



Approximation Techniques for Bayesian Logistic Regression

Mauro Camara Escudero

Supervisor: Jonathan J. Forster

School of Mathematics, University of Southampton

Table of contents

1. Motivation and Problem Formulation
2. Stochastic Approximations
3. Deterministic Approximations
4. Results

Motivation and Problem Formulation

The Integration Problem

Many tasks in Bayesian Statistics can be seen as performing **expectation** with respect to a posterior distribution

$$\mathbb{E}_{p(\boldsymbol{\theta}|\mathbf{x})} [f(\boldsymbol{\theta})] = \int_{\Theta} f(\boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{x}) d\boldsymbol{\theta}$$

The resulting integration is often **computationally intractable**.

Overview of Approximate Inference

STOCHASTIC APPROXIMATIONS

- *Monte Carlo*: Averages independent samples drawn from desired distribution.
- *Markov Chain Monte Carlo*: Monte Carlo on dependent samples produced by suitable Markov Chain.

Overview of Approximate Inference

STOCHASTIC APPROXIMATIONS

- *Monte Carlo*: Averages independent samples drawn from desired distribution.
- *Markov Chain Monte Carlo*: Monte Carlo on dependent samples produced by suitable Markov Chain.

DETERMINISTIC APPROXIMATIONS

- *Laplace Approximation*: Approximates $p(\boldsymbol{\theta} \mid \mathbf{x})$ with normal distribution centered at the mode.
- *Variational Inference*: Approximates $p(\boldsymbol{\theta} \mid \mathbf{x})$ with the closest distribution $q(\boldsymbol{\theta})$ according to some objective function.

Bayesian Logistic Regression

Relationship between explanatory variables and response modelled with a **Generalized Linear Model**:

- Exponential Family distribution: $Y_i \sim \text{Bernoulli}(\pi_i)$
- Link Function: Logit $g(\pi_i) = \ln \left(\frac{\pi_i}{1-\pi_i} \right) = \mathbf{x}_i^\top \boldsymbol{\beta}$

Choose a **normal prior** $\boldsymbol{\beta} \sim \mathcal{N}(\mu_0, \Sigma_0)$.

Stochastic Approximations

- PSEUDO-RANDOM NUMBER GENERATORS allow us to sample from a **uniform** distribution $\mathcal{U}(0, 1)$.

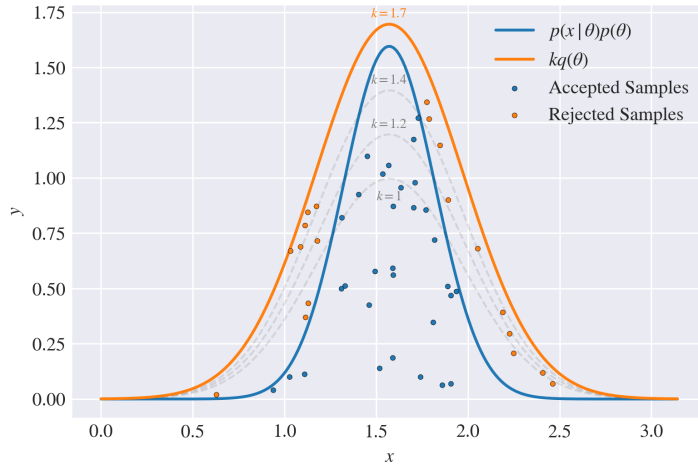
Sampling Methods

- PSEUDO-RANDOM NUMBER GENERATORS allow us to sample from a uniform distribution $\mathcal{U}(0, 1)$.
- Samples from non-uniform distributions can be obtained by transforming uniform samples but not always feasible, and computationally demanding.

Sampling Methods

- PSEUDO-RANDOM NUMBER GENERATORS allow us to sample from a **uniform** distribution $\mathcal{U}(0, 1)$.
- Samples from non-uniform distributions can be obtained by transforming uniform samples but not always feasible, and computationally demanding.
- REJECTION SAMPLING draws uniform samples from a *function* whose graph is always above that of $p(\boldsymbol{\theta} \mid \mathbf{x})$ and accepts them as samples from $p(\boldsymbol{\theta} \mid \mathbf{x})$ if they are also under its graph.

