.0334	1.0008
		0.7	1.7321	0.9676	0.9472	0.9570	0.9610
				1.7031	1.6748	1.5024	1.4329
		0.8	2.2361	0.9584	0.9340	0.9480	0.9630
				2.6343	2.5552	2.5707	2.3715
		0.9	3.3166	0.9776	0.9854	0.9700	0.9530
				3.7193	3.1759	6.0620	5.1867
30	10	0.1	0.4714	0.9080	0.9160	0.9180	0.9320
				0.3381	0.3365	0.2993	0.2860
		0.2	0.6124	0.9366	0.9378	0.9380	0.9350
				0.4364	0.4352	0.4016	0.3909
		0.3	0.7559	0.9512	0.9438	0.9360	0.9420
				0.5327	0.5308	0.4956	0.4848
		0.4	0.9129	0.9580	0.9460	0.9590	0.9530
				0.6518	0.6474	0.6116	0.5981
		0.5	1.0954	0.9648	0.9450	0.9500	0.9370
				0.8244	0.8155	0.7619	0.7427
30	10	0.6	1.3229	0.9682	0.9488	0.9520	0.9420
				1.1110	1.0921	0.9959	0.9644
		0.7	1.6330	0.9670	0.9442	0.9450	0.9560
				1.6557	1.6214	1.4402	1.3750
		0.8	2.1213	0.9640	0.9368	0.9330	0.9550
				2.6070	2.5219	2.4750	2.2797
		0.9	3.1623	0.9748	0.9850	0.9720	0.9640
				3.7544	3.1944	5.8542	5.0215

30	15	0.1	0.4303	0.9060	0.9194	0.9300	0.9380
				0.3443	0.3423	0.3144	0.3004
		0.2	0.5774	0.9426	0.9362	0.9410	0.9390
				0.4410	0.4402	0.4153	0.4052
		0.3	0.7237	0.9538	0.9448	0.9420	0.9550
				0.5334	0.5311	0.5014	0.4914
		0.4	0.8819	0.9570	0.9466	0.9570	0.9420
				0.6490	0.6433	0.6084	0.5954
		0.5	1.0646	0.9660	0.9472	0.9490	0.9440
				0.8183	0.8067	0.7664	0.7474
		0.6	1.2910	0.9704	0.9494	0.9530	0.9480
				1.1020	1.0815	0.9838	0.9541
		0.7	1.5986	0.9680	0.9444	0.9500	0.9530
				1.6464	1.6093	1.3800	1.3217
		0.8	2.0817	0.9658	0.9398	0.9380	0.9600
				2.6058	2.5152	2.4471	2.2607
		0.9	3.1091	0.9754	0.9872	0.9770	0.9690
				3.7698	3.2049	5.6525	4.8553
30	20	0.1	0.4082	0.9098	0.9228	0.9300	0.9340
				0.3515	0.3491	0.3262	0.3124
		0.2	0.5590	0.9448	0.9358	0.9370	0.9360
				0.4459	0.4453	0.4191	0.4095
		0.3	0.7071	0.9584	0.9454	0.9420	0.9500
				0.5354	0.5326	0.5073	0.4975
		0.4	0.8660	0.9556	0.9472	0.9530	0.9350
				0.6489	0.6418	0.6119	0.5994
		0.5	1.0488	0.9660	0.9456	0.9540	0.9440
				0.8163	0.8026	0.7628	0.7442
		0.6	1.2748	0.9716	0.9486	0.9540	0.9520
				1.0985	1.0763	0.9979	0.9682
		0.7	1.5811	0.9690	0.9422	0.9480	0.9550
				1.6430	1.6040	1.4081	1.3463
		0.8	2.0616	0.9660	0.9404	0.9360	0.9600
				2.6066	2.5125	2.4042	2.2215
		0.9	3.0822	0.9754	0.9896	0.9750	0.9700
				3.7784	3.2102	5.6780	4.8721
30	25	0.1	0.3944	0.9178	0.9234	0.9320	0.9400
				0.3580	0.3552	0.3327	0.3190
		0.2	0.5477	0.9454	0.9340	0.9430	0.9460
				0.4500	0.4496	0.4288	0.4193
		0.3	0.6969	0.9610	0.9448	0.9510	0.9540
				0.5372	0.5341	0.5123	0.5026
		0.4	0.8563	0.9550	0.9466	0.9530	0.9320
				0.6492	0.6410	0.6196	0.6072
		0.5	1.0392	0.9658	0.9448	0.9510	0.9380
				0.8155	0.8001	0.7562	0.7381
		0.6	1.2649	0.9722	0.9466	0.9410	0.9390
				1.0968	1.0731	0.9952	0.9649
		0.7	1.5706	0.9698	0.9382	0.9470	0.9530
				1.6413	1.6008	1.3961	1.3350
		0.8	2.0494	0.9658	0.9402	0.9300	0.9670
				2.6075	2.5110	2.4566	2.2672
		0.9	3.0659	0.9758	0.9916	0.9690	0.9540
				3.7836	3.2135	5.7897	4.9572

