

Laplace Approximation

Advanced Statistical Inference

Simone Rossi

Where are we?

Given a likelihood $p(\mathbf{y} \mid \boldsymbol{\theta})$ and a prior $p(\boldsymbol{\theta})$, the goal is to compute the posterior distribution $p(\boldsymbol{\theta} \mid \mathbf{y})$

$$p(\boldsymbol{\theta} \mid \mathbf{y}) = \frac{p(\mathbf{y} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{y})}$$

...

We have seen that exact inference is possible if the prior and likelihood are *conjugate* distributions (e.g. for Bayesian linear regression with Gaussian likelihood and Gaussian prior).

...

But in most cases, the Bayesian approach is intractable.

- We cannot compute the posterior distribution $p(\boldsymbol{\theta} \mid \mathbf{y})$ in closed form.
- We cannot make predictions for a new point \mathbf{y}_\star in closed form $p(\mathbf{y}_\star \mid \mathbf{y}) = \int p(\mathbf{y}_\star \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid \mathbf{y}) d\boldsymbol{\theta}$.

Last Week

We have seen that we can use **Monte Carlo** methods to approximate the posterior distribution $p(\boldsymbol{\theta} \mid \mathbf{y})$.

- **Rejection sampling**: sample from a simple distribution and accept/reject samples
- **Markov Chain Monte Carlo (MCMC)**: sample from a Markov chain that converges to the target distribution.
 - *Metropolis-Hastings*: do random walk in the parameter space and accept/reject samples based on the ratio of the target distribution in the current and next state.

- *Hamiltonian Monte Carlo (HMC)*: explore the parameter space by simulating a particle with a potential given by negative log posterior and a random momentum.

...

This Week

In this week, we will see another way to approximate the posterior distribution $p(\theta \mid \mathbf{y})$.

Can we use a simpler distribution to approximate the posterior?

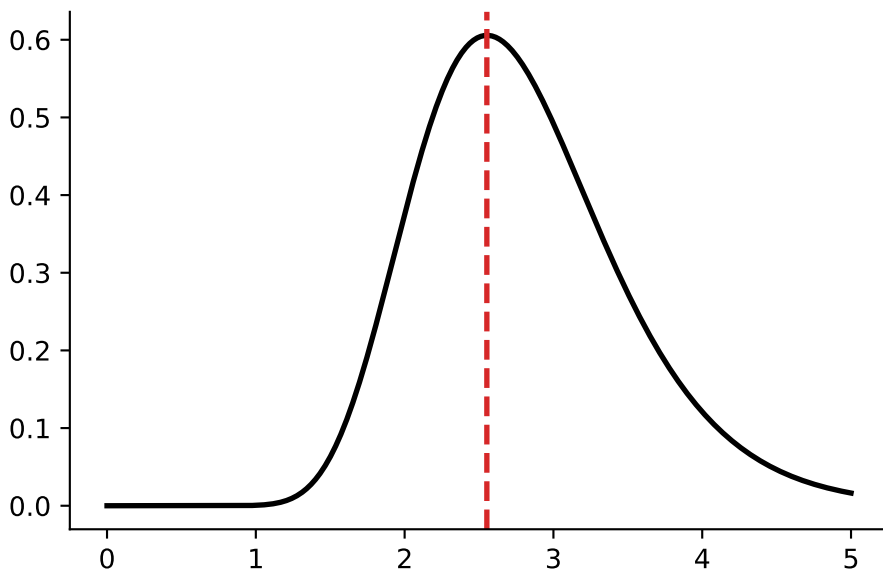
And if so, **how?**

Laplace Approximation

Laplace Approximation

Idea:

- Let's find a way to have a local approximation of the posterior distribution $p(\theta \mid \mathbf{y})$ around the mode $\hat{\theta}$.



Derivation of Laplace Approximation

Given a likelihood $p(\mathbf{y} \mid \boldsymbol{\theta})$ and a prior $p(\boldsymbol{\theta})$, the posterior distribution is given by

$$p(\boldsymbol{\theta} \mid \mathbf{y}) = \frac{1}{Z} \exp(-\mathcal{U}(\boldsymbol{\theta}))$$

where Z is the normalizing constant and $\mathcal{U}(\boldsymbol{\theta}) = -\log p(\mathbf{y} \mid \boldsymbol{\theta}) - \log p(\boldsymbol{\theta})$ is the (un-normalized) negative log posterior.

...

