

# Laplace Approximation

---

Zeel B Patel, Nipun Batra

August 27, 2023

IIT Gandhinagar

# Outline

Taylor Series Expansion

ND Taylor Series

Laplace Approximation



Brook Taylor



Pierre-Simon Laplace

## Overall idea

- Posterior distribution  $p(\boldsymbol{\theta}|\mathcal{D})$  might be intractable but we can compute the MAP estimate.
- We know that posterior would be in form:  $p(\boldsymbol{\theta}|\mathcal{D}) = \frac{1}{Z}p(\mathcal{D}, \boldsymbol{\theta})$ , where  $Z$  is the normalizing constant.
- We can approximate this posterior using Taylor series expansion around the MAP estimate and it turns out that, after making a few assumptions, the resulting distribution is a Gaussian:  
 $p(\boldsymbol{\theta}|\mathcal{D}) \approx \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\theta}_{MAP}, H^{-1})$ , where  $H$  is the Hessian matrix of the log joint evaluated at  $\boldsymbol{\theta}_{MAP}$ .

# Taylor Series Expansion

---

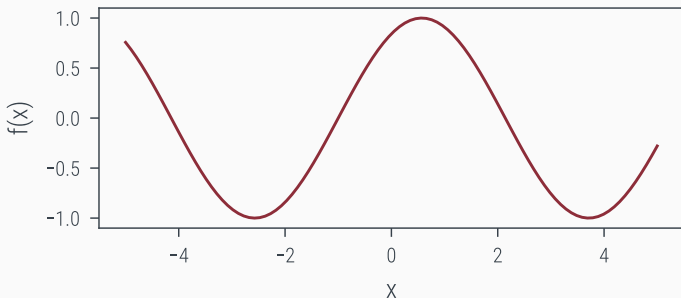
# 1D Taylor Series

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \frac{f'''(x_0)}{3!}(x - x_0)^3 + \dots$$

# Taylor Approximation of a 1D Function

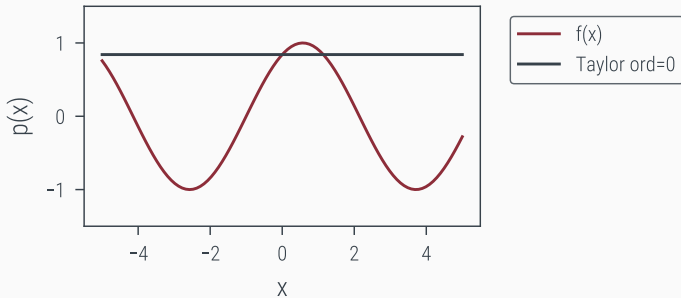
Consider the following function:

$$f(x) = \sin(1 + x)$$



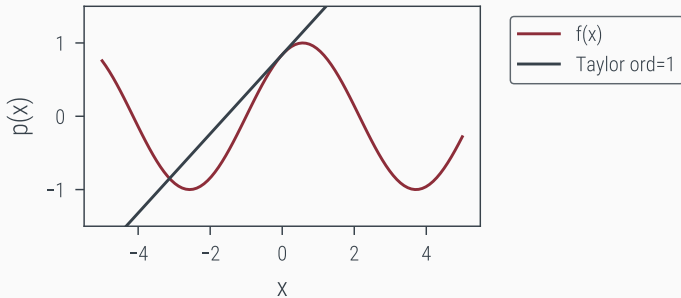
# Taylor Approximation of a 1D Function

Taylor approximation at  $x_0 = 0$ :



# Taylor Approximation of a 1D Function

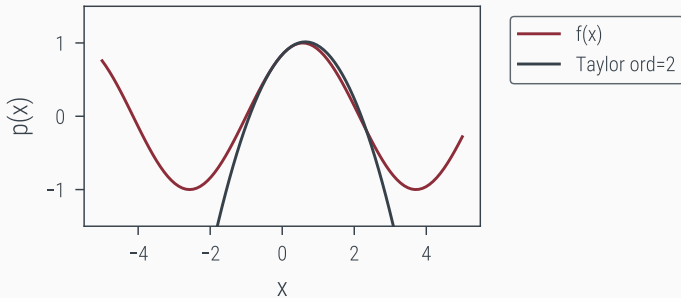
Taylor approximation at  $x_0 = 0$ :





