

Sampling (wrap-up)

CS772A: Probabilistic Machine Learning

Piyush Rai

Plan Today

- Wrapping up sampling methods and posterior approximation
 - Hamiltonian Monte Carlo
 - Some other methods for posterior approximation for Bayesian neural nets
 - Parallel/Distributed MCMC



Hamiltonian/Hybrid Monte Carlo (HMC)

- HMC (Neal, 1996) is an “auxiliary variable sampler” and incorporates gradient info
- Uses the idea of simulating a Hamiltonian Dynamics of a physical system
- Consider the target posterior $p(\theta|\mathcal{D}) \propto \exp(-U(\theta))$
- Think of θ as “position” then $U(\theta) = -\log p(\mathcal{D}|\theta)p(\theta)$ is like “potential energy”
- Let’s introduce an auxiliary variable - the momentum \mathbf{r} of the system
- Can now define a joint distribution over the position and momentum as

$$p(\theta, \mathbf{r}|\mathcal{D}) \propto \exp\left(-U(\theta) - \frac{1}{2}\mathbf{r}^\top M^{-1}\mathbf{r}\right) \propto p(\theta|\mathcal{D})p(\mathbf{r})$$

- The total energy (potential + kinetic) or the Hamiltonian of the system

Constant w.r.t. time $\rightarrow H(\theta, \mathbf{r}) = U(\theta) - \frac{1}{2}\mathbf{r}^\top M^{-1}\mathbf{r} = U(\theta) + K(\mathbf{r})$

- Given a sample (θ, \mathbf{r}) from $p(\theta, \mathbf{r})$, ignoring \mathbf{r} , θ will be a sample from $p(\theta|\mathcal{D})$



Generating Samples in HMC

- Given an initial (θ, \mathbf{r}) , Hamiltonian Dynamics defines how (θ, \mathbf{r}) changes w.r.t. time t

$$\begin{aligned} \frac{\partial \theta}{\partial t} &= \frac{\partial H}{\partial \mathbf{r}} = \frac{\partial K}{\partial \mathbf{r}} \\ \frac{\partial \mathbf{r}}{\partial t} &= -\frac{\partial H}{\partial \theta} = -\frac{\partial U}{\partial \theta} \end{aligned} \quad (H(\theta, \mathbf{r}) = U(\theta) + \frac{1}{2} \mathbf{r}^\top M^{-1} \mathbf{r} = U(\theta) + K(\mathbf{r}))$$

- We can use these equations to update $(\theta, \mathbf{r}) \rightarrow (\theta^*, \mathbf{r}^*)$ by discretizing time

- For $s = 1:S$, sample as follows

- Initialize $\theta_0 = \theta^{(s-1)}$, $\mathbf{r}_* \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\mathbf{r}_0 = \mathbf{r}_* - \frac{\rho}{2} \frac{\partial U}{\partial \theta} |_{\theta_0}$
- Do L “leapfrog” steps with learning rates $\rho_\ell = \rho$ for $\ell < L$ and $\rho_L = \rho/2$
 - For $\ell = 1:L$

$$\theta_\ell = \theta_{\ell-1} + \rho \frac{\partial K}{\partial \mathbf{r}} |_{\mathbf{r}_{\ell-1}}$$

$$\mathbf{r}_\ell = \mathbf{r}_{\ell-1} - \rho_\ell \frac{\partial U}{\partial \theta} |_{\theta_\ell}$$

- Perform MH accept/reject test on (θ_L, \mathbf{r}_L) . If accepted $\theta^{(s)} = \theta_L$

- The momentum forces exploring different regions instead of getting driven to regions where the MAP solution is

Reason: Getting analytical solutions for the above requires integrals which is in general intractable

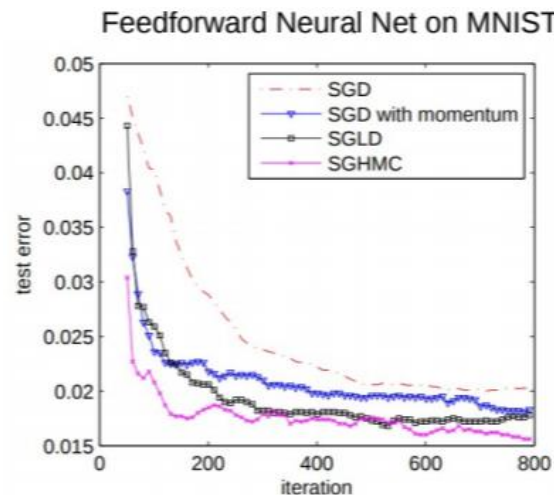
L usually set to 5 and learning rate tuned to make acceptance rate around 90%

