

Bayesian model selection method and its application in modeling Covid-19 pandemic development

Truong-Vinh Hoang

MATH4UQ, RWTH Aachen
July 23, 2020

- Presenting Bayesian approach for model selection
reference: Chapter 5, Ando, Tomohiro. *Bayesian model selection and statistical modeling*. CRC Press, 2010.
- Apply for modeling the Covid-19 evolution

1 Introduction

2 Theoretical background

- General framework
- Evaluation of marginal likelihood
 - Analytical method: Exact calculation of marginal likelihood
 - Laplace's method for computing the marginal likelihood
 - Bayesian information criterion (BIC)
- Expected predictive likelihood approach

3 Examples with Covid19 data

General framework

Posterior, marginal likelihood

Setting: Select a model from a set of candidate models: M_1, \dots, M_r

- Model M_k is characterized by probability density $f_k(\mathbf{x}|\boldsymbol{\theta}_k)$ where $\boldsymbol{\theta}_k$ is a p_k -dimensional vector of unknown parameter
- $\pi_k(\boldsymbol{\theta}_k)$ is the prior distribution for parameter vector $\boldsymbol{\theta}_k$
- Posterior probability of the model M_k given data set $\mathbf{X}_n = \{x_1, \dots, x_n\}$ of i.i.d samples

$$P(M_k|\mathbf{X}_n) = \frac{P(M_k) \int f_k(\mathbf{X}_n|\boldsymbol{\theta}_k) \pi_k(\boldsymbol{\theta}_k) d\boldsymbol{\theta}_k}{\sum_{\alpha=1}^r P(M_\alpha) \int f_\alpha(\mathbf{X}_n|\boldsymbol{\theta}_\alpha) \pi_\alpha(\boldsymbol{\theta}_\alpha) d\boldsymbol{\theta}_\alpha}$$

- $P(M_k)$ is the prior: e.g.
 - Uniform
 - Poisson prior: $P(M_k) \propto \lambda^{p_k} \exp(-\lambda)$, $\lambda > 0 \rightarrow$ higher probability on simpler models
 - $P(M_k) \propto \prod_{j=1}^p \pi_j^{\gamma_{k,j}} (1 - \pi_j)^{1-\gamma_{k,j}}$, where π_j are the prior probability that the j -th predictor is included in the model
- $\int f_k(\mathbf{X}_n|\boldsymbol{\theta}_k) \pi_k(\boldsymbol{\theta}_k) d\boldsymbol{\theta}_k$ is the **marginal likelihood**.

General framework

Bayes factor

Bayes factor (BF): comparing models, for testing hypotheses.

$$\text{BF} (M_k, M_j) \equiv \frac{P(\mathbf{X}_n | M_k)}{P(\mathbf{X}_n | M_j)}$$

Measures the evidence for M_k vs. M_j based on the data information
[posterior odds = Bayes factor \times the prior odds]

$$\frac{P(M_k | \mathbf{X}_n)}{P(M_j | \mathbf{X}_n)} = \frac{P(\mathbf{X}_n | M_k)}{P(\mathbf{X}_n | M_j)} \times \frac{P(M_k)}{P(M_j)}$$

Bayes factor can be reduced to classical likelihood ratio

1 Introduction

2 Theoretical background

- General framework
- Evaluation of marginal likelihood
 - Analytical method: Exact calculation of marginal likelihood
 - Laplace's method for computing the marginal likelihood
 - Bayesian information criterion (BIC)
- Expected predictive likelihood approach

3 Examples with Covid19 data

Methods to evaluate marginal likelihood $f(\mathbf{X}_n|\boldsymbol{\theta}_k, M_k)$

- Analytical method
 - E.g. Binomial model with conjugate prior
- Laplace's approximation
 - Bayesian information criterion (BIC)
 - Bayesian information criterion vs. Akaike's Information criterion
 - Generalized BIC
- Numerical method

Exact calculation of marginal likelihood

Binomial model with conjugate prior

- Observations: n independent samples
 $\mathbf{X}_n = \{x_1, x_2, \dots, x_n\} \sim \text{Bin}(n, p)$ with unknown p ,
- Conjugate prior: beta distribution with parameter α, β for p

Exact calculation of marginal likelihood

Binomial model with conjugate prior

- Observations: n independent samples
 $\mathbf{X}_n = \{x_1, x_2, \dots, x_n\} \sim \text{Bin}(n, p)$ with unknown p ,
- Conjugate prior: beta distribution with parameter α, β for p
- Let $y_n = \sum_{i=1}^n x_i$, the marginal likelihood:

$$\begin{aligned} f(\mathbf{X}_n | \alpha, \beta) &= \\ & \int_0^1 \left[\binom{n}{y_n} p^{y_n} (1-p)^{n-y_n} \times \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1} \right] dp \\ &= \binom{n}{y_n} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(y_n + \alpha) \Gamma(n + \beta - y_n)}{\Gamma(n + \alpha + \beta)} \end{aligned}$$

Exact calculation of marginal likelihood

Binomial model with conjugate prior

- Observations: n independent samples
 $\mathbf{X}_n = \{x_1, x_2, \dots, x_n\} \sim \text{Bin}(n, p)$ with unknown p ,
- Conjugate prior: beta distribution with parameter α, β for p
- Let $y_n = \sum_{i=1}^n x_i$, the marginal likelihood:

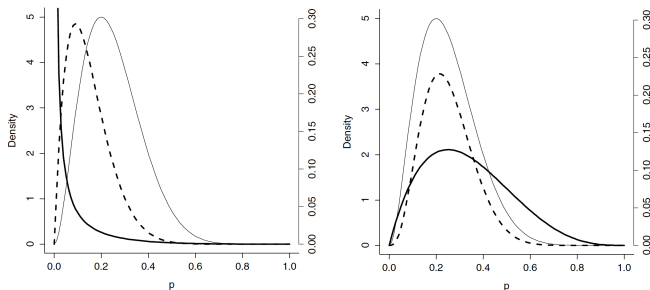
$$\begin{aligned} f(\mathbf{X}_n | \alpha, \beta) &= \\ & \int_0^1 \left[\binom{n}{y_n} p^{y_n} (1-p)^{n-y_n} \times \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1} \right] dp \\ &= \binom{n}{y_n} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(y_n + \alpha) \Gamma(n + \beta - y_n)}{\Gamma(n + \alpha + \beta)} \end{aligned}$$

- Example: Observation: $n = 10, y_n = 2$
- Model selection: $M_1 : p \sim \text{Beta}(0.1, 4); M_2 : p \sim \text{Beta}(2, 4); M_3 : p \sim \text{Beta}(4, 4); M_4 : p \sim \text{Beta}(8, 4)$

Exact calculation of marginal likelihood

Binomial model with conjugate prior

Models: $M_1 : p \sim \text{Beta}(0.1, 4)$; $M_2 : p \sim \text{Beta}(2, 4)$; $M_3 : p \sim \text{Beta}(4, 4)$; $M_4 : p \sim \text{Beta}(8, 4)$



left M_1 , right M_2

likelihood function (thin line), prior (thick line), posterior (dash line)

$f(\mathbf{X}_{10}|M_1) = 0.0277$, $f(\mathbf{X}_{10}|M_2) = 0.1648$, $\rightarrow M_2$ has the best BF

$f(\mathbf{X}_{10}|M_3) = 0.0848$, $f(\mathbf{X}_{10}|M_4) = 0.0168$.

Laplace's method

Asymptotic evaluation of the marginal likelihood

Assumption: the posterior density $\pi(\boldsymbol{\theta}|\mathbf{X}_n)$ is sufficiently well-behaved (highly peaked at the posterior mode $\hat{\boldsymbol{\theta}}_n$). Let

$$s_n(\boldsymbol{\theta}) = \log \{f(\mathbf{X}_n|\boldsymbol{\theta}) \pi(\boldsymbol{\theta})\}$$

Taylor series expansion of s_n yields (first order derivative at $\hat{\boldsymbol{\theta}}_n = 0$):

Laplace's method

Asymptotic evaluation of the marginal likelihood

Assumption: the posterior density $\pi(\boldsymbol{\theta}|\mathbf{X}_n)$ is sufficiently well-behaved (highly peaked at the posterior mode $\hat{\boldsymbol{\theta}}_n$). Let

$$s_n(\boldsymbol{\theta}) = \log \{f(\mathbf{X}_n|\boldsymbol{\theta}) \pi(\boldsymbol{\theta})\}$$

Taylor series expansion of s_n yields (first order derivative at $\hat{\boldsymbol{\theta}}_n = 0$):

$$s_n(\boldsymbol{\theta}) \approx s_n(\hat{\boldsymbol{\theta}}_n) - \frac{n}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_n)^T S_n(\hat{\boldsymbol{\theta}}_n) (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_n)$$

$$\text{where } S_n(\hat{\boldsymbol{\theta}}_n) = - \frac{1}{n} \frac{\partial^2 \log \{f(\mathbf{X}_n|\boldsymbol{\theta}) \pi(\boldsymbol{\theta})\}}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_n}$$

$$\begin{aligned} \exp \{s_n(\boldsymbol{\theta})\} &\approx \exp \{s_n(\hat{\boldsymbol{\theta}}_n)\} \times \exp \left\{ -\frac{n}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_n)^T S_n(\hat{\boldsymbol{\theta}}_n) (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_n) \right\} \\ &\propto \mathcal{N} \left(\hat{\boldsymbol{\theta}}_n, n^{-1} S_n(\hat{\boldsymbol{\theta}}_n)^{-1} \right) \end{aligned}$$

Laplace's method

Asymptotic evaluation of the marginal likelihood

Assumption: the posterior density $\pi(\boldsymbol{\theta}|\mathbf{X}_n)$ is sufficiently well-behaved (highly peaked at the posterior mode $\hat{\boldsymbol{\theta}}_n$). Let

$$s_n(\boldsymbol{\theta}) = \log \{f(\mathbf{X}_n|\boldsymbol{\theta}) \pi(\boldsymbol{\theta})\}$$

Taylor series expansion of s_n yields (first order derivative at $\hat{\boldsymbol{\theta}}_n = 0$):

$$s_n(\boldsymbol{\theta}) \approx s_n(\hat{\boldsymbol{\theta}}_n) - \frac{n}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_n)^T S_n(\hat{\boldsymbol{\theta}}_n) (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_n)$$

$$\text{where } S_n(\hat{\boldsymbol{\theta}}_n) = - \frac{1}{n} \frac{\partial^2 \log \{f(\mathbf{X}_n|\boldsymbol{\theta}) \pi(\boldsymbol{\theta})\}}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_n}$$

$$\begin{aligned} \exp \{s_n(\boldsymbol{\theta})\} &\approx \exp \{s_n(\hat{\boldsymbol{\theta}}_n)\} \times \exp \left\{ -\frac{n}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_n)^T S_n(\hat{\boldsymbol{\theta}}_n) (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_n) \right\} \\ &\propto \mathcal{N} \left(\hat{\boldsymbol{\theta}}_n, n^{-1} S_n(\hat{\boldsymbol{\theta}}_n)^{-1} \right) \end{aligned}$$

