# Automated Posterior Interval Evaluation for Inference in Probabilistic Programming

## Edward Kao, Michael Yee | MIT Lincoln Laboratory

## ABSTRACT

In probabilistic inference, credible intervals constructed from posterior samples provide ranges of likely values for continuous parameters of interest. Intuitively, an inference procedure is optimal if it produces the most precise posterior intervals that cover the true parameter value with the expected frequency in repeated experiments. We present theories and methods for automating posterior interval evaluation of inference performance in probabilistic programming using two metrics: 1.) truth coverage, and 2.) ratio of the empirical over the ideal interval widths. Demonstrating with inference on popular regression and state-space models, we show how the metrics provide effective comparisons between different inference procedures, and capture the effects of collinearity and model misspecification. Overall, we claim such automated interval evaluation can accelerate the robust design and comparison of probabilistic inference programs by directly diagnosing how accurately and precisely they can estimate parameters of interest.

## THEORY

Based on the statistical principle of evaluating Bayesian inference with frequentist properties (1,2), we compute two metrics for inference output in repeated experiments: 1.) posterior credible interval coverage of the true parameter value (90% intervals should cover the truth 90% of the time), and 2.) ratio of the empirical over the ideal interval widths (ratio of 1 indicates precise inference). The ideal interval width can be computed based on the asymptotic theorem:

**The Bernstein-von Mises Theorem** (3). For regular models, the posterior distribution of continuous parameters in finite dimensions converges asymptotically, with increasing data, in distribution to Normal with mean at the true value $\theta^*$ and covariance equal to the inverse of the Fisher information matrix $\mathcal{I}$ evaluated at $\theta^*$:

$$\theta \xrightarrow{D} \text{Normal}(\theta^*, \mathcal{I}(\theta^*)^{-1}) \qquad [1]$$

The diagonal terms of the asymptotic covariance $\mathcal{I}(\theta^*)^{-1}$ provide the ideal interval width for each parameter in the model. In the non-asymptotic regime, the ideal interval width can be computed using the Laplace approximation where $q''(\theta^*)$ is the Hessian of the log posterior distribution evaluated at $\theta^*$:

$$\theta \approx \text{Normal}(\theta^*, -q''(\theta^*)^{-1}) \qquad [2]$$

## AUTOMATION

Computing the proposed metrics based on posterior intervals can be automated for any probabilistic programming systems (4) that **simulate** data $\mathcal{D}$ and parameters $\theta$ based on statistical models $M$ and priors, and **infer** the posterior distribution such that the likelihood function and the unnormalized posterior distribution are accessible. The Fisher information matrix can be computed via the hessian function on the log likelihood, through auto-differentiation. For non-asymptotic cases using the Laplace approximation, simply replace the likelihood with the unnormalized posterior distribution. We implement and demonstrate this automated evaluation in Gen (5).

```
function simulate(M)
    θ ~ p_M(·)
    D|θ ~ p_M(·|θ)
    return (θ, D)
end
```

```
function ll-hessian(M, θ, D)
    return [ ∂ log p_M(D|θ) / ∂θ_u θ_v ]_{u,v ∈ θ×θ}
end
```

### Algorithm for Computing the Truth Coverage Rate

**Input** : Model $M$, inference program `infer`, interval probability $ci$, #simulations $S$, #inference samples $T$

**Output** : Coverage rates for $S$ simulated datasets

```
1  for s ← 1 to S do
2      (θ*, D) ← simulate(M)
3      outcomes ← [ ]
4      θ^{1:T} ← infer(M, D, T)
5      for u ∈ θ do
6          (θ_lo, θ_hi) ← empirical-interval(θ_u^{1:T}, ci)
7          outcomes.append(θ_lo ≤ θ_u* ≤ θ_hi)
8      end
9  end
10 return mean(outcomes)
```

### Algorithm for Computing the Interval Width Ratios

**Input** : Model $M$, inference program `infer`, interval probability $ci$, #simulations $S$, #inference samples $T$

**Output** : Empirical to ideal width ratios

```
   ratios ← [ ]
1  for s ← 1 to S do
2      (θ*, D) ← simulate(M)
3      I(θ*) ← -ll-hessian(M, θ*, D)
4      σ* ← sqrt(diag(I(θ*)^{-1}))
5      θ^{1:T} ← infer(M, D, T)
6      for u ∈ θ do
7          (θ_lo, θ_hi) ← empirical-interval(θ_u^{1:T}, ci)
8          (θ_lo*, θ_hi*) ← ideal-interval(θ_u*, σ_u*, ci)
9          ratios.append((θ_hi - θ_lo)/(θ_hi* - θ_lo*))
10     end
11 end
12 return ratios
```

## DEMONSTRATION ON BAYESIAN LINEAR REGRESSION

$$w|\alpha \sim \text{Normal}(m_0, \alpha^{-1}I)$$
$$y|X, w, \beta \sim \text{Normal}(Xw, \beta^{-1}I)$$

where $w$ and $m_0 \in \mathbb{R}^d$, $y \in \mathbb{R}^N$, and $X$ is an $N \times d$ covariate matrix. Here, $d = 2$, $N = 30$ and the covariates $X$ are generated independently, or with collinearity of 0.9 correlation.