100	1	0.1	1.1055	0.9424	0.9430	0.9770	0.9790
				0.3806	0.3790	0.3124	0.3071
		0.2	1.2247	0.9458	0.9456	0.9680	0.9660
				0.4362	0.4343	0.3835	0.3778
		0.3	1.3628	0.9460	0.9444	0.9580	0.9550
				0.5066	0.5042	0.4620	0.4556
		0.4	1.5275	0.9544	0.9470	0.9580	0.9580
				0.6000	0.5970	0.5523	0.5443
		0.5	1.7321	0.9540	0.9502	0.9530	0.9510
				0.7317	0.7272	0.6814	0.6707
		0.6	2.0000	0.9560	0.9504	0.9540	0.9510
				0.9373	0.9302	0.8739	0.8580
		0.7	2.3805	0.9626	0.9470	0.9550	0.9630
				1.3090	1.2937	1.2187	1.1902
		0.8	3.0000	0.9592	0.9512	0.9470	0.9610
				2.1455	2.1001	1.9203	1.8553
		0.9	4.3589	0.9494	0.9308	0.9310	0.9480
				4.6522	4.5566	4.7023	4.3162
100	5	0.1	0.5774	0.9396	0.9406	0.9430	0.9420
				0.1911	0.1903	0.1575	0.1547
		0.2	0.7071	0.9470	0.9426	0.9540	0.9500
				0.2392	0.2382	0.2141	0.2116
		0.3	0.8452	0.9476	0.9436	0.9320	0.9320
				0.2894	0.2883	0.2679	0.2651
		0.4	1.0000	0.9524	0.9480	0.9380	0.9410
				0.3498	0.3485	0.3295	0.3260
		0.5	1.1832	0.9544	0.9488	0.9420	0.9380
				0.4327	0.4311	0.4078	0.4033
		0.6	1.4142	0.9588	0.9454	0.9490	0.9440
				0.5555	0.5531	0.5273	0.5205
		0.7	1.7321	0.9584	0.9476	0.9350	0.9390
				0.7656	0.7609	0.7233	0.7125
		.8	2.2361	0.9684	0.9450	0.9410	0.9480
				1.2255	1.2119	1.1298	1.1054
		0.9	3.3166	0.9694	0.9422	0.9510	0.9500
				2.9347	2.8818	2.4366	2.3347
100	10	0.1	0.4714	0.9412	0.9420	0.9390	0.9310
				0.1839	0.1831	0.1675	0.1646
		0.2	0.6124	0.9506	0.9450	0.9560	0.9440
				0.2321	0.2312	0.2198	0.2175
		0.3	0.7559	0.9486	0.9462	0.9510	0.9540
				0.2793	0.2782	0.2693	0.2668
		0.4	0.9129	0.9510	0.9446	0.9510	0.9540
				0.3355	0.3343	0.3266	0.3234
		0.5	1.0954	0.9578	0.9482	0.9510	0.9550
				0.4132	0.4116	0.4003	0.3961
		0.6	1.3229	0.9574	0.9434	0.9470	0.9520
				0.5289	0.5267	0.5125	0.5063
		0.7	1.6330	0.9630	0.9478	0.9580	0.9580
				0.7267	0.7225	0.6973	0.6867
		0.8	2.1213	0.9708	0.9448	0.9380	0.9380
				1.1637	1.1510	1.0718	1.0491
		0.9	3.1623	0.9726	0.9358	0.9550	0.9540
				2.8396	2.7804	2.3995	2.2973