Taylor expand $\mathcal{U}(\boldsymbol{\theta})$ around the mode $\hat{\boldsymbol{\theta}}$ to the second order

$$\mathcal{U}(\boldsymbol{\theta}) \approx \mathcal{U}(\hat{\boldsymbol{\theta}}) + (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \mathbf{g} + \frac{1}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \mathbf{H}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})$$

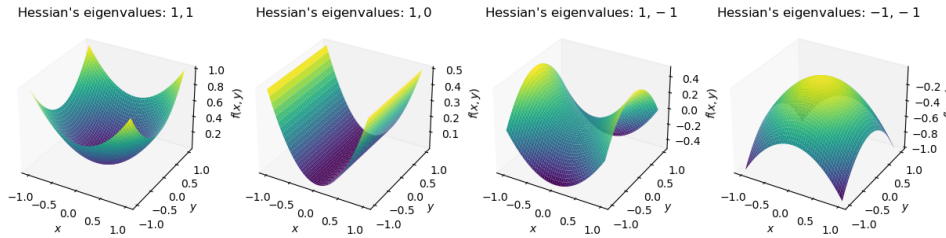
where:

- $\mathbf{g} = \nabla \mathcal{U}(\hat{\boldsymbol{\theta}})$ is the gradient
- $\mathbf{H} = \nabla^2 \mathcal{U}(\hat{\boldsymbol{\theta}})$ is the Hessian matrix

Hessians

The Hessian matrix \mathbf{H} is a square matrix $D \times D$ of second-order partial derivatives of $\mathcal{U}(\boldsymbol{\theta})$ computed in $\boldsymbol{\theta}$.

$$\mathbf{H} = \nabla^2 \mathcal{U}(\boldsymbol{\theta}) \begin{bmatrix} \frac{\partial^2 \mathcal{U}}{\partial \theta_1^2}(\boldsymbol{\theta}) & \frac{\partial^2 \mathcal{U}}{\partial \theta_1 \partial \theta_2}(\boldsymbol{\theta}) & \cdots & \frac{\partial^2 \mathcal{U}}{\partial \theta_1 \partial \theta_D}(\boldsymbol{\theta}) \\ \frac{\partial^2 \mathcal{U}}{\partial \theta_2 \partial \theta_1}(\boldsymbol{\theta}) & \frac{\partial^2 \mathcal{U}}{\partial \theta_2^2}(\boldsymbol{\theta}) & \cdots & \frac{\partial^2 \mathcal{U}}{\partial \theta_2 \partial \theta_D}(\boldsymbol{\theta}) \\ \cdots & \cdots & \ddots & \cdots \\ \frac{\partial^2 \mathcal{U}}{\partial \theta_D \partial \theta_1}(\boldsymbol{\theta}) & \frac{\partial^2 \mathcal{U}}{\partial \theta_D \partial \theta_2}(\boldsymbol{\theta}) & \cdots & \frac{\partial^2 \mathcal{U}}{\partial \theta_D^2}(\boldsymbol{\theta}) \end{bmatrix}$$



Derivation of Laplace Approximation

$$\mathcal{U}(\boldsymbol{\theta}) \approx \mathcal{U}(\hat{\boldsymbol{\theta}}) + (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \mathbf{g} + \frac{1}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \mathbf{H}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})$$

- The mode $\hat{\boldsymbol{\theta}}$ is the value of $\boldsymbol{\theta}$ that maximizes $p(\boldsymbol{\theta} \mid \mathbf{y})$ or equivalently minimizes $\mathcal{U}(\boldsymbol{\theta})$.
- Because we are expanding around the mode, $\mathbf{g} = \nabla \mathcal{U}(\hat{\boldsymbol{\theta}}) = 0$.

Approximate the posterior distribution with this expansion:

$$\begin{aligned} p(\boldsymbol{\theta} \mid \mathbf{y}) &\approx \frac{1}{Z} \exp \left(-\mathcal{U}(\hat{\boldsymbol{\theta}}) - \frac{1}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \mathbf{H}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \right) \\ &= \frac{1}{Z} \exp \left(-\mathcal{U}(\hat{\boldsymbol{\theta}}) \right) \exp \left(-\frac{1}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \mathbf{H}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \right) \end{aligned}$$

Derivation of Laplace Approximation

$$p(\boldsymbol{\theta} \mid \mathbf{y}) \approx \frac{1}{Z} \exp \left(-\mathcal{U}(\hat{\boldsymbol{\theta}}) \right) \exp \left(-\frac{1}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \mathbf{H}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \right)$$

Compare this to the form of a Gaussian distribution:

$$\mathcal{N}(\boldsymbol{\theta} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{2} (2\pi)^{D/2} \det(\boldsymbol{\Sigma})^{-1/2} \exp \left(-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\theta} - \boldsymbol{\mu}) \right)$$

Taylor expansion of the log posterior \Rightarrow Gaussian approximation of the posterior.

- The mean of the Gaussian is the mode of the posterior $\hat{\boldsymbol{\theta}}$.
- The covariance of the Gaussian is the inverse of the Hessian \mathbf{H}^{-1} .
- The normalizing constant is $Z = (2\pi)^{D/2} \det(\mathbf{H})^{-1/2} \exp(-\mathcal{U}(\hat{\boldsymbol{\theta}}))$.

Practical Considerations

1. The mode $\hat{\theta}$ can be found by gradient-based optimization

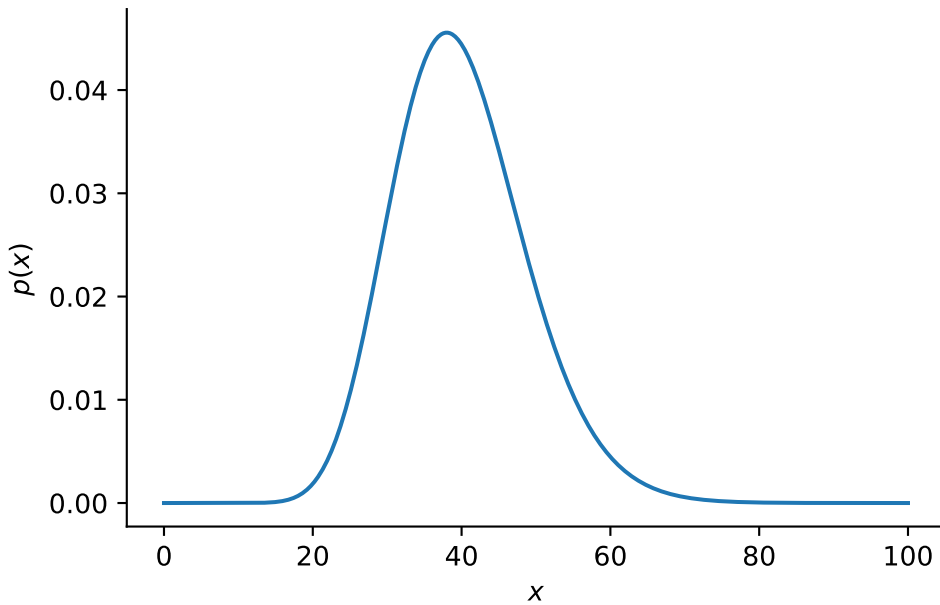
$$\theta \leftarrow \theta - \eta \nabla \mathcal{U}(\theta) \quad \text{until convergence}$$

2. The exact Hessian \mathbf{H} can be expensive to compute and invert for large D .
 - Use a diagonal approximation of the Hessian or a block-diagonal approximation.
3. The Laplace approximation is only valid in a neighborhood of the mode.

Laplace Approximation: Example

Consider a Gamma distribution with shape parameter α and rate parameter β .

$$p(x \mid \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} \exp(-\beta x)$$



Laplace Approximation: Example

1. Write the negative log posterior $\mathcal{U}(x)$.

$$\mathcal{U}(x) = -(\alpha - 1) \log x + \beta x$$

2. Compute the mode \hat{x} by setting $\frac{d\mathcal{U}(x)}{dx} = 0$.

$$\frac{d\mathcal{U}(x)}{dx} = -\frac{\alpha - 1}{x} + \beta = 0 \quad \Rightarrow \quad \hat{x} = \frac{\alpha - 1}{\beta}$$

3. Compute the Hessian $\frac{d^2\mathcal{U}(x)}{dx^2}$ in $x = \hat{x}$.

$$\frac{d^2\mathcal{U}(x)}{dx^2} = \frac{\alpha - 1}{x^2} \Big|_{x=\hat{x}} = \frac{\alpha - 1}{(\alpha - 1)^2 / \beta^2} = \frac{\beta^2}{\alpha - 1}$$

4. Set up the Gaussian approximation of the posterior.

$$p(x \mid \alpha, \beta) \approx \mathcal{N}\left(x \mid \frac{\alpha - 1}{\beta}, \frac{\alpha - 1}{\beta^2}\right)$$

Final considerations

Benefits

1. It is relatively simple to implement and understand.

Limitations

1. It only captures local curvature around the mode $\hat{\theta}$.
2. It assumes the posterior is approximately Gaussian, which may not hold in practice (e.g., multimodal posteriors).
3. It requires expensive computations:
 - Finding the mode $\hat{\theta}$ via optimization (if not available in closed form).
 - Computing and inverting the Hessian \mathbf{H} to get the covariance (at least cubic in the number of parameters).

i In practice

1. Laplace approximation is useful if a “pre-trained” model is available (e.g., a neural network shared on Huggingface) and we want to perform Bayesian inference around the learned parameters.
2. The Hessian is often approximated with a diagonal or block-diagonal matrix to reduce computational costs.