Laplace's method for computing the marginal likelihood

Marginal likelihood

- Marginal likelihood

$$\begin{aligned} P(\mathbf{X}_n|M) &= \int \exp\{s_n(\boldsymbol{\theta})\} d\boldsymbol{\theta} \\ &\approx \exp\left\{s_n\left(\hat{\boldsymbol{\theta}}_n\right)\right\} \times \int \exp\left\{-\frac{n}{2}\left(\boldsymbol{\theta}-\hat{\boldsymbol{\theta}}_n\right)^T S_n\left(\hat{\boldsymbol{\theta}}_n\right)\left(\boldsymbol{\theta}-\hat{\boldsymbol{\theta}}_n\right)\right\} d\boldsymbol{\theta} \\ &= f\left(\mathbf{X}_n|\hat{\boldsymbol{\theta}}_n\right) \pi\left(\hat{\boldsymbol{\theta}}_n\right) \times \frac{\left(2 \pi\right)^{\frac{p}{2}}}{n^{\frac{p}{2}}\left|S_n\left(\hat{\boldsymbol{\theta}}_n\right)\right|^{\frac{1}{2}}} \end{aligned}$$

where p is the dimension of vector $\boldsymbol{\theta}$

- $$\text{BF}\left(M_k, M_j\right) \approx \frac{f_k\left(\mathbf{X}_n|\hat{\boldsymbol{\theta}}_{kn}\right)}{f_j\left(\mathbf{X}_n|\hat{\boldsymbol{\theta}}_{jn}\right)} \times \frac{\pi_k\left(\hat{\boldsymbol{\theta}}_{kn}\right)}{\pi_j\left(\hat{\boldsymbol{\theta}}_{jn}\right)} \times \frac{\left|S_{jn}\left(\hat{\boldsymbol{\theta}}_{jn}\right)\right|^{\frac{1}{2}}}{\left|S_{kn}\left(\hat{\boldsymbol{\theta}}_{kn}\right)\right|^{\frac{1}{2}}} \times \left(\frac{2 \pi}{n}\right)^{\frac{p_k-p_j}{2}}$$

Laplace's method for computing the marginal likelihood

Marginal likelihood

- Marginal likelihood

$$\begin{aligned} P(\mathbf{X}_n|M) &= \int \exp\{s_n(\boldsymbol{\theta})\} d\boldsymbol{\theta} \\ &\approx \exp\left\{s_n\left(\hat{\boldsymbol{\theta}}_n\right)\right\} \times \int \exp\left\{-\frac{n}{2}\left(\boldsymbol{\theta}-\hat{\boldsymbol{\theta}}_n\right)^T S_n\left(\hat{\boldsymbol{\theta}}_n\right)\left(\boldsymbol{\theta}-\hat{\boldsymbol{\theta}}_n\right)\right\} d\boldsymbol{\theta} \\ &= f\left(\mathbf{X}_n|\hat{\boldsymbol{\theta}}_n\right) \pi\left(\hat{\boldsymbol{\theta}}_n\right) \times \frac{\left(2 \pi\right)^{\frac{p}{2}}}{n^{\frac{p}{2}}\left|S_n\left(\hat{\boldsymbol{\theta}}_n\right)\right|^{\frac{1}{2}}} \end{aligned}$$

where p is the dimension of vector $\boldsymbol{\theta}$

- $\text{BF}\left(M_k, M_j\right) \approx \frac{f_k\left(\mathbf{X}_n|\hat{\boldsymbol{\theta}}_{kn}\right)}{f_j\left(\mathbf{X}_n|\hat{\boldsymbol{\theta}}_{jn}\right)} \times \frac{\pi_k\left(\hat{\boldsymbol{\theta}}_{kn}\right)}{\pi_j\left(\hat{\boldsymbol{\theta}}_{jn}\right)} \times \frac{\left|S_{jn}\left(\hat{\boldsymbol{\theta}}_{jn}\right)\right|^{\frac{1}{2}}}{\left|S_{kn}\left(\hat{\boldsymbol{\theta}}_{kn}\right)\right|^{\frac{1}{2}}} \times \left(\frac{2 \pi}{n}\right)^{\frac{p_k-p_j}{2}}$
- Deterministic approaches: first two terms
- Third term: prefer model with less sensitivity
- Last term: complexity of the model (p_k)

Laplace's method for computing the marginal likelihood

Bayesian information criterion

We will consider two cases:

- $\log \pi(\theta) = O_p(1)$, $n \gg 1 \rightarrow$ prior information is ignored
- $\log \pi(\theta) = O_p(n)$

O_p : stochastic boundedness:

if $X_n = O_p(a_n)$ For any $\epsilon > 0$, there exists a finite $m > 0$ and a finite $N > 0$ such that, $P(|X_n/a_n| > m) < \epsilon$, $\forall n > N$.

Bayesian information criterion

Assumption: $\log \pi(\boldsymbol{\theta}) = O_p(1)$, $n \gg 1 \rightarrow$ prior information is ignored

The marginal likelihood is approximated as:

$$P(\mathbf{X}_n|M) \approx f(\mathbf{X}_n|\hat{\boldsymbol{\theta}}_{\text{MLE}}) \pi(\hat{\boldsymbol{\theta}}_{\text{MLE}}) \times \frac{(2\pi)^{\frac{p}{2}}}{n^{\frac{p}{2}} |J_n(\hat{\boldsymbol{\theta}}_{\text{MLE}})|^{\frac{1}{2}}}$$

J_n is minus the Hessian matrix of function $\frac{1}{n} \log f(\mathbf{X}_n|\hat{\boldsymbol{\theta}})$ at $\hat{\boldsymbol{\theta}}_{\text{MLE}}$

Bayesian information criterion

Assumption: $\log \pi(\boldsymbol{\theta}) = O_p(1)$, $n \gg 1 \rightarrow$ prior information is ignored

The marginal likelihood is approximated as:

$$P(\mathbf{X}_n|M) \approx f(\mathbf{X}_n|\hat{\boldsymbol{\theta}}_{\text{MLE}}) \pi(\hat{\boldsymbol{\theta}}_{\text{MLE}}) \times \frac{(2\pi)^{\frac{p}{2}}}{n^{\frac{p}{2}} |J_n(\hat{\boldsymbol{\theta}}_{\text{MLE}})|^{\frac{1}{2}}}$$