For purpose of demonstration, we infer $w$ using Gibbs sampling, even though the posterior has a closed-form expression. Evaluation cases: 1.) Regular Gibbs sampling, 2.) Covariates generated with collinearity, and 3.) Prior location is misplaced.



Regular Case | Collinear Covariates | Misplaced Prior

Empirical to Ideal Interval Width Ratio / Truth Coverage / Number of samples after burn-in

**Posterior interval evaluation identifies the expected effects of collinearity and misplaced prior on inference**

## DEMONSTRATION ON BAYESIAN LOGISTIC REGRESSION

$$w|\alpha \sim \text{Normal}(m_0, \alpha^{-1}I)$$
$$y|X, w \sim \text{Bernoulli}(\text{sigmoid}(Xw))$$

where $w$ and $m_0 \in \mathbb{R}^d$, $y \in \{0,1\}^N$, and $X$ is an $N \times d$ covariate matrix. Here, $d = 10$, $N = 100$ and the covariates $X$ are generated with collinearity of 0.95 correlation

We infer $w$ using Random Walk Metropolis-Hastings with two multivariate normal proposals: $w' \sim \text{Normal}(w, \Sigma)$:

1. Scott proposal (approximation to asymptotically optimal proposal):
$$\Sigma_{Scott} = \left(V_0^{-1} + \frac{6}{\pi^2}X^TX\right)^{-1},$$
where $V_0$ is the covariance of multivariate normal prior on $w$

2. Naïve proposal (diagonal covariance matrix): $\Sigma = 0.2\,I$



Empirical to Ideal Interval Width Ratio / Truth Coverage / Number of samples after burn-in

Proposal w/ Scott covariance
Proposal w/ diagonal covariance

**Evaluation quantifies how much faster the Scott proposal converges than the naïve proposal, under collinearity**

## DEMONSTRATION ON NONLINEAR STATE-SPACE MODEL

$$x_t = \frac{x_{t-1}}{2} + 25\frac{x_{t-1}}{1+x_{t-1}^2} + 8\cos(0.1t) + \delta_{t-1}$$
$$y_t = x_t + \epsilon_t$$

$$x_1 \sim \text{Normal}(\mu, \nu^2)$$
$$\delta_{t-1} \sim \text{Normal}(0, \omega^2)$$
$$\epsilon_t \sim \text{Normal}(0, \sigma^2)$$

A popular nonlinear state-space model in the literature. Here we generate 100 steps in time to capture both the periodic motion and nonlinear drift.

Infer the states $x$ given the observations $y$ using particle filters. Compare standard particle filters against one with rejuvenation moves on past states.