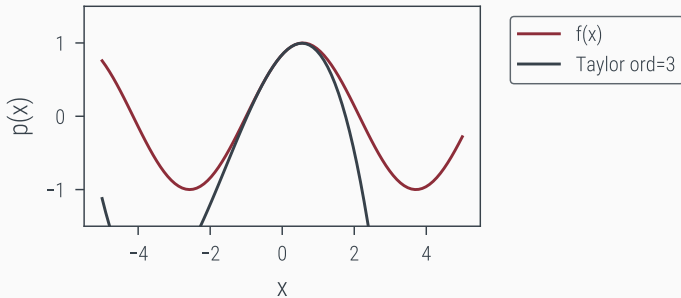
# Taylor Approximation of a 1D Function

Taylor approximation at  $x_0 = 0$ :



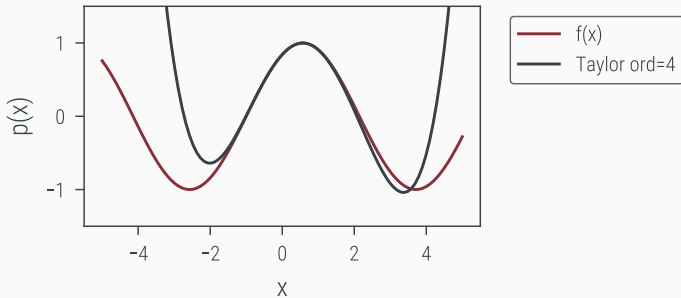
# Taylor Approximation of a 1D Function

Taylor approximation at  $x_0 = 0$ :



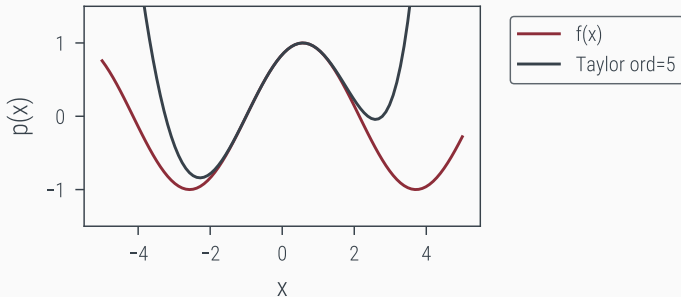
# Taylor Approximation of a 1D Function

Taylor approximation at  $x_0 = 0$ :



# Taylor Approximation of a 1D Function

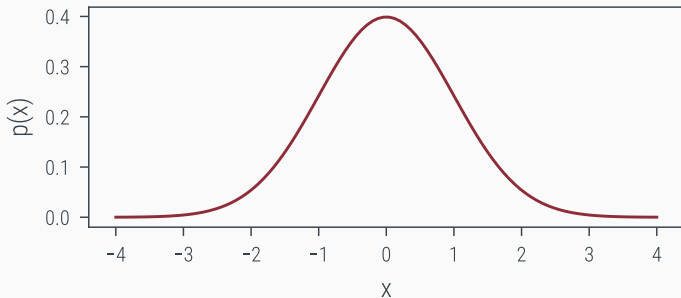
Taylor approximation at  $x_0 = 0$ :



# Taylor Approximation of a 1D Gaussian Function

Consider the standard normal distribution:  $p(x) \sim \mathcal{N}(x|0,1)$

$$p(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$



## Taylor Approximation of a 1D Gaussian Function

We will try to approximate log of the above function using Taylor series expansion around  $x_0 = 0$ .

## Taylor Approximation of a 1D Gaussian Function

We will try to approximate log of the above function using Taylor series expansion around  $x_0 = 0$ .

$$f(x) = \log p(x) = \log \left( \frac{1}{\sqrt{2\pi}} \right) - \frac{x^2}{2}$$

## Taylor Approximation of a 1D Gaussian Function

We will try to approximate log of the above function using Taylor series expansion around  $x_0 = 0$ .

$$f(x) = \log p(x) = \log \left( \frac{1}{\sqrt{2\pi}} \right) - \frac{x^2}{2}$$

$$f(0) = \log \left( \frac{1}{\sqrt{2\pi}} \right)$$



## Taylor Approximation of a 1D Gaussian Function

We will try to approximate log of the above function using Taylor series expansion around  $x_0 = 0$ .

$$f(x) = \log p(x) = \log \left( \frac{1}{\sqrt{2\pi}} \right) - \frac{x^2}{2}$$

$$f(0) = \log \left( \frac{1}{\sqrt{2\pi}} \right)$$

$$f'(0) \cdot x = -x_0 \cdot x = 0$$

## Taylor Approximation of a 1D Gaussian Function

We will try to approximate log of the above function using Taylor series expansion around  $x_0 = 0$ .

