MCMC Sampling (wrap-up), Deep Generative Models

CS772A: Probabilistic Machine Learning Piyush Rai

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 <u>auxiliary variable</u> the momentum $m{r}$ of the system
- Can now define a joint distribution over the position and momentum as $p(\theta, \boldsymbol{r} | \mathcal{D}) \propto \exp\