Empirical to Ideal Interval Width Ratio / Truth Coverage / Number of Particles

Particle filter with rejuvenation
Particle filter

**Evaluation shows with the addition of rejuvenation, same performance is reached with far fewer particles**

## FUTURE WORK

- Apply the proposed evaluation to real-world scenarios with a single data realization and unknown truth.
- Extend approach to general models with a mixture of regular and irregular parameters through conditioning and exploring generalizations of the Bernstein-von Mises Theorem.
- Explore simple and automated indicators for the adequacy of the normal assumption on the true posterior.
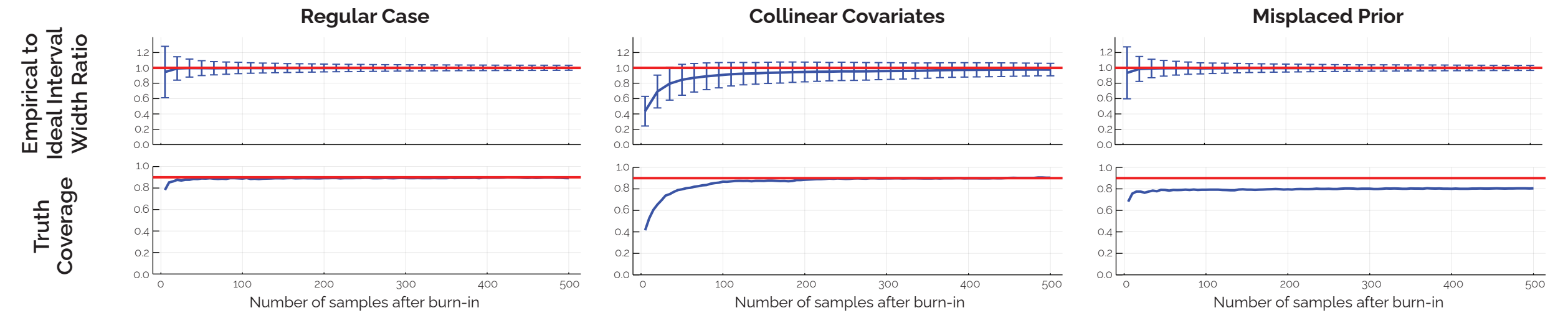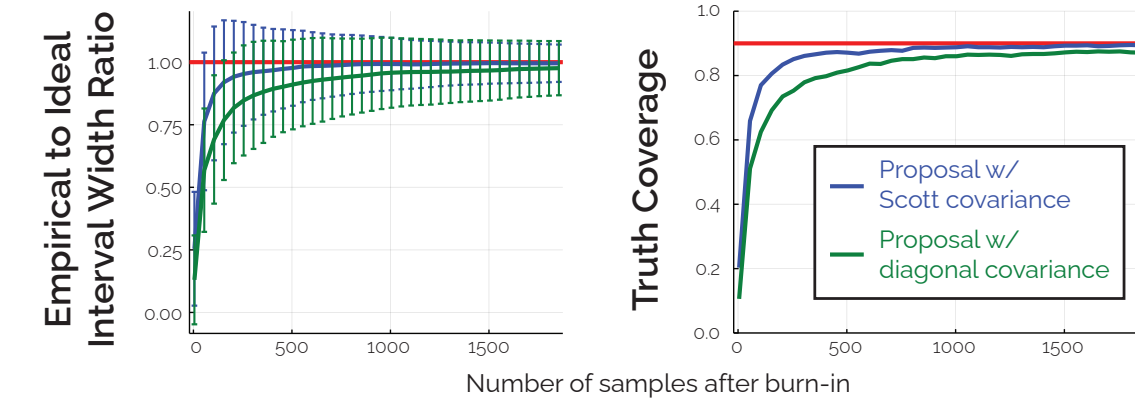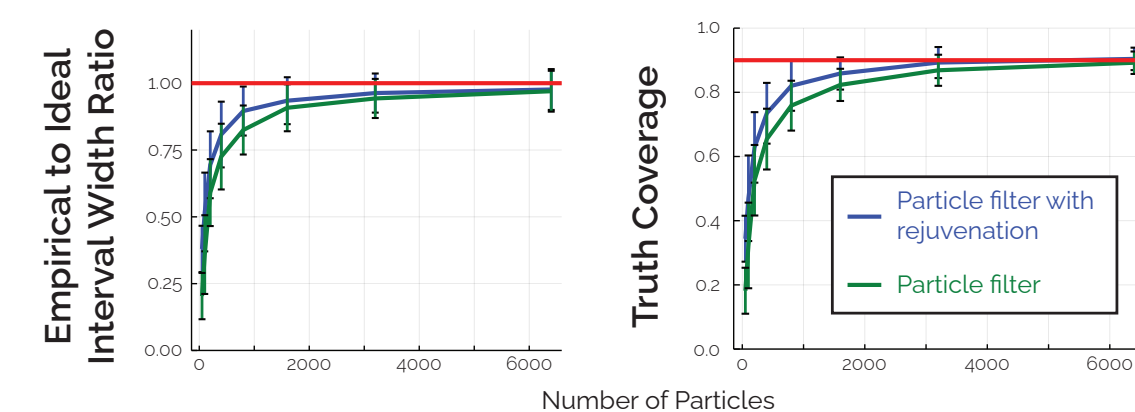
1. Rubin, D.B. Bayesianly justifiable and relevant frequency calculations for the applied statistician. The Annals of Statistics (1984).
2. Gelman, A., Carlin, J.B., Stern, H.S., and Rubin, D.B. Bayesian data analysis. Taylor & Francis (2014).
3. Van der Vaart, A.W. Asymptotic statistics, 10.2 Bernstein–von Mises Theorem. Cambridge University Press (1998).
4. Gordon, A.D., Henzinger, T.A., Nori, A.V., and Rajamani, S.K. Probabilistic programming. ACM (2014).
5. Cusumano-Towner, M.F., Saad, F.A., Lew, A.K., and Mansinghka V.K. Gen: a general-purpose probabilistic programming system with programmable inference. ACM (2019).

**LINCOLN LABORATORY**
**MASSACHUSETTS INSTITUTE OF TECHNOLOGY**