$$f(x) = \log p(x) = \log \left( \frac{1}{\sqrt{2\pi}} \right) - \frac{x^2}{2}$$

$$f(0) = \log \left( \frac{1}{\sqrt{2\pi}} \right)$$

$$f'(0) \cdot x = -x_0 \cdot x = 0$$

$$f''(0) \cdot \frac{x^2}{2} = -1 \cdot \frac{x^2}{2}$$

## Taylor Approximation of a 1D Gaussian Function

We will try to approximate log of the above function using Taylor series expansion around  $x_0 = 0$ .

$$f(x) = \log p(x) = \log \left( \frac{1}{\sqrt{2\pi}} \right) - \frac{x^2}{2}$$

$$f(0) = \log \left( \frac{1}{\sqrt{2\pi}} \right)$$

$$f'(0) \cdot x = -x_0 \cdot x = 0$$

$$f''(0) \cdot \frac{x^2}{2} = -1 \cdot \frac{x^2}{2}$$

$$\text{Taylor approximated } \tilde{f}(x) = \log \left( \frac{1}{\sqrt{2\pi}} \right) - \frac{x^2}{2}$$

## ND Taylor Series

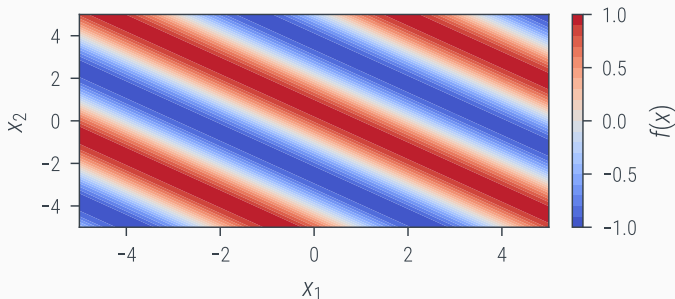
---

$$f(\mathbf{x}) = f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^T \nabla^2 f(\mathbf{x}_0) (\mathbf{x} - \mathbf{x}_0) + \dots$$

# Approximate a 2d function

We take the following function:

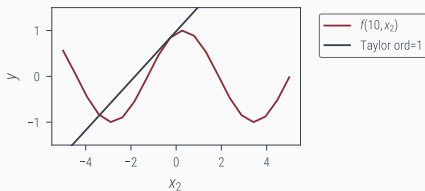
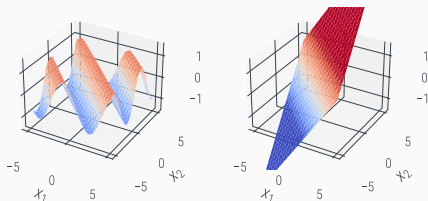
$$f(x_1, x_2) = \sin(1 + x_1 + x_2)$$



# Approximate a 2d function

Taylor approximation at  $x_0 = (0, 0)$ :

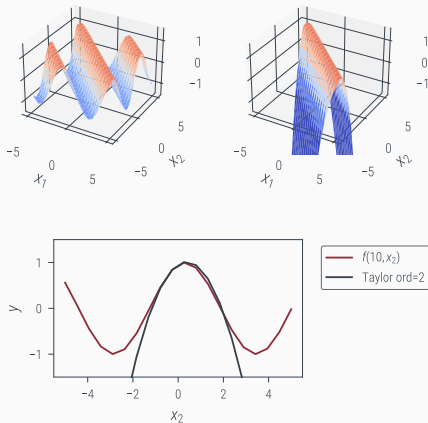
Taylor ord=1



# Approximate a 2d function

Taylor approximation at  $x_0 = (0, 0)$ :

Taylor ord=2





# Laplace Approximation

---

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{1}{Z}p(\mathcal{D}, \boldsymbol{\theta})$$

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{1}{Z} p(\mathcal{D}, \boldsymbol{\theta})$$

We can rewrite this as:

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{1}{Z} e^{-f(\boldsymbol{\theta})}$$
$$f(\boldsymbol{\theta}) = -\log p(\mathcal{D}, \boldsymbol{\theta})$$

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{1}{Z}p(\mathcal{D}, \boldsymbol{\theta})$$

We can rewrite this as:

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{1}{Z}e^{-f(\boldsymbol{\theta})}$$
$$f(\boldsymbol{\theta}) = -\log p(\mathcal{D}, \boldsymbol{\theta})$$

Note that  $f(\boldsymbol{\theta})$  is the negative log joint which is used as a loss function to estimate  $\boldsymbol{\theta}_{MAP}$ .

# Laplace Approximation

- Highest mass is concentrated around  $\theta_{MAP}$  and hence it makes sense to get Taylor approximation around that point.