100	15	0.1	0.4303	0.9392	0.9414	0.9640	0.9610
				0.1877	0.1869	0.1783	0.1757
		0.2	0.5774	0.9504	0.9468	0.9510	0.9490
				0.2339	0.2331	0.2266	0.2245
		0.3	0.7237	0.9480	0.9470	0.9450	0.9420
				0.2790	0.2779	0.2722	0.2697
		0.4	0.8819	0.9506	0.9452	0.9520	0.9520
				0.3334	0.3322	0.3266	0.3236
		0.5	1.0646	0.9600	0.9484	0.9460	0.9540
				0.4091	0.4075	0.4005	0.3961
		0.6	1.2910	0.9562	0.9434	0.9390	0.9430
				0.5224	0.5203	0.5090	0.5031
		0.7	1.5986	0.9636	0.9474	0.9410	0.9410
				0.7166	0.7125	0.6919	0.6814
		0.8	2.0817	0.9726	0.9438	0.9530	0.9540
				1.1472	1.1341	1.0847	1.0608
		0.9	3.1091	0.9734	0.9368	0.9550	0.9550
				2.8175	2.7542	2.3798	2.2705
100	20	0.1	0.4082	0.9412	0.9406	0.9490	0.9360
				0.1919	0.1911	0.1849	0.1824
		0.2	0.5590	0.9524	0.9482	0.9520	0.9490
				0.2361	0.2352	0.2297	0.2276
		0.3	0.7071	0.9478	0.9422	0.9490	0.9480
				0.2797	0.2787	0.2741	0.2716
		0.4	0.8660	0.9498	0.9442	0.9520	0.9510
				0.3330	0.3318	0.3274	0.3245
		0.5	1.0488	0.9584	0.9464	0.9470	0.9530
				0.4076	0.4060	0.3987	0.3948
		0.6	1.2748	0.9558	0.9436	0.9540	0.9580
				0.5197	0.5174	0.5078	0.5017
		0.7	1.5811	0.9638	0.9474	0.9510	0.9610
				0.7121	0.7081	0.6901	0.6800
		0.8	2.0616	0.9724	0.9430	0.9450	0.9570
				1.1399	1.1264	1.0753	1.0525
		0.9	3.0822	0.9756	0.9356	0.9450	0.9350
				2.8081	2.7419	2.3505	2.2437
100	25	0.1	0.3944	0.9402	0.9394	0.9480	0.9360
				0.1955	0.1947	0.1879	0.1854
		0.	0.5477	0.9516	0.9478	0.9520	0.9520
				0.2378	0.2370	0.2331	0.2310
		0.3	0.6969	0.9474	0.9418	0.9430	0.9480
				0.2804	0.2794	0.2764	0.2740
		0.4	0.8563	0.9506	0.9440	0.9600	0.9560
				0.3329	0.3318	0.3276	0.3247
		0.5	1.0392	0.9588	0.9468	0.9400	0.9430
				0.4069	0.4053	0.3951	0.3912
		0.6	1.2649	0.9564	0.9432	0.9480	0.9480
				0.5183	0.5160	0.5068	0.5010
		0.7	1.5706	0.9636	0.9476	0.9510	0.9550
				0.7097	0.7056	0.6859	0.6761
		0.8	2.0494	0.9724	0.9432	0.9440	0.9450
				1.1357	1.1220	1.0660	1.0432
		0.9	3.0659	0.9766	0.9342	0.9590	0.9520
				2.8031	2.7350	2.3184	2.2180

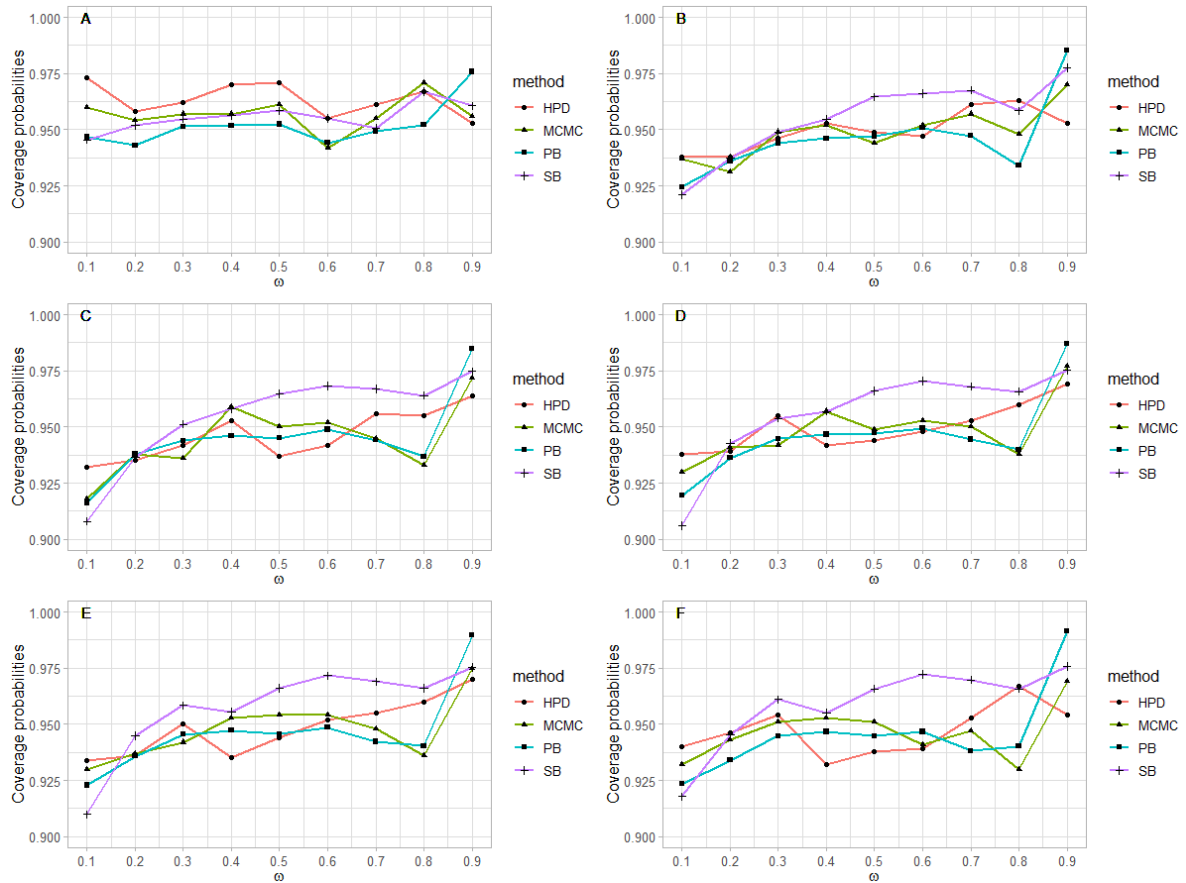


Figure 3. CP performances of the four methods for the 95% confidence intervals of the CV of a ZIP distribution for $n = 30$: (A) $\lambda = 1$, (B) $\lambda = 5$, (C) $\lambda = 10$, (D) $\lambda = 15$, (E) $\lambda = 20$, and (F) $\lambda = 25$.