Rejection Sampling Example



Monte Carlo Methods

- MONTE CARLO method uses $\theta_1, \dots, \theta_N$ independent samples from $p(\theta \mid \mathbf{x})$ to approximate the expectation

$$\mathbb{E}_{p(\theta \mid \mathbf{x})} [f(\theta)] \approx \frac{1}{N} \sum_{i=1}^N f(\theta_i)$$

this approximation converges by the law of large numbers.

Monte Carlo Methods

- MONTE CARLO method uses $\theta_1, \dots, \theta_N$ independent samples from $p(\theta | \mathbf{x})$ to approximate the expectation

$$\mathbb{E}_{p(\theta|\mathbf{x})} [f(\theta)] \approx \frac{1}{N} \sum_{i=1}^N f(\theta_i)$$

this approximation converges by the law of large numbers.

- UNIFORM SAMPLING: Rather than $\theta_i \sim p(\theta | \mathbf{x})$, draw θ_i uniformly in Θ and use a weighted average of $f(\theta_i)p(\theta_i | \mathbf{x})$ instead. Inefficient because most time spent in regions of Θ with negligible mass.

Monte Carlo Methods

- MONTE CARLO method uses $\theta_1, \dots, \theta_N$ independent samples from $p(\theta | \mathbf{x})$ to approximate the expectation

$$\mathbb{E}_{p(\theta|\mathbf{x})} [f(\theta)] \approx \frac{1}{N} \sum_{i=1}^N f(\theta_i)$$

this approximation converges by the law of large numbers.

- UNIFORM SAMPLING: Rather than $\theta_i \sim p(\theta | \mathbf{x})$, draw θ_i uniformly in Θ and use a weighted average of $f(\theta_i)p(\theta_i | \mathbf{x})$ instead. Inefficient because most time spent in regions of Θ with negligible mass.
- IMPORTANCE SAMPLING: Draw samples from a proposal distribution $\theta_i \sim q(\theta)$ resembling $p(\theta | \mathbf{x})$ and uses a weighted average with whose weights w_i compensate the error introduced by sampling from wrong distribution. Still inefficient in high dimensions.

Markov Chain Monte Carlo

- MARKOV CHAIN MONTE CARLO (MCMC) trades off the costly independence and exactness of the samples for **approximate, dependent** samples that are **cheaper** to compute.

Markov Chain Monte Carlo

- MARKOV CHAIN MONTE CARLO (MCMC) trades off the costly independence and exactness of the samples for **approximate, dependent** samples that are **cheaper** to compute.
- $\theta_1, \dots, \theta_N$ generated by a Markov Chain whose **equilibrium distribution** is $p(\theta \mid \mathbf{x})$.

Markov Chain Monte Carlo

- MARKOV CHAIN MONTE CARLO (MCMC) trades off the costly independence and exactness of the samples for **approximate, dependent** samples that are **cheaper** to compute.
- $\theta_1, \dots, \theta_N$ generated by a Markov Chain whose **equilibrium distribution** is $p(\theta \mid \mathbf{x})$.
- RANDOM WALK METROPOLIS-HASTINGS (RWMH) iteratively draws a sample from a **symmetric** proposal distribution $\theta^* \sim q(\theta^* \mid \theta_i)$ depending only on the current sample value θ_i . Accepts θ^* with a probability that makes the chain converge to $p(\theta \mid \mathbf{x})$.

Deterministic Approximations

Laplace Approximation

Approximates $p(\boldsymbol{\theta} \mid \mathbf{x})$ with a multivariate normal distribution

$$q(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} \mid \boldsymbol{\theta}_0, -H(\boldsymbol{\theta}_0)^{-1})$$

centered at the mode $\boldsymbol{\theta}_0$ and with variance-covariance matrix

$$-\nabla^2 \ln(p(\mathbf{x} \mid \boldsymbol{\theta})p(\boldsymbol{\theta}))^{-1} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}$$

the negative inverse Hessian matrix of $\ln(p(\mathbf{x} \mid \boldsymbol{\theta})p(\boldsymbol{\theta}))$.