J_n is minus the Hessian matrix of function $\frac{1}{n} \log f(\mathbf{X}_n|\hat{\boldsymbol{\theta}})$ at $\hat{\boldsymbol{\theta}}_{\text{MLE}}$

Derive $-2 \log(P(\mathbf{X}_n|M))$, ignore constant terms when $n \gg 1 \rightarrow$ Schwarz's BIC

$$\text{BIC} = -2 \log f(\mathbf{X}_n|\hat{\boldsymbol{\theta}}_{\text{MLE}}) + p \log n$$

BIC gives a rough approximation to $\log(\text{BF})$

$$\log(\text{BF}(M_k, M_j)) = \log[P(\mathbf{X}_n|M_k)] - \log[P(\mathbf{X}_n|M_j)] \approx (\text{BIC}_j - \text{BIC}_k)/2$$

Bayesian information criterion

Assumption: $\log \pi(\boldsymbol{\theta}) = O_p(1)$, $n \gg 1 \rightarrow$ prior information is ignored

The marginal likelihood is approximated as:

$$P(\mathbf{X}_n|M) \approx f(\mathbf{X}_n|\hat{\boldsymbol{\theta}}_{\text{MLE}}) \pi(\hat{\boldsymbol{\theta}}_{\text{MLE}}) \times \frac{(2\pi)^{\frac{p}{2}}}{n^{\frac{p}{2}} \left| J_n(\hat{\boldsymbol{\theta}}_{\text{MLE}}) \right|^{\frac{1}{2}}}$$

J_n is minus the Hessian matrix of function $\frac{1}{n} \log f(\mathbf{X}_n|\hat{\boldsymbol{\theta}})$ at $\hat{\boldsymbol{\theta}}_{\text{MLE}}$

Derive $-2 \log(P(\mathbf{X}_n|M))$, ignore constant terms when $n \gg 1 \rightarrow$ Schwarz's BIC

$$\text{BIC} = -2 \log f(\mathbf{X}_n|\hat{\boldsymbol{\theta}}_{\text{MLE}}) + p \log n$$

BIC gives a rough approximation to $\log(\text{BF})$

$$\log(\text{BF}(M_k, M_j)) = \log[P(\mathbf{X}_n|M_k)] - [P(\mathbf{X}_n|M_j)] \approx (\text{BIC}_j - \text{BIC}_k)/2$$

$$\lim_{n \rightarrow \infty} \frac{(\text{BIC}_j - \text{BIC}_k)/2 - \log(\text{BF}(M_k, M_j))}{\log(\text{BF}(M_k, M_j))} = 0$$

BIC does not involve J_n

BIC vs Akaike information criterion (AIC)

Model selection by minimizing ICs: $IC_k = -2 \log f_k(\mathbf{X}_n | \hat{\boldsymbol{\theta}}_{\text{MLE } k}) + c_{n,p_k}$
 $c_{n,p}$: penalty to encourage the selection of a **parsimonious model** (fewer parameter dimension)

- BIC: $c_{n,p} = p \log(n)$, AIC: $c_{n,p} = 2p$

¹discussed later in Expected predictive likelihood approach

BIC vs Akaike information criterion (AIC)

Model selection by minimizing ICs: $IC_k = -2 \log f_k(\mathbf{X}_n | \hat{\boldsymbol{\theta}}_{\text{MLE } k}) + c_{n,p_k}$
 $c_{n,p}$: penalty to encourage the selection of a **parsimonious model** (fewer parameter dimension)

- BIC: $c_{n,p} = p \log(n)$, AIC: $c_{n,p} = 2p$

The best model $f_0(z|\boldsymbol{\theta}_0)$ that has the lowest Kullback-Leibler divergence¹ from true model $G(z)$, or equivalently maximum expected log-likelihood:

$$\int \log f_0(z|\boldsymbol{\theta}_0) dG(z) = \max_k \left\{ \sup_{\boldsymbol{\theta}_k} \int \log f_k(z|\boldsymbol{\theta}_k) dG(z) \right\}$$

¹discussed later in Expected predictive likelihood approach

BIC vs Akaike information criterion (AIC)

Model selection by minimizing ICs: $IC_k = -2 \log f_k(\mathbf{X}_n | \hat{\boldsymbol{\theta}}_{\text{MLE } k}) + c_{n,p_k}$
 $c_{n,p}$: penalty to encourage the selection of a **parsimonious model** (fewer parameter dimension)

- BIC: $c_{n,p} = p \log(n)$, AIC: $c_{n,p} = 2p$

The best model $f_0(z|\boldsymbol{\theta}_0)$ that has the lowest Kullback-Leibler divergence¹ from true model $G(z)$, or equivalently maximum expected log-likelihood:

$$\int \log f_0(z|\boldsymbol{\theta}_0) dG(z) = \max_k \left\{ \sup_{\boldsymbol{\theta}_k} \int \log f_k(z|\boldsymbol{\theta}_k) dG(z) \right\}$$

Consistency in picking model $\log f_0(z|\boldsymbol{\theta}_0)$ requires for all non-optimal f_k :

$$\liminf_{n \rightarrow \infty} \left[\frac{1}{n} \int \log f_0(\mathbf{X}_n | \boldsymbol{\theta}_0) dG(\mathbf{X}_n) - \frac{1}{n} \int \log f_k(\mathbf{X}_n | \boldsymbol{\theta}_{k0}) dG(\mathbf{X}_n) \right] > 0$$

¹discussed later in Expected predictive likelihood approach

BIC vs Akaike information criterion (AIC)

Model selection by minimizing ICs: $IC_k = -2 \log f_k(\mathbf{X}_n | \hat{\boldsymbol{\theta}}_{\text{MLE } k}) + c_{n,p_k}$
 $c_{n,p}$: penalty to encourage the selection of a **parsimonious model** (fewer parameter dimension)