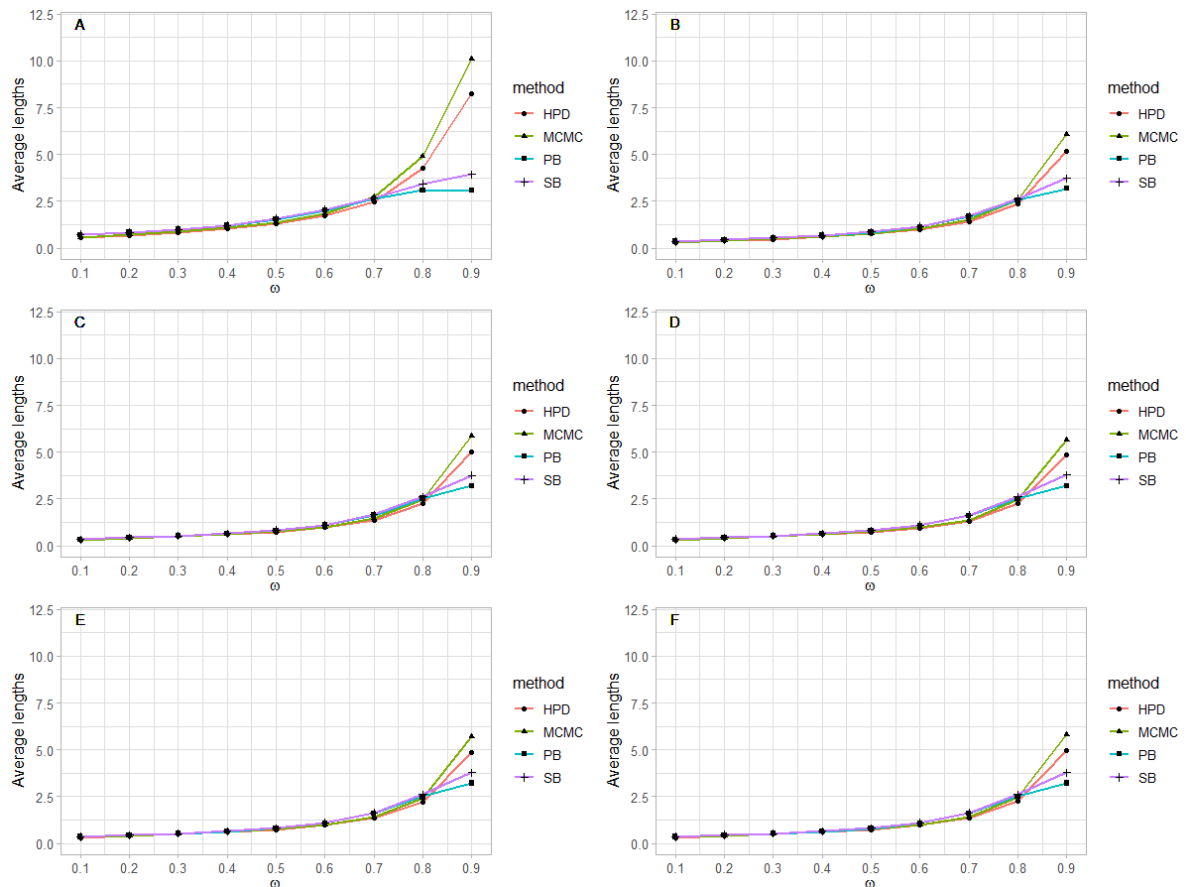


Figure 4. ALs performances of the four methods for the 95% confidence intervals of the CV of a ZIP distribution for $n = 30$: (A) $\lambda = 1$, (B) $\lambda = 5$, (C) $\lambda = 10$, (D) $\lambda = 15$, (E) $\lambda = 20$, and (F) $\lambda = 25$.