Laplace Approximation for Bayesian Logistic Regression

In Bayesian Logistic Regression the multivariate normal distribution is given by

$$q(\boldsymbol{\beta}) = \mathcal{N}(\boldsymbol{\beta} \mid \boldsymbol{\beta}_0, -H(\boldsymbol{\beta}_0)^{-1})$$

with

$$-H(\boldsymbol{\beta}_0) = \boldsymbol{\Sigma}_0^{-1} + \sum_{i=1}^n \pi_i(\boldsymbol{\beta}_0)(1 - \pi_i(\boldsymbol{\beta}_0)) \mathbf{x}_i \mathbf{x}_i^{\top}$$

- Define a family of distributions \mathcal{D} .

- Define a family of distributions \mathcal{D} .
- Define **objective function** measuring distance between $p(\boldsymbol{\theta} \mid \mathbf{x})$ and $q(\boldsymbol{\theta}) \in \mathcal{D}$. Usually **KL divergence**.

$$KL(q(\boldsymbol{\theta})||p(\boldsymbol{\theta} \mid \mathbf{x})) = \mathbb{E}_{q(\boldsymbol{\theta})} \left[\ln \left(\frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta} \mid \mathbf{x})} \right) \right]$$

- Define a family of distributions \mathcal{D} .
- Define **objective function** measuring distance between $p(\boldsymbol{\theta} \mid \mathbf{x})$ and $q(\boldsymbol{\theta}) \in \mathcal{D}$. Usually **KL divergence**.

$$KL(q(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta} \mid \mathbf{x})) = \mathbb{E}_{q(\boldsymbol{\theta})} \left[\ln \left(\frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta} \mid \mathbf{x})} \right) \right]$$

- Choose distribution minimizing KL divergence, or equivalently choose distribution *maximizing* **Evidence Lower Bound** (ELBO)

$$\text{elbo}(q(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta} \mid \mathbf{x})) = \mathbb{E}_{q(\boldsymbol{\theta})} [\ln (p(\mathbf{x} \mid \boldsymbol{\theta}))] - KL(q(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta}))$$

- Bound each log success probability $\ln(\pi_i)$ in the likelihood with a **quadratic function** of $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$ depending on variational parameters ξ_i .

Local Variational Methods for Bayesian Logistic Regression

- Bound each log success probability $\ln(\pi_i)$ in the likelihood with a **quadratic function** of $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$ depending on variational parameters ξ_i .
- This results in a **Gaussian** bound for the posterior.

$$q(\boldsymbol{\beta}) = \mathcal{N}(\boldsymbol{\beta} \mid \boldsymbol{\mu}_v, \boldsymbol{\Sigma}_v)$$

with $\boldsymbol{\mu}_v$ and $\boldsymbol{\Sigma}_v$ depending on the ξ_i 's.

Local Variational Methods for Bayesian Logistic Regression

- Bound each log success probability $\ln(\pi_i)$ in the likelihood with a **quadratic function** of $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$ depending on variational parameters ξ_i .
- This results in a **Gaussian** bound for the posterior.

$$q(\boldsymbol{\beta}) = \mathcal{N}(\boldsymbol{\beta} \mid \boldsymbol{\mu}_v, \boldsymbol{\Sigma}_v)$$

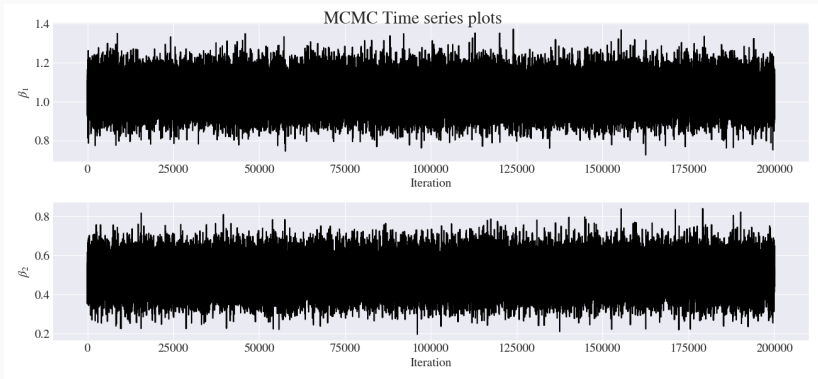
with $\boldsymbol{\mu}_v$ and $\boldsymbol{\Sigma}_v$ depending on the ξ_i 's.

- Maximize this bound using the **EM algorithm**

Results

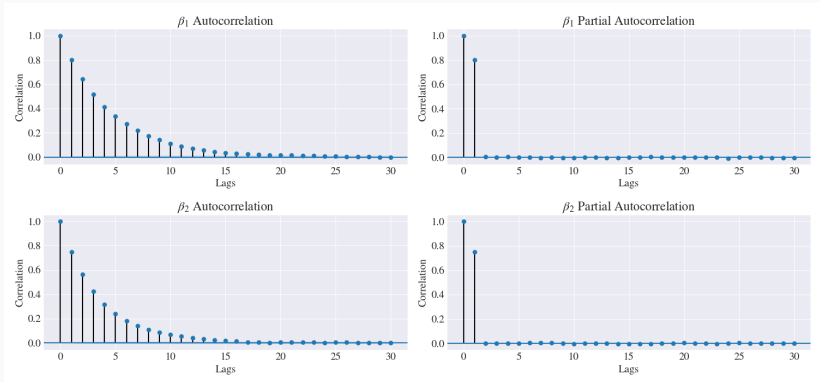
Models with Explanatory Variables

RWMH convergence diagnosed via **trace plots**.

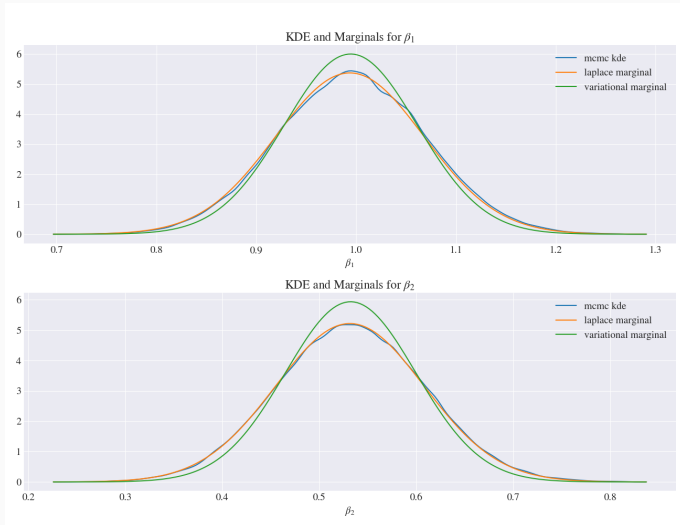


Models with Explanatory Variables

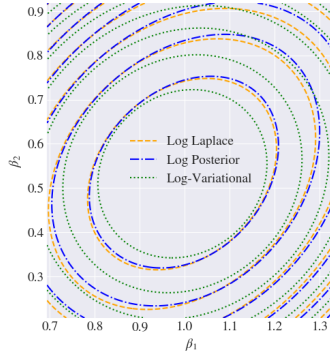
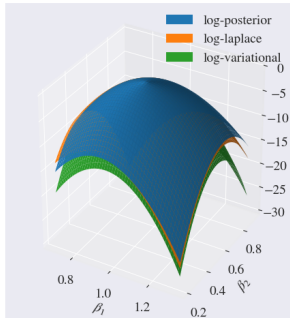
Markov Property diagnosed with Auto-Correlation and Partial Auto-Correlation plots.



Models with Explanatory Variables



Models with Explanatory Variables



Conclusion

Laplace outperforms Variational, which underestimates the tails.