A single sample generated by taking L steps



HMC in Practice

- HMC typically has very low rejection rate (that too, primarily due to discretization error)
- Performance can be sensitive to L (no. of leapfrog steps) and step-sizes, tuning hard
- A lot of renewed interest in HMC (you may check out NUTS - No U-turn Sampler – doesn't require setting L)
 - Prob. Prog. packages e.g., Tensorflow Prob., Stan, etc, contain implementations of HMC
- Can also do HMC on minibatches (Stochastic Gradient HMC - Chen et al, 2014)
- An illustration: SGHMC vs other methods on MNIST classification (Bayesian neural net)



(Figure: Stochastic Gradient Hamiltonian Monte Carlo (Chen et al, 2014))

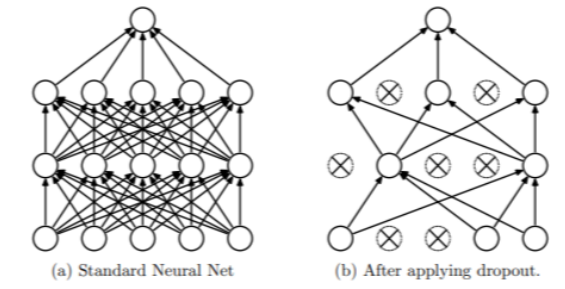


Efficient Posterior approx. for Bayesian Neural Nets⁶

- Monte Carlo Dropout is another popular and efficient way

- Standard Dropout

- Drop some weights randomly (with some “drop” probability) during training
- At test time, multiply each weight by the “keep” probability
- Note: Dropout applied only at training time



- Monte Carlo Dropout* (dropout also at test time)

Can be seen as learning a variational approximation of the weights (see paper for details, if interested)

$$p(y_* | x_*, \mathcal{D}) \approx \frac{1}{S} \sum_{s=1}^S p(y_* | x_*, \theta^{(s)})$$



$$p(y_* | x_*, \mathcal{D}) \approx \frac{1}{S} \sum_{s=1}^S p(y_* | x_*, \theta^{(s)})$$

where $\theta^{(s)} \sim p(\theta | \mathcal{D})$

where $\theta^{(s)} = \epsilon^{(s)} \odot \hat{\theta}$

Vector of Bernoulli or Gaussian noise

Elementwise product

Point estimate



*Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning (Gal and Ghahramani, 2016)

Efficient Posterior approx. for Bayesian Neural Nets⁷

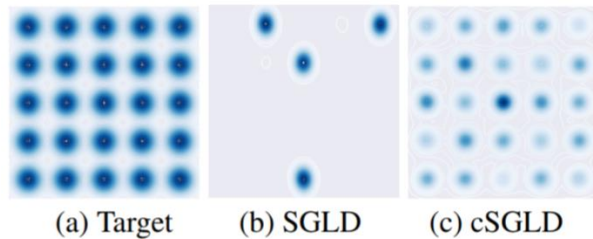
- SGMCMC methods like SGLD and SGHMC are also used nowadays (very efficient)

$$\theta^{(t)} = \theta^{(t-1)} + \frac{\eta_t}{2} \nabla_{\theta} [\log p(\mathcal{D}|\theta) + \log p(\theta)]|_{\theta^{(t-1)}} + \epsilon_t$$

- Recently, SGMCMC with **cyclic step sizes (cSGLD)** was proposed (Zhang et al, 2020)
 - Use big steps to explore different modes
 - Use small steps later to sample once a mode is localized