- BIC: $c_{n,p} = p \log(n)$, AIC: $c_{n,p} = 2p$

The best model $f_0(z|\boldsymbol{\theta}_0)$ that has the lowest Kullback-Leibler divergence¹ from true model $G(z)$, or equivalently maximum expected log-likelihood:

$$\int \log f_0(z|\boldsymbol{\theta}_0) dG(z) = \max_k \left\{ \sup_{\boldsymbol{\theta}_k} \int \log f_k(z|\boldsymbol{\theta}_k) dG(z) \right\}$$

Consistency in picking model $\log f_0(z|\boldsymbol{\theta}_0)$ requires for all non-optimal f_k :

$$\liminf_{n \rightarrow \infty} \left[\frac{1}{n} \int \log f_0(\mathbf{X}_n | \boldsymbol{\theta}_0) dG(\mathbf{X}_n) - \frac{1}{n} \int \log f_k(\mathbf{X}_n | \boldsymbol{\theta}_{k0}) dG(\mathbf{X}_n) \right] > 0$$

if $f_0(z|\boldsymbol{\theta}_0)$ is unique, the consistency of selection holds;

if $f_0(z|\boldsymbol{\theta}_0)$ is not-unique (e.g. f_k, f_j), only BIC satisfies the consistency in

picking the simplest model: $\log \left(\frac{f_k(\mathbf{X}_n | \boldsymbol{\theta}_{k0})}{f_j(\mathbf{X}_n | \boldsymbol{\theta}_{j0})} \right) = O_p(1)$, and $c_{n,p} \rightarrow \infty$

¹discussed later in Expected predictive likelihood approach

Generalized Bayesian information criterion

Consider $\log \pi(\boldsymbol{\theta}) = O_p(n)$. The marginal likelihood is approximated as:

$$P(\mathbf{X}_n|M) \approx f(\mathbf{X}_n|\hat{\boldsymbol{\theta}}_n) \pi(\hat{\boldsymbol{\theta}}_n) \times \frac{(2\pi)^{\frac{p}{2}}}{n^{\frac{p}{2}} |S_n(\hat{\boldsymbol{\theta}}_n)|^{\frac{1}{2}}}$$

$S_n(\hat{\boldsymbol{\theta}})$ is minus the Hessian matrix of function $\frac{1}{n} \log f(\mathbf{X}_n|\boldsymbol{\theta})\pi(\boldsymbol{\theta})$ at the mode of posterior distribution $\hat{\boldsymbol{\theta}}_n$

Generalized Bayesian information criterion:

$$\text{GBIC} = -2 \log f(\mathbf{X}_n|\hat{\boldsymbol{\theta}}) - 2 \log \pi(\hat{\boldsymbol{\theta}}) + p \log n + \log |S_n(\hat{\boldsymbol{\theta}})| - p \log(2\pi)$$

Choosing model with the largest posterior \equiv the model that minimizes the criterion.

Difficulty: require to compute the Hessian matrix

1 Introduction

2 Theoretical background

- General framework
- Evaluation of marginal likelihood
 - Analytical method: Exact calculation of marginal likelihood
 - Laplace's method for computing the marginal likelihood
 - Bayesian information criterion (BIC)
- Expected predictive likelihood approach

3 Examples with Covid19 data

Expected predictive likelihood approach

Motivation: the main weakness of BF is its sensitivity to the prior

Let g be the true density of observations, $\mathbf{Z}_n = \{z_1, \dots, z_n\}$ is a set of i.i.d. unseen future observations

$f(\mathbf{Z}_n|\mathbf{X}_n, M)$: Bayesian predictive distribution

$$f(\mathbf{Z}_n|\mathbf{X}_n, M) = \int f(\mathbf{Z}_n|\mathbf{X}_n, \theta) \pi(\theta|\mathbf{X}_n) d\theta$$

Predictive ability of a given model M using Kullback-Leibler divergence

$$\begin{aligned} KL(g||f) &= \int \left[\log \frac{g(\mathbf{Z}_n)}{f(\mathbf{Z}_n|\mathbf{X}_n, M)} \right] g(\mathbf{Z}_n) d\mathbf{Z}_n \\ &= \int \log g(\mathbf{Z}_n) g(\mathbf{Z}_n) d\mathbf{Z}_n - \int \log (f(\mathbf{Z}_n|\mathbf{X}_n, M)) g(\mathbf{Z}_n) d\mathbf{Z}_n \end{aligned}$$

Expected predictive likelihood approach

Motivation: the main weakness of BF is its sensitivity to the prior

Let g be the true density of observations, $\mathbf{Z}_n = \{z_1, \dots, z_n\}$ is a set of i.i.d. unseen future observations

$f(\mathbf{Z}_n|\mathbf{X}_n, M)$: Bayesian predictive distribution

$$f(\mathbf{Z}_n|\mathbf{X}_n, M) = \int f(\mathbf{Z}_n|\mathbf{X}_n, \theta) \pi(\theta|\mathbf{X}_n) d\theta$$

Predictive ability of a given model M using Kullback-Leibler divergence

$$\begin{aligned} KL(g||f) &= \int \left[\log \frac{g(\mathbf{Z}_n)}{f(\mathbf{Z}_n|\mathbf{X}_n, M)} \right] g(\mathbf{Z}_n) d\mathbf{Z}_n \\ &= \int \log g(\mathbf{Z}_n) g(\mathbf{Z}_n) d\mathbf{Z}_n - \int \log (f(\mathbf{Z}_n|\mathbf{X}_n, M)) g(\mathbf{Z}_n) d\mathbf{Z}_n \end{aligned}$$

Expected log-predictive likelihood \rightarrow a general approach for evaluating the goodness of fit (validation step)

$$\eta(M) = \int \log f(\mathbf{Z}_n|\mathbf{X}_n, M) g(\mathbf{Z}_n) d\mathbf{Z}_n$$

1 Introduction

2 Theoretical background

- General framework
- Evaluation of marginal likelihood
 - Analytical method: Exact calculation of marginal likelihood
 - Laplace's method for computing the marginal likelihood
 - Bayesian information criterion (BIC)
- Expected predictive likelihood approach

3 Examples with Covid19 data

Problem setting

We will compare SIR models with different setting in modeling of Covid-19 pandemic evolution in Germany, Italy, Uruguay, and Saudi Arabia.