Model without Explanatory Variables

- Posterior distribution

$$p(\beta | \mathbf{y}) = \frac{1}{B(n\bar{y}, n - n\bar{y})} \frac{e^{\beta n\bar{y}}}{(1 + e^{\beta})^n}$$

Model without Explanatory Variables

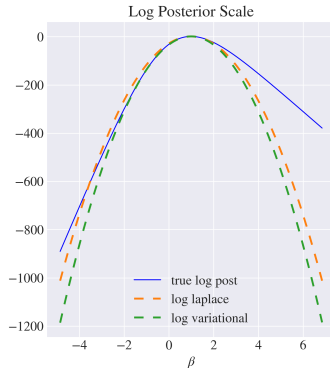
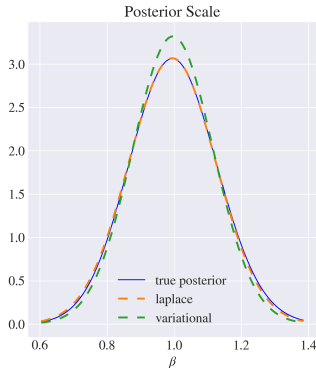
- Posterior distribution

$$p(\beta | \mathbf{y}) = \frac{1}{B(n\bar{y}, n - n\bar{y})} \frac{e^{\beta n\bar{y}}}{(1 + e^{\beta})^n}$$

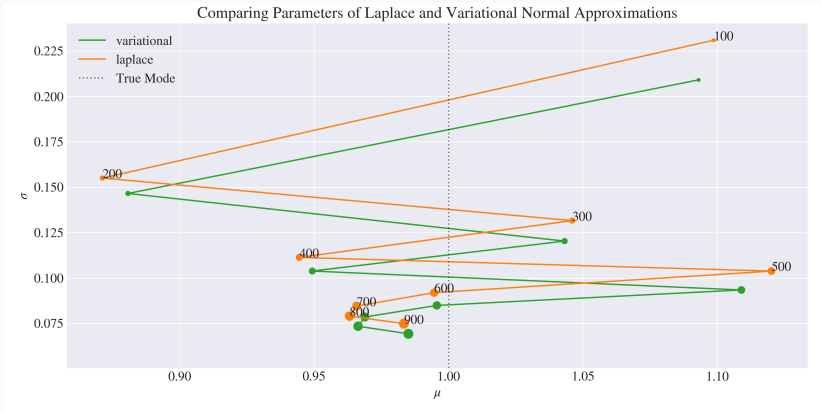
- Laplace approximation

$$q_l(\beta) = \mathcal{N} \left(\ln \left(\frac{\bar{y}}{1 - \bar{y}} \right), \frac{1}{n\bar{y}(1 - \bar{y})} \right)$$

Model without Explanatory Variables



Model without Explanatory Variables



Conclusion

Variational Mean closer to true population mode β , especially for small data set sizes.

- Multimodal posterior distributions.

- Multimodal posterior distributions.
- Different variational bound.

- Multimodal posterior distributions.
- Different variational bound.
- Relationship between Laplace and Variational distributions.

Questions?

Backup Slide - Variational EM

Algorithm 1: Variational Approximation

- 1 Initialize $\boldsymbol{\xi} = (\xi_1^{(1)}, \dots, \xi_n^{(1)})^\top \in \mathbb{R}^n$ randomly.
- 2 **for** $j = 1, 2, \dots, \Delta$:
- 3 **for** $i = 1, 2, \dots, n$ **do**:
- 4 Find mean and variance-covariance matrix and update variational parameters.

$$\Sigma_v = \left(2 \sum_{i=1}^n \lambda(\xi_i^{(j)}) \mathbf{x}_i \mathbf{x}_i^\top + \Sigma_0^{-1} \right)^{-1}$$

$$\boldsymbol{\mu}_v = \Sigma_v \left(\sum_{i=1}^n \mathbf{x}_i \left(y_i - \frac{1}{2} \right) + \Sigma_0^{-1} \boldsymbol{\mu}_0 \right)$$

$$\xi_i^{(j+1)} = \sqrt{\mathbf{x}_i^\top (\Sigma_v + \boldsymbol{\mu}_v \boldsymbol{\mu}_v^\top) \mathbf{x}_i}$$

end

5 end