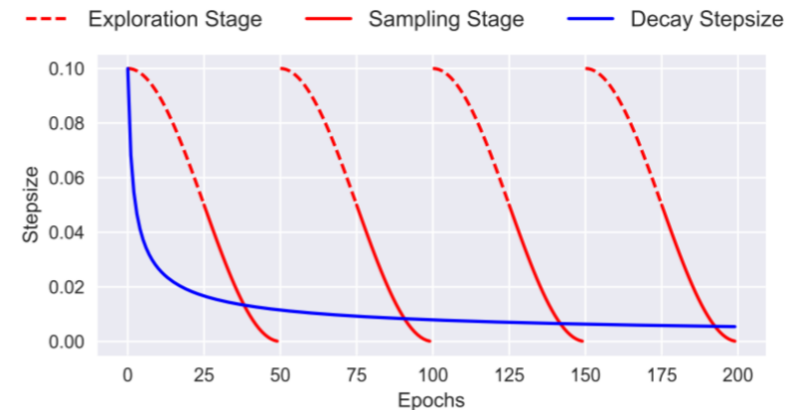
Step size in iteration k

$$\alpha_k = \frac{\alpha_0}{2} \left[\cos \left(\frac{\pi \text{mod}(k-1, \lceil K/M \rceil)}{\lceil K/M \rceil} \right) + 1 \right]$$



A complex mixture of Gaussian distributions

K is the total number of iterations and M is the number of cycles



	CIFAR-10	CIFAR-100
SGD	5.29±0.15	23.61±0.09
SGDM	5.17±0.09	22.98±0.27
Snapshot-SGD	4.46±0.04	20.83±0.01
Snapshot-SGDM	4.39±0.01	20.81±0.10
SGLD	5.20±0.06	23.23±0.01
cSGLD	4.29±0.06	20.55±0.06
SGHMC	4.93±0.1	22.60±0.17
cSGHMC	4.27±0.03	20.50±0.11



Deep Ensembles

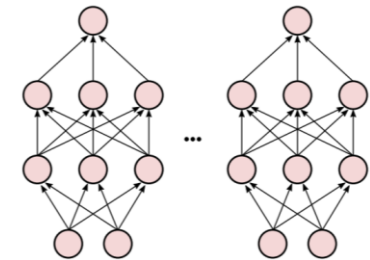
- Most inference methods tend to produce local approximations only
 - VI methods typically learn an approximation around one of the modes
 - Sampling methods may give most samples near one of the modes (though in principle they may explore other modes as well)
 - Thus the uncertainties may be underestimated in general

Both VI and Sampling may be prone to capturing only a single "Basin of attraction"

- Deep Ensembles* is a method that tries to address this issue
 - Train the network M times with different seeds and permutations of training data
 - Denote the learned weights by $\theta_1, \theta_2, \dots, \theta_M$ (assuming these are M modes)
 - Approximate the posterior by the following

$$p(\theta|\mathcal{D}) = \frac{1}{M} \sum_{m=1}^M \delta_{\theta_m}(\theta)$$

Akin to Bayesian Model Averaging using M models



- This approach is considered non-Bayesian but often performs better (in terms of more diversity in the set of parameters learned) than other inference methods



Parallel/Distributed MCMC

- Suppose our goal is to compute the posterior of $\theta \in \mathbb{R}^D$ (assuming N is very large)

$$p(\theta|\mathbf{X}) \propto p(\theta)p(\mathbf{X}|\theta) = p(\theta) \prod_{n=1}^N p(\mathbf{x}_n|\theta)$$

- Suppose we have J machines with data partitioned as $\mathbf{X} = \{\mathbf{X}^{(j)}\}_{j=1}^J$
- Let's assume that the posterior $p(\theta|\mathbf{X})$ factorizes as

$$p(\theta|\mathbf{X}) = \prod_{i=1}^J p^{(j)}(\theta|\mathbf{X}^{(j)})$$

- Here $p^{(j)}(\theta|\mathbf{X}^{(j)}) \propto p(\theta)^{1/J} \prod_{\mathbf{x}_n \in \mathbf{X}^{(j)}} p(\mathbf{x}_n|\theta)$ is known as the “subset posterior”
- Assume the j^{th} machine generates T MCMC samples $\{\theta_{j,t}\}_{t=1}^T$
- We need a way to combine these subset posteriors using a “consensus”

$$\hat{\theta}_1, \dots, \hat{\theta}_T = \text{CONSENSUSSAMPLES}(\{\theta_{j,1}, \dots, \theta_{j,T}\}_{j=1}^J)$$



Parallel/Distributed MCMC

- Many ways to compute the consensus samples. Let's look at two of them
- Approach 1: **Weighted Average**: $\hat{\theta}_t = \sum_{j=1}^J W_j \theta_{j,t}$ where W_j can be learned as follows
 - Assuming Gaussian likelihood and Gaussian prior on θ

$$\begin{aligned}\bar{\Sigma}_j &= \text{sample covariance of } \{\theta_{j,1}, \dots, \theta_{j,T}\} \\ \Sigma &= (\Sigma_0^{-1} + \sum_{j=1}^J \bar{\Sigma}_j^{-1})^{-1} \quad (\Sigma_0 \text{ is the prior's covariance}) \\ W_j &= \Sigma(\Sigma_0^{-1}/J + \bar{\Sigma}_j^{-1})\end{aligned}$$

These approaches can also be used to make VI parallel/distributed



- Approach 2: Fit J Gaussians, one for each $\{\theta_{j,t}\}_{t=1}^T$ and take their product

$$\begin{aligned}\bar{\mu}_j &= \text{sample mean of } \{\theta_{j,1}, \dots, \theta_{j,T}\}, \quad \bar{\Sigma}_j = \text{sample covariance of } \{\theta_{j,1}, \dots, \theta_{j,T}\} \\ \hat{\Sigma}_J &= (\sum_{j=1}^J \bar{\Sigma}_j^{-1})^{-1}, \quad \hat{\mu}_J = \hat{\Sigma}_J (\sum_{j=1}^J \bar{\Sigma}_j^{-1} \bar{\mu}_j) \quad (\text{cov and mean of prod. of Gaussians}) \\ \hat{\theta}_t &\sim \mathcal{N}(\hat{\mu}_J, \hat{\Sigma}_J), t = 1, \dots, T \quad (\text{the final consensus samples})\end{aligned}$$

- For detailed proofs and other approaches, may refer to the reference below

Approximate Inference: VI vs Sampling

- VI approximates a posterior distribution $p(\mathbf{Z}|\mathbf{X})$ by another distribution $q(\mathbf{Z}|\phi)$
- Sampling uses S samples $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(S)}$ to approximate $p(\mathbf{Z}|\mathbf{X})$
- Sampling can be used within VI (ELBO approx using Monte-Carlo)
- In terms of “comparison” between VI and sampling, a few things to be noted
 - **Convergence:** VI only has local convergence, sampling (in theory) can give exact posterior
 - **Storage:** Sampling based approx needs to storage all samples, VI only needs var. params ϕ
 - **Prediction Cost:** Sampling always requires Monte-Carlo avg for posterior predictive; with VI, sometimes we can get closed form posterior predictive

PPD if using sampling:

$$p(x_*|X) = \int p(x_*|Z)p(Z|X)dZ \approx \frac{1}{S} \sum_{s=1}^S p(x_*|Z^{(s)})$$

Closed form if integral is tractable (otherwise Monte Carlo avg still needed for PPD)

PPD if using VI:

$$p(x_*|X) = \int p(x_*|Z)p(Z|X)dZ \approx \int p(x_*|Z)q(Z|\phi)dZ$$

Compressing the S samples into something more compact

- There is some work on “compressing” sampling-based approximations*

Inference Methods: Summary

- MLE/MAP: Straightforward for differentiable models (can even use automatic diff.)
- Conjugate models with one “main” parameter: Straightforward posterior updates
- MLE-II/MAP-II: Often useful for estimating the hyperparameters
- EM: If we want to do MLE/MAP for models with latent variables
 - Very general algorithm, can also be made online
 - Used when we want point estimates for some unknowns and posterior over others
 - Can use it for hyperparameter estimation as well
 - Often better than using direct gradient methods
- VI and sampling methods can be used to get full posterior for complex models
 - Quite easy if we have local conjugacy (VI has closed form updates, Gibbs sampler is easy to derive)
 - In other cases, we have general VI with Monte-Carlo gradients, MH sampling
 - MCMC can also make use of gradient info (LD/SGLD)
- For large-scale problems, online/distributed VI/MCMC, or SGD based posterior approx