Model selection is performed in three steps

- Introduction of SIR models
- Parameter identification using Bayesian approach
- Evaluation of model marginal likelihood and Bayes factor

Modeling of Covid-19 pandemic evolution

SIR models

- **SIR model**

$$\begin{cases} \frac{dS}{dt} = -\beta \frac{IS}{N} \\ \frac{dI}{dt} = \beta \frac{IS}{N} - \gamma I \\ \frac{dR}{dt} = \gamma I \end{cases}$$

- S , I , R and N refer to the susceptible, infected, recovered and total populations, respectively
- $\theta = \{\beta, \gamma, S_0, I_0\}$ are parameter to be identified

Modeling of Covid-19 pandemic evolution

SIR models

- **SIR model**

$$\begin{cases} \frac{dS}{dt} = -\beta \frac{IS}{N} \\ \frac{dI}{dt} = \beta \frac{IS}{N} - \gamma I \\ \frac{dR}{dt} = \gamma I \end{cases}$$

- S , I , R and N refer to the susceptible, infected, recovered and total populations, respectively
- $\theta = \{\beta, \gamma, S_0, I_0\}$ are parameter to be identified

We consider two models $P(M_1) = P(M_2) = 0.5$:

- M_1 : β is varying with time while γ is constant ($p_1 = 7$)
- M_2 : both β and γ are time-dependent ($p_2 = 12$)

β in M_1 , M_2 , and γ in M_2 are represented as piecewise linear functions with a time interval of 20 days (except for Uruguay with 10 days interval).

Bayesian approach for parameter identification and model selection

- Let π_k be the prior densities of the parameters θ_k of model M_k
- \mathbf{X} : observations of new infected and recovered cases during time intervals of **14** days
 - reduce the sensitivity due to data errors

- Using Bayesian approach, the posterior of the parameters θ_i is obtained as

$$\pi_i(\theta_i|\mathbf{X}) \propto f(\mathbf{X}|\theta_i, M_i)\pi_i(\theta_i)$$

where f is the likelihood function

- Computational method: MH-MCMC
- Marginal predictive likelihood of model M (estimated using MCMC samples)

$$f(\mathbf{Z}|\mathbf{X}, M) = \int f(\mathbf{Z}|\theta_i, M_i)\pi_i(\theta_i|\mathbf{X})d\theta_i$$

- \mathbf{Z} : observations of new infected and recovered cases during time intervals of **20** days

Bayesian parameter identification

Likelihood function

Construct likelihood function is a difficult task

Motivation: the number of infected cases $>$ that confirmed

- using the conventional Gaussian likelihood function \rightarrow Probability of (predicted cases are smaller than confirmed case) ~ 0.5
- Seroprevalence studies of anti-SARS-CoV-2 IgG antibodies suggest 200% in Spain and 1100% in Switzerland (estimate from small population samples)
- **Assumption:** number of infections/recovers is up to 120% of that confirmed (belief based likelihood)
- (not ideal assumption)

Bayesian parameter identification

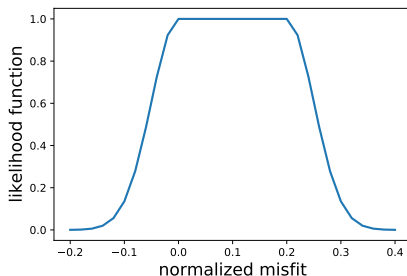
Likelihood function

Assumption: number of infections/recovers is up to 120% of that confirmed

Let $\overline{\text{err}}$ be the normalized misfit of the model prediction: $\overline{\text{err}} = \frac{x_i - y_i(\theta_k)}{x_i}$

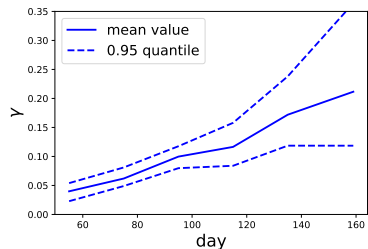
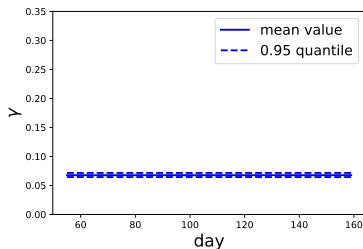
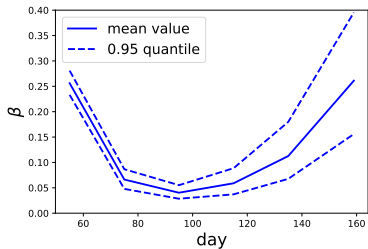
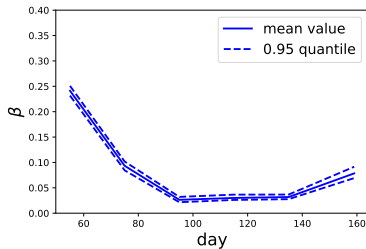
where x_i is an observation, and y_i is the corresponding model prediction

$$f(x_i|\theta_k) = \begin{cases} 1 & \text{if } 0 \leq \overline{\text{err}} \leq 0.2 \\ \mathcal{N}(0, 0.05^2) * \sqrt{2\pi}0.05 & \text{if } \overline{\text{err}} < 0 \\ \mathcal{N}(0.2, 0.05^2) * \sqrt{2\pi}0.05 & \text{if } \overline{\text{err}} > 0.2 \end{cases}$$



Results: Model 1 (left) Model 2 (right)

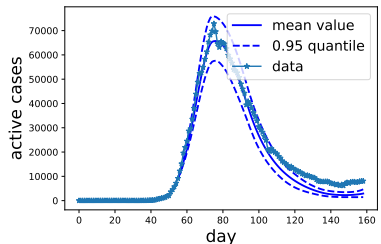
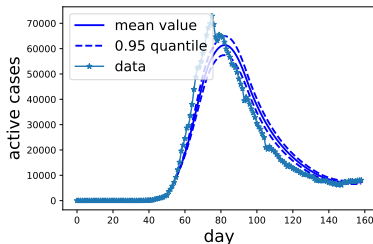
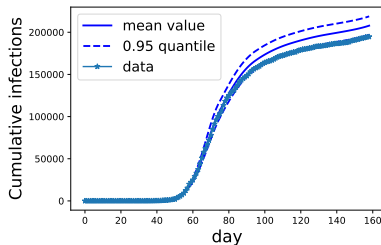
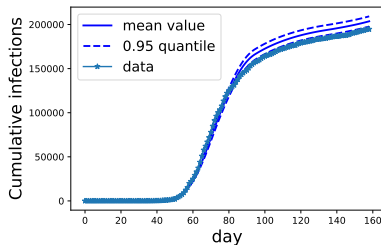
Germany (day 60: 22.03.2020, day 160: 28.06.2020)



MH-MCMC with 1 million samples

Results: Model 1 (left) Model 2 (right)

Germany, day 60: 22.03.2020, day 159: 28.06.2020



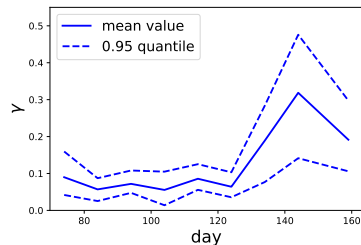
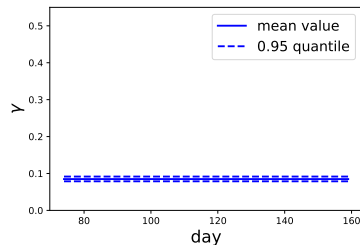
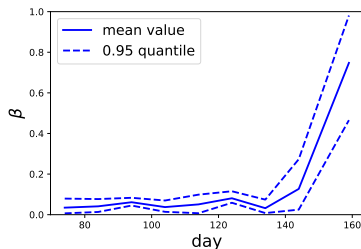
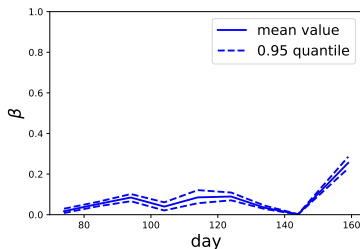
$$f(\mathbf{Z}|\mathbf{X}, M_1) = 1.e - 8$$

$$f(\mathbf{Z}|\mathbf{X}, M_2) = 0.19$$

(Current setting: $f = 1$ indicate best models)

Results: Model 1 (left) Model 2 (right)

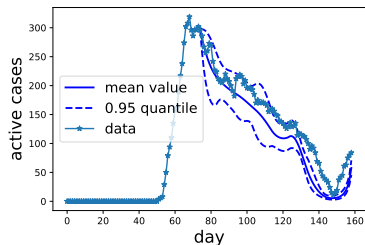
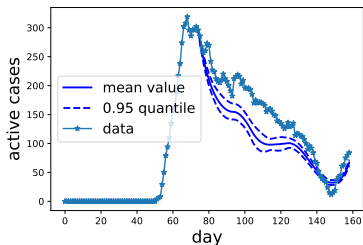
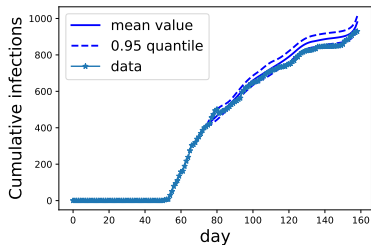
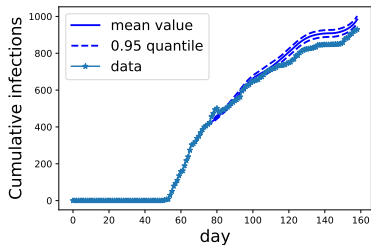
Uruguay, day 70: 01.04.2020, day 159: 28.06.2020



β and γ in Model 2 are represented as piecewise linear functions with time interval 10 days.

Results: Model 1 (left) Model 2 (right)

Uruguay, day 70: 01.04.2020, day 159: 28.06.2020

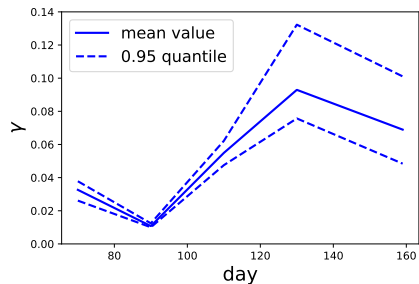
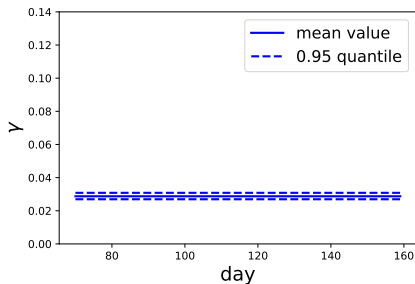
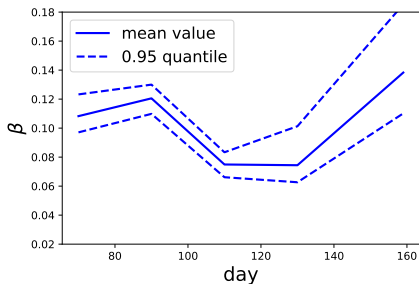
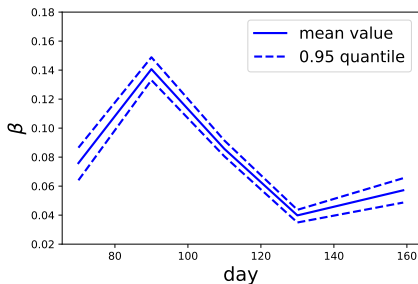


$$f(\mathbf{Z}|\mathbf{X}, M_1) = 3e - 11$$

$$f(\mathbf{Z}|\mathbf{X}, M_2) = 0.02$$

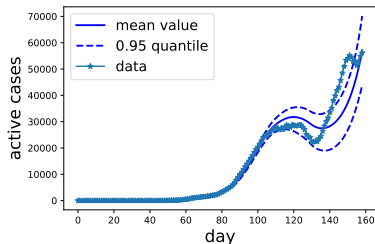
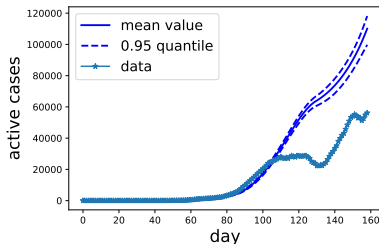
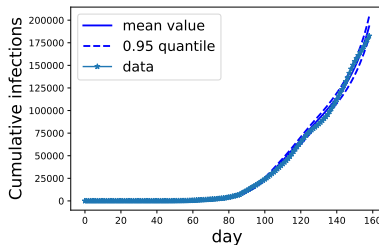
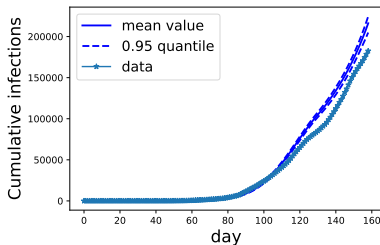
Results: Model 1 (left) Model 2 (right)

Saudi Arabia, day 70: 01.04.2020, day 159: 28.06.2020



Results: Model 1 (left) Model 2 (right)

Saudi Arabia, day 70: 01.04.2020, day 159: 28.06.2020

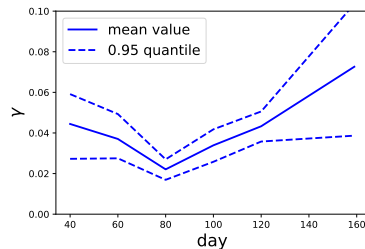
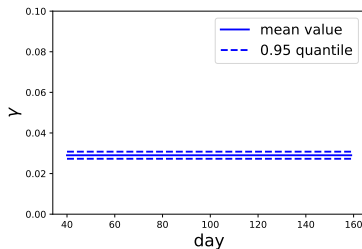
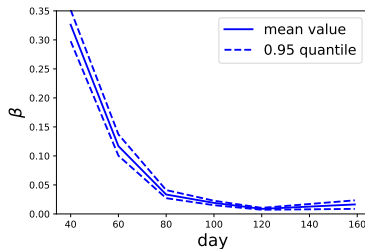
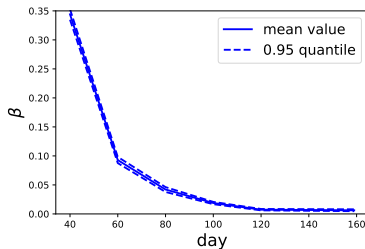


$$f(\mathbf{Z}|\mathbf{X}, M_1) = 3e - 4$$

$$f(\mathbf{Z}|\mathbf{X}, M_2) = 0.05$$

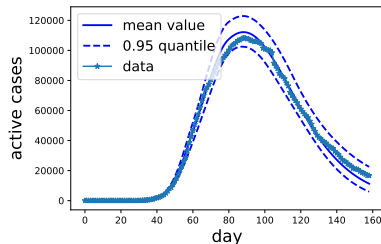
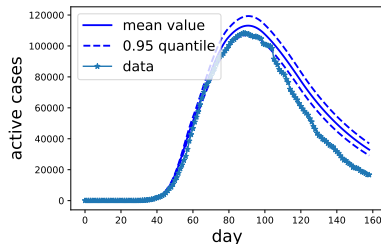
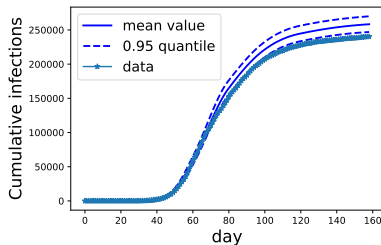
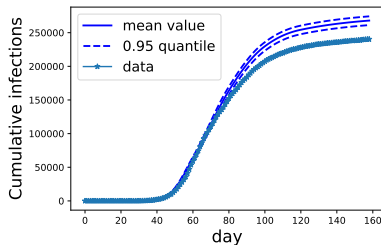
Results: Model 1 (left) Model 2 (right)

Italy, day 40: 02.03.2020, day 159: 28.06.2020



Results: Model 1 (left) Model 2 (right)

Italy, day 40: 02.03.2020, day 159: 28.06.2020



$$f(\mathbf{Z}|\mathbf{X}, M_1) = 0.07$$

$$f(\mathbf{Z}|\mathbf{X}, M_2) = 0.45$$

Summary

	$f(\mathbf{Z} \mathbf{X}, M_1)$	$f(\mathbf{Z} \mathbf{X}, M_2)$
Germany	1e-8	0.19
Uruguay	3e-11	0.02
Saudi Arabia	3e-4	0.05
Italy	0.07	0.45

There is strong evidence supporting Model 2

Conclusion

- An overview of Bayesian approach for model selection
- Different information criteria (BIC, GBIC, AIC)
- Predictive Bayesian distribution
- Apply for modeling the evolution of Covid-19 pandemic

Open problems:

- How to construct likelihood function
- Compare with SEIR, (E: exposed)
- Efficient method for identification