# Laplace Approximation

CS698X: Topics in Probabilistic Modeling and Inference

Piyush Rai

# Laplace Approximation of Posterior Distribution

- Consider a posterior distribution that is intractable to compute

Unknowns of the model

Data

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

- Laplace approximation approximates the above using a Gaussian distribution

$$p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta_{MAP}, \mathbf{H}^{-1})$$

Laplace Approx. Gaussian

Target posterior

$\theta_{MAP}$

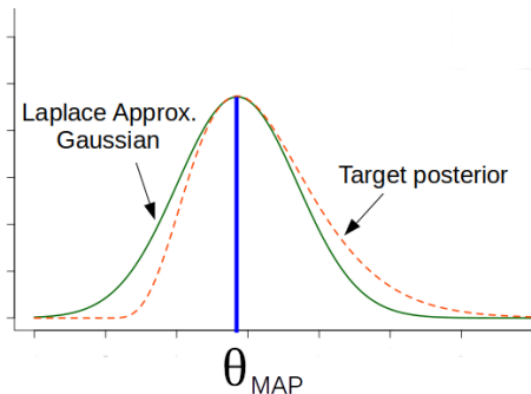$$\theta_{MAP} = \arg\max_{\theta} p(\theta|\mathcal{D}) = \arg\max_{\theta} p(\mathcal{D},\theta)$$

$$= \arg\max_{\theta} p(\mathcal{D}|\theta)p(\theta)$$

$$= \arg\max_{\theta} [\log p(\mathcal{D}|\theta) + \log p(\theta)]$$

$$\mathbf{H} = -\nabla^2 \log p(\theta|\mathcal{D})\big|_{\theta=\theta_{MAP}} = -\nabla^2 \log p(\mathcal{D},\theta)\big|_{\theta=\theta_{MAP}}$$

$$= -\nabla^2 [\log p(\mathcal{D}|\theta) + \log p(\theta)]\big|_{\theta=\theta_{MAP}}$$

- Why is the above Gaussian a reasonable approximation to the posterior?

# Derivation of the Laplace Approximation

$-\mathbf{H}$ Recall that Hessian is the second derivative of the negative of log-joint

$$\approx \frac{1}{2}(\theta - \theta_{MAP})^\top \boxed{\nabla^2 \log p(\mathcal{D}, \theta_{MAP})}(\theta - \theta_{MAP}) + \text{const}$$

- Let's write the Bayes rule as

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}, \theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D}, \theta)}{\int p(\mathcal{D}, \theta)d\theta} = \frac{e^{\log p(\mathcal{D}, \theta)}}{\int e^{\log p(\mathcal{D}, \theta)}d\theta}$$

Aha! This is a Gaussian!

Comparing with a Gaussian PDF
Mean = $\theta_{MAP}$
Cov. Matrix = $\mathbf{H}^{-1}$

- Approximating $\log p(\mathcal{D}, \theta)$ by a quadratic function of $\theta$ will make it a Gaussian

- Consider the second-order Taylor approx of a function $f(\theta)$ around some $\theta_0$

$$f(\theta) \approx f(\theta_0) + (\theta - \theta_0)^\top \nabla f(\theta_0) + \frac{1}{2}(\theta - \theta_0)^\top \nabla^2 f(\theta_0)(\theta - \theta_0)$$

- Assuming $f(\theta) = \log p(\mathcal{D}, \theta)$ and $\theta_0 = \theta_{MAP}$, $\nabla f(\theta_{MAP}) = \nabla \log p(\mathcal{D}, \theta_{MAP}) = 0$

Constant w.r.t. $\theta$

$$\log p(\mathcal{D}, \theta) \approx \log p(\mathcal{D}, \theta_{MAP}) + \frac{1}{2}(\theta - \theta_{MAP})^\top \nabla^2 \log p(\mathcal{D}, \theta_{MAP})(\theta - \theta_{MAP})$$
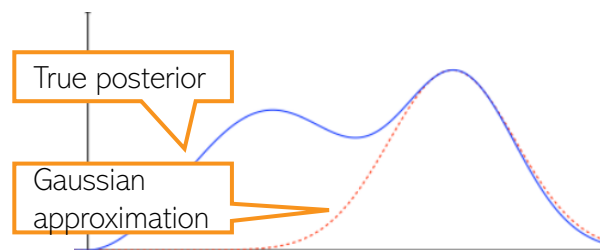
- Thus Laplace approx. is based on a second-order Taylor approx. of the posterior

# Properties of Laplace Approximation

- Usually straightforward if derivatives (first and second) can be computed easily

- Expensive if parameter $\boldsymbol{\theta}$ is very high dimensional

  > E.g., a deep neural network, or even in simpler models (e.g., logistic reg with a very large number of features

  - Reason: We need to invert the Hessian whose size is $D \times D$ ($D$ is the # of params)

- Can do badly if the (true) posterior is multimodal

  True posterior

  Gaussian approximation

- Applicable only when $\boldsymbol{\theta}$ is real-valued (won't if, say, it is positive, binary etc)

- Note: Even if we have a <u>non-probabilistic</u> model (loss function + regularization), we can obtain an approx "posterior" for that model using the Laplace approximation

  - Optima of the regularized loss function will be Gaussian's mean
  - Second derivative of the regularized loss function will be the Hessian

# Laplace Approx. for High-Dimensional Problems

- When $\boldsymbol{\theta}$ is very high dim, one option is to approximate the Hessian itself

- One such approx. of the Hessian is a diagonal approximation

Assuming a discriminative model with parameters $\boldsymbol{\theta}$

Fisher Information Matrix (FIM)

$$\mathbf{H} \approx \mathrm{diag}(\mathbf{F})$$

$$\mathbf{F} = \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x},\mathbf{y})} \left[ \nabla \log p(\mathbf{y}|\mathbf{x},\boldsymbol{\theta}) \nabla \log p(\mathbf{y}|\mathbf{x},\boldsymbol{\theta})^\top \right]$$
$$\approx \mathbb{E}_{p_D(\mathbf{x},\mathbf{y})} \left[ \nabla \log p(\mathbf{y}|\mathbf{x},\boldsymbol{\theta}) \nabla \log p(\mathbf{y}|\mathbf{x},\boldsymbol{\theta})^\top \right]$$
$$= \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x},\mathbf{y}) \in \mathcal{D}} \nabla \log p(\mathbf{y}|\mathbf{x},\boldsymbol{\theta}) \nabla \log p(\mathbf{y}|\mathbf{x},\boldsymbol{\theta})^\top$$

Example: A Bayesian neural net for regression/classification ($\boldsymbol{\theta}$ denotes the weights of the network)

FIM is easily computable in auto-diff frameworks used in deep learning

- The diagonal approx. of Hessian may be too crude ☹
  - Ignores covariances among params and treats them as being independent of each other

- A block-diagonal approx. proposed recently (in the context of deep neural nets)
  - Treats params across layers to be independent but correlated within the same layer
  - The approach known as Kronecker-Product Factored (KFAC) Laplace approximation

# Coming Up

- Generalized Linear Models (GLM)
  - Models of the form $p(y|x)$ where $p(y|x)$ is some exponential family distribution
  - Note: Prob. linear regression and logistic regression were also examples of GLMs
- Generative Classification