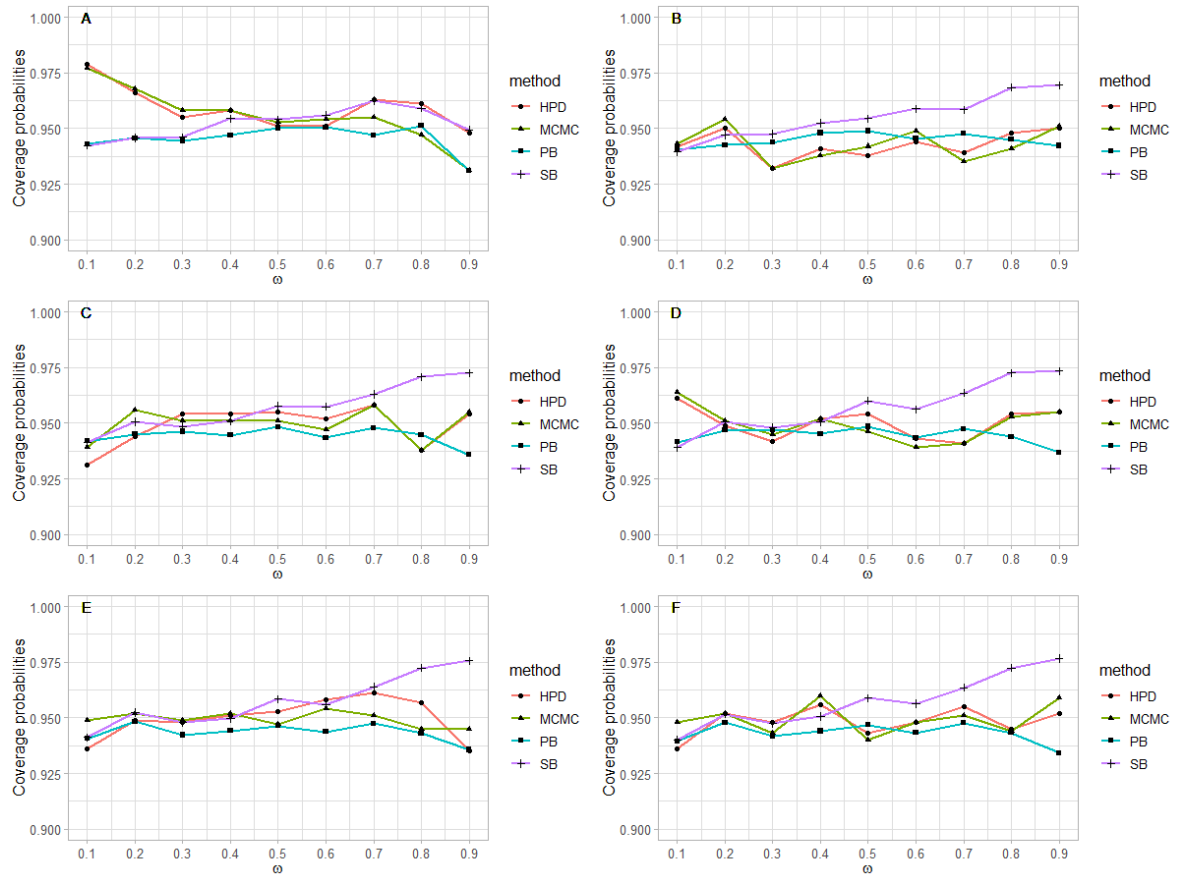


Figure 5. CP performances of the four methods for the 95% confidence intervals of the CV of a ZIP distribution for $n = 100$: (A) $\lambda = 1$, (B) $\lambda = 5$, (C) $\lambda = 10$, (D) $\lambda = 15$, (E) $\lambda = 20$, and (F) $\lambda = 25$.

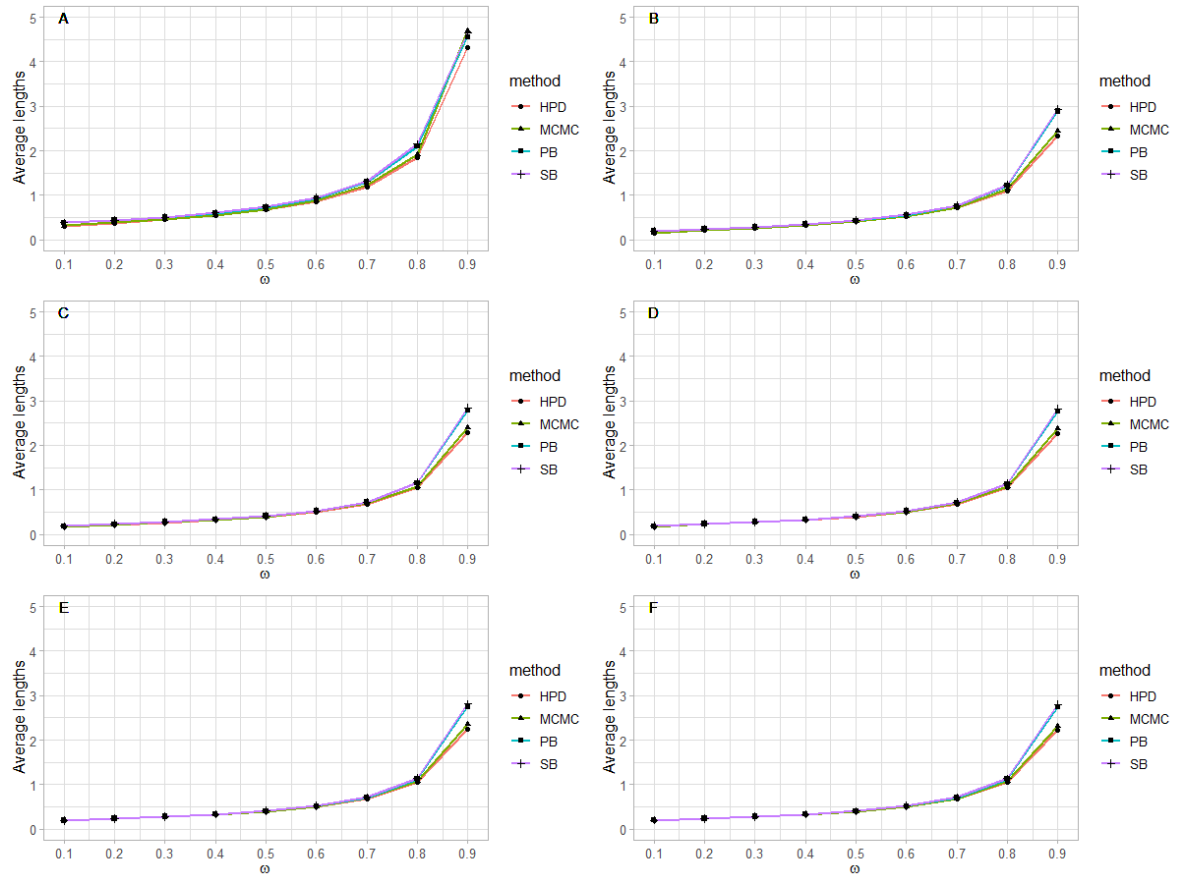


Figure 6. AL performances of the four methods for the 95% confidence intervals of the CV of a ZIP distribution for $n = 100$: (A) $\lambda = 1$, (B) $\lambda = 5$, (C) $\lambda = 10$, (D) $\lambda = 15$, (E) $\lambda = 20$, and (F) $\lambda = 25$.

5- Application of the Four Methods to the Daily COVID-19 Deaths Data for Thailand

Data of the daily covid-19 deaths in Thailand were retrieved using the R package 'utils' version 3.6.3 as a datasheet from the European Centre for Disease Prevention and Control (ECDC; <https://opendata.ecdc.europa.eu/covid19/casedistribution/csv>). The downloadable data file is updated daily and contains the latest available public data on COVID-19. Each row/entry contains the number of new cases reported per day and per country. Data from 3/12/2019 to 30/06/2020 for Thailand comprising 176 observations were used in the analysis. A histogram of the data is shown in Figure 1. First, whether the model for the data is appropriate was checked by comparing the Akaike information criterion (AIC) and the Bayesian information criterion (BIC) values for six distributions: ZIP, zero-inflated negative binomial (ZINB), Poisson, negative binomial (NB), geometric, and Gaussian. As reported in Table 2, the AIC and BIC values for ZIP are very similar (239.1898 and 245.5308) and the lowest recorded, thereby inferring that it provided the most efficient model.

Table 2. AIC and BIC values for six related models.

Models	ZIP	ZINB	Poisson	NB	Geometric	Gaussian
AIC	239.1898	241.1899	299.2355	244.5589	264.0685	448.6331
BIC	245.5308	250.7014	302.4060	250.8998	267.2390	454.9740

The mean is 0.3295 and the variance is 0.7365, thus it can clearly be seen that the variance exceeds the mean and there is overdispersion. Since the data contain many days of zero deaths, this makes the variance higher than the mean and provides a high CV of $\sqrt{0.7365}/0.3295 = 2.6045$. In this study, we used the *pscl* package to construct a ZIP model that provided count model coefficients for a Poisson distribution with log link = 0.466 ($\hat{\lambda} = e^{0.466} = 1.5936$) and ZI model coefficients for a binomial distribution with logit link = 1.344 ($\hat{\omega} = e^{1.344}/(1 + e^{1.344}) = 0.7931$). Hence, the estimator for the CV is $\hat{\theta} = \sqrt{1 + (0.7931)(1.5936)/(1 - 0.7931)(1.5936)} = 2.6203$. The 95% confidence intervals for the CV using the four methods cover the point estimator (Table 3). According to the simulation results for $n = 200$, $\lambda = 1$, and $\omega = 0.8$, the HPD method is the best for constructing the confidence interval for the CV because it provided a CP of more than 0.95 and the shortest AL.

Table 3. Estimation of the number of daily COVID-19 deaths in Thailand.

Method	CV Estimation	
	95% CI	Length of CI
SB	(2.0984 3.1879)	1.0895
PB	(2.1826 3.2990)	1.1164
MCMC	(2.1749 3.1820)	1.0072
HPD	(2.1215 3.1165)	0.9950

6- Discussion

In this study, four methods: SB, PB, MCMC, and HPD were applied to analyze and construct confidence intervals for the CV of a ZIP distribution. From the simulation results, it can be seen that the HPD method did not always provide CPs of more than 0.95 because the zeros can be from two sources: (1) real zeros from the Bernoulli component and (2) zeros in the Poisson distribution. This led to the sample from the ZIP distribution with a low λ having a higher proportion of zeros than one with a high λ , so the sample tended to have a greater number of zero items than the actual proportion of zeros (ω) in the distribution. In other words, the proportion of zeros in the sample becomes closer to ω as λ increases. Thus, the estimate of ω can lead to erroneous conclusions, especially when the sample is generated from a ZIP distribution with low λ and high ω . For example, when $n = 30$, $\lambda = 1$, and $\omega = 0.8-0.9$, the CPs of all of the methods were close to 0.95 but the ALs were too wide, especially for MCMC and HPD (see Figures 3 and 4). However, when the sample size was large, the ALs of the methods were similar for all cases. Moreover, we also conducted simulations with sample sizes $n = 50$ and 200 to investigate the trends in the CPs and ALs of the four methods. The simulation results for $n = 50$ were quite similar to those for $n = 30$ (i.e., a small sample size) and the simulation results for $n = 200$ were similar to those for $n = 100$ (i.e., a large sample size).

We can see from the results in Figures 3 and 5 that the methods involving bootstrapping performed well, especially SB, which provided CPs higher than 0.95 in all cases. However, the bootstrap method is not suitable for the ZI count data since its procedure requires a replacement sample. This makes the sample become inflated with zeros, especially when the proportion of zeros is high, which can also lead to erroneous estimation.

7- Conclusion

In this study, confidence intervals for the CV of a ZIP distribution were constructed by applying four methods, namely SB, PB, MCMC, and HPD. Since these methods do not require determining the variance of the CV of a ZIP distribution, which has a complex form with two parameters and is difficult to estimate, they are more convenient than the maximum likelihood estimation approach. When data are overdispersed with excess zeros, as is the case for the number of daily COVID-19 deaths in Thailand, the ZIP distribution is the best choice to estimate parameters for constructing confidence intervals for the CV. The results in Table 3 show that although the 95% confidence intervals for the CV using the four methods covered the true parameter, the HPD interval provided the shortest length. Similar to the simulation results, even though the CPs of four methods of the confidence intervals were close to the nominal confidence level, the HPD interval provided the shortest average length for almost all cases. Hence, we recommend the HPD method to construct the confidence interval of the CV of a ZIP distribution because it provided CPs of 0.95 or better and the shortest ALs in all of the scenarios investigated in this study.

8- Declarations

8-1-Author Contributions

Conceptualization, S.A.N. and S.N.; methodology, S.J. and S.A.N.; software, S.J.; validation, S.J., S.A.N., and S.N.; formal analysis, S.J. and S.A.N.; writing—original draft preparation, S.J.; writing—review and editing, S.A.N. and S.N.; supervision, S.N.; funding acquisition, S.A.N. All authors have read and agreed to the published version of the manuscript.

8-2-Data Availability Statement

The data presented in this study are openly available in the R package “utils” version 3.6.3 as a datasheet from the European Centre for Disease Prevention and Control (ECDC; <https://opendata.ecdc.europa.eu/covid19/casedistribution/csv>).

8-3-Funding

The first author received funding from Science Achievement Scholarship of Thailand. The second author received funding from King Mongkut's University of Technology North Bangkok. Grant number: KMUTNB-FF-65-22.

8-4-Conflicts of Interest

The authors declare that there is no conflict of interests regarding the publication of this manuscript. In addition, the ethical issues, including plagiarism, informed consent, misconduct, data fabrication and/or falsification, double publication and/or submission, and redundancies have been completely observed by the authors.

9- References

- [1] Valencia, Damian N. “Brief Review on COVID-19: The 2020 Pandemic Caused by SARS-CoV-2.” *Cureus* (March 24, 2020). doi:10.7759/cureus.7386.
- [2] Cameron, A. Colin, and Pravin K. Trivedi. “Regression analysis of count data.” Vol. 53. Cambridge University Press, (2013).
- [3] Lambert, Diane. “Zero-Inflated Poisson Regression, with an Application to Defects in Manufacturing.” *Technometrics* 34, no. 1 (February 1992): 1-14. doi:10.2307/1269547.
- [4] Böhning, D., E. Dietz, P. Schlattmann, L. Mendonça, and U. Kirchner. “The Zero-inflated Poisson Model and the Decayed, Missing and Filled Teeth Index in Dental Epidemiology.” *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 162, no. 2 (January 1999): 195–209. doi:10.1111/1467-985x.00130.
- [5] Lee, Cheol-Eung, and Sang Kim. “Applicability of Zero-Inflated Models to Fit the Torrential Rainfall Count Data with Extra Zeros in South Korea.” *Water* 9, no. 2 (February 16, 2017): 123. doi:10.3390/w9020123.
- [6] Boucher, Jean-Philippe, Michel Denuit, and Montserrat Guillen. “Number of Accidents or Number of Claims? An Approach with Zero-Inflated Poisson Models for Panel Data.” *Journal of Risk and Insurance* 76, no. 4 (December 2009): 821–846. doi:10.1111/j.1539-6975.2009.01321.x.
- [7] Kusuma, Rahmiani Dwinta, and Yogo Purwono. “Zero-Inflated Poisson Regression Analysis On Frequency Of Health Insurance Claim PT. XYZ.” *Proceedings of the 12th International Conference on Business and Management Research (ICBMR 2018)* (2019). doi:10.2991/icbmr-18.2019.52.
- [8] Rodrigues, Josemar. “Bayesian Analysis of Zero-Inflated Distributions.” *Communications in Statistics - Theory and Methods* 32, no. 2 (January 3, 2003): 281–289. doi:10.1081/sta-120018186.
- [9] Xu, Hai-yan, Min Xie, and Thong Ngee Goh. “Objective Bayes Analysis of Zero-Inflated Poisson Distribution with Application to Healthcare Data.” *IIE Transactions* 46, no. 8 (May 2014): 843–852. doi:10.1080/0740817x.2013.770190.

- [10] Unhapipat, Suntaree, Montip Tiensuwan, and Nabendu Pal. "Bayesian Predictive Inference for Zero-Inflated Poisson (ZIP) Distribution with Applications." *American Journal of Mathematical and Management Sciences* 37, no. 1 (October 20, 2017): 66–79. doi:10.1080/01966324.2017.1380545.
- [11] Wagh, Yogita S., and Kirtee K. Kamalja. "Zero-Inflated Models and Estimation in Zero-Inflated Poisson Distribution." *Communications in Statistics-Simulation and Computation* 47, no. 8 (August 4, 2017): 2248–2265. doi:10.1080/03610918.2017.1341526.
- [12] Srisuradetchai, P., and Junnumtuam, S. "Wald Confidence Intervals for the Parameter in a Bernoulli Component of Zero-Inflated Poisson and Zero-Altered Poisson Models with Different Link Functions" *Science & Technology Asia* 25, no.2 (2020) doi:10.14456/scitechasia.2020.16.
- [13] Waguespack, Dustin, K. Krishnamoorthy, and Meesook Lee. "Tests and Confidence Intervals for the Mean of a Zero-Inflated Poisson Distribution." *American Journal of Mathematical and Management Sciences* 39, no. 4 (June 22, 2020): 383–390. doi:10.1080/01966324.2020.1777914.
- [14] Junnumtuam, Sunisa, Sa-Aat Niwitpong, and Suparat Niwitpong. "The Bayesian Confidence Interval for the Mean of the Zero-Inflated Poisson Distribution." *Integrated Uncertainty in Knowledge Modelling and Decision Making. IUKM 2020. Lecture Notes in Computer Science* 12482. Springer, Cham. (November, 2020): 419–430. doi:10.1007/978-3-030-62509-2_35.
- [15] Zou, Y., Hannig, J., and Young, D.S. "Generalized fiducial inference on the mean of zero-inflated Poisson and Poisson hurdle models." *J Stat Distrib App* 8, no.5 (2021) doi: 10.1186/s40488-021-00117-0.
- [16] Vangel, Mark G. "Confidence Intervals for a Normal Coefficient of Variation." *The American Statistician* 50, no. 1 (February 1996): 21–26. doi:10.1080/00031305.1996.10473537.
- [17] Panichkitkosolkul, W. "A simulation comparison of new confidence intervals for the coefficient of variation of Poisson distribution." *Silpakorn University Science and Technology Journal* 4, no. 2 (2010):14-20. doi:10.14456/sustj.2010.7.
- [18] Panichkitkosolkul, W. "Asymptotic confidence interval for the coefficient of variation of Poisson distribution: a simulation study." *Maejo International Journal of Science and Technology* 4, no. 1 (2010): 1-7.
- [19] Niwitpong, Sa-aat. "Confidence Intervals for Coefficient of Variation of Lognormal Distribution with Restricted Parameter Space." *Applied Mathematical Sciences* 7 (2013): 3805–3810. doi:10.12988/ams.2013.35251.
- [20] Yosboonruang, Noppadon, Suparat Niwitpong, and Sa-Aat Niwitpong. "Confidence Intervals for Coefficient of Variation of Three Parameters Delta-Lognormal Distribution." *Studies in Computational Intelligence* (November 24, 2018): 352–363. doi:10.1007/978-3-030-04263-9_27.
- [21] Efron, B., and R. Tibshirani. "Bootstrap Methods for Standard Errors, Confidence Intervals, and Other Measures of Statistical Accuracy." *Statistical Science* 1, no. 1 (February 1, 1986): 54-75. doi:10.1214/ss/1177013815.
- [22] Tanner, Martin A., and Wing Hung Wong. "The calculation of posterior distributions by data augmentation." *Journal of the American statistical Association* 82, no. 398 (June, 1987): 528-540. doi:10.2307/2289457.
- [23] Chen, Ming-Hui, and Qi-Man Shao. "Monte Carlo Estimation of Bayesian Credible and HPD Intervals." *Journal of Computational and Graphical Statistics* 8, no. 1 (March 1999): 69-92. doi:10.2307/1390921.
- [24] Box, George E.P., and George C. Tiao. "Bayesian Inference in Statistical Analysis", New York, Wiley (April 6, 1992). doi:10.1002/9781118033197.