# Laplace Approximation

- Highest mass is concentrated around  $\theta_{MAP}$  and hence it makes sense to get Taylor approximation around that point.
- In other words, if our approximation is bad where we have low probability mass, it doesn't matter much.

# Laplace Approximation

- Highest mass is concentrated around  $\theta_{MAP}$  and hence it makes sense to get Taylor approximation around that point.
- In other words, if our approximation is bad where we have low probability mass, it doesn't matter much.
- Thus, we expand  $f(\theta)$  around  $\theta_{MAP}$  using Taylor series expansion up to second derivative:

$$\begin{aligned} f(\theta) \approx & f(\theta_{MAP}) + \nabla f(\theta_{MAP})^T (\theta - \theta_{MAP}) \\ & + \frac{1}{2} (\theta - \theta_{MAP})^T \nabla^2 f(\theta_{MAP}) (\theta - \theta_{MAP}) \end{aligned}$$

$$\begin{aligned} f(\boldsymbol{\theta}) &\approx f(\boldsymbol{\theta}_{MAP}) + \nabla f(\boldsymbol{\theta}_{MAP})^T (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP}) \\ &\quad + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP})^T \nabla^2 f(\boldsymbol{\theta}_{MAP}) (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP}) \end{aligned}$$



# Laplace Approximation

$$\begin{aligned} f(\boldsymbol{\theta}) &\approx f(\boldsymbol{\theta}_{MAP}) + \nabla f(\boldsymbol{\theta}_{MAP})^T (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP}) \\ &\quad + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP})^T \nabla^2 f(\boldsymbol{\theta}_{MAP}) (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP}) \end{aligned}$$

Since,  $\boldsymbol{\theta}_{MAP}$  is minima of  $f(\boldsymbol{\theta})$ ,  $\nabla f(\boldsymbol{\theta}_{MAP}) = 0$ .

# Laplace Approximation

$$\begin{aligned} f(\boldsymbol{\theta}) &\approx f(\boldsymbol{\theta}_{MAP}) + \nabla f(\boldsymbol{\theta}_{MAP})^T (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP}) \\ &\quad + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP})^T \nabla^2 f(\boldsymbol{\theta}_{MAP}) (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP}) \end{aligned}$$

Since,  $\boldsymbol{\theta}_{MAP}$  is minima of  $f(\boldsymbol{\theta})$ ,  $\nabla f(\boldsymbol{\theta}_{MAP}) = 0$ .

$$f(\boldsymbol{\theta}) \approx f(\boldsymbol{\theta}_{MAP}) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP})^T \nabla^2 f(\boldsymbol{\theta}_{MAP}) (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP})$$

# Laplace Approximation

$$\begin{aligned} f(\boldsymbol{\theta}) &\approx f(\boldsymbol{\theta}_{MAP}) + \nabla f(\boldsymbol{\theta}_{MAP})^T (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP}) \\ &\quad + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP})^T \nabla^2 f(\boldsymbol{\theta}_{MAP}) (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP}) \end{aligned}$$

Since,  $\boldsymbol{\theta}_{MAP}$  is minima of  $f(\boldsymbol{\theta})$ ,  $\nabla f(\boldsymbol{\theta}_{MAP}) = 0$ .

$$\begin{aligned} f(\boldsymbol{\theta}) &\approx f(\boldsymbol{\theta}_{MAP}) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP})^T \nabla^2 f(\boldsymbol{\theta}_{MAP}) (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP}) \\ &= f(\boldsymbol{\theta}_{MAP}) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP})^T H (\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP}) \end{aligned}$$

where  $H$  is the Hessian matrix of  $f(\boldsymbol{\theta})$  evaluated at  $\boldsymbol{\theta}_{MAP}$ .

# Laplace Approximation

Plugging this back to the posterior equation:

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{1}{Z} e^{-f(\boldsymbol{\theta})}$$

# Laplace Approximation

Plugging this back to the posterior equation:

$$\begin{aligned} p(\boldsymbol{\theta}|\mathcal{D}) &= \frac{1}{Z} e^{-f(\boldsymbol{\theta})} \\ &\approx \frac{1}{Z} e^{-f(\boldsymbol{\theta}_{MAP})} e^{-\frac{1}{2}(\boldsymbol{\theta}-\boldsymbol{\theta}_{MAP})^T H(\boldsymbol{\theta}-\boldsymbol{\theta}_{MAP})} \end{aligned}$$

# Laplace Approximation

Plugging this back to the posterior equation:

$$\begin{aligned} p(\boldsymbol{\theta}|\mathcal{D}) &= \frac{1}{Z} e^{-f(\boldsymbol{\theta})} \\ &\approx \frac{1}{Z} e^{-f(\boldsymbol{\theta}_{MAP})} e^{-\frac{1}{2}(\boldsymbol{\theta}-\boldsymbol{\theta}_{MAP})^T H(\boldsymbol{\theta}-\boldsymbol{\theta}_{MAP})} \\ &\sim \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\theta}_{MAP}, H^{-1}) \end{aligned}$$

# Laplace Approximation

Plugging this back to the posterior equation:

$$\begin{aligned} p(\boldsymbol{\theta}|\mathcal{D}) &= \frac{1}{Z} e^{-f(\boldsymbol{\theta})} \\ &\approx \frac{1}{Z} e^{-f(\boldsymbol{\theta}_{MAP})} e^{-\frac{1}{2}(\boldsymbol{\theta}-\boldsymbol{\theta}_{MAP})^T H(\boldsymbol{\theta}-\boldsymbol{\theta}_{MAP})} \\ &\sim \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\theta}_{MAP}, H^{-1}) \\ Z &= p(\mathcal{D}, \boldsymbol{\theta}_{MAP}) \cdot (2\pi)^{D/2} \cdot |H|^{-\frac{1}{2}} \end{aligned}$$

# Laplace Approximation

Plugging this back to the posterior equation:

$$\begin{aligned} p(\boldsymbol{\theta}|\mathcal{D}) &= \frac{1}{Z} e^{-f(\boldsymbol{\theta})} \\ &\approx \frac{1}{Z} e^{-f(\boldsymbol{\theta}_{MAP})} e^{-\frac{1}{2}(\boldsymbol{\theta}-\boldsymbol{\theta}_{MAP})^T H(\boldsymbol{\theta}-\boldsymbol{\theta}_{MAP})} \\ &\sim \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\theta}_{MAP}, H^{-1}) \\ Z &= p(\mathcal{D}, \boldsymbol{\theta}_{MAP}) \cdot (2\pi)^{D/2} \cdot |H|^{-\frac{1}{2}} \end{aligned}$$

Note that this result is not specific to Bayesian inference and can be used to approximate any intractable function.



# Pros and Cons of Laplace Approximation

- Pros:
  - Simple to implement
  - Computationally efficient
  - Can be used to approximate any intractable function

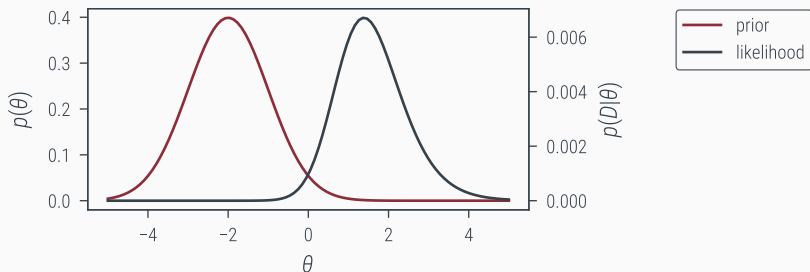
# Pros and Cons of Laplace Approximation

- Pros:
  - Simple to implement
  - Computationally efficient
  - Can be used to approximate any intractable function
- Cons:
  - It can give bad approximation when posterior is not unimodal
  - Gaussian assumption can be too restrictive at times
  - Hessian matrix inversion can be numerically unstable and expensive. A diagonal or block-wise approximation can be applied to resolve this. Checkout [Laplace-Redux](#) for more details.

# Normal Prior for Coin Toss

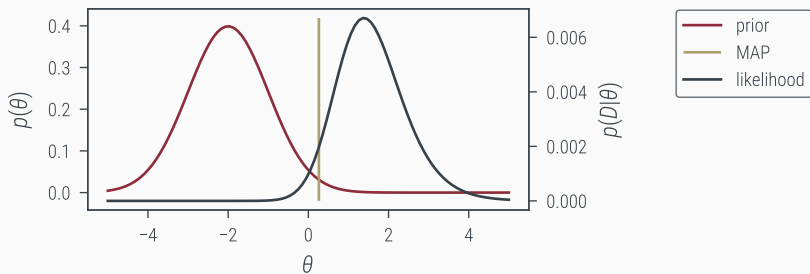
Consider the following coin toss experiment scenario:

- $\mathcal{D} = \{1, 1, 1, 1, 1, 1, 1, 0, 0\}$
- $p(\theta) = \mathcal{N}(\theta | -2, 1)$
- $h = \sigma(\theta)$
- $p(y|\theta) = h^y(1 - h)^{1-y}$



# Normal Prior for Coin Toss

First, we find the MAP estimate.



# Normal Prior for Coin Toss

Now, according to the Laplace Approximation, the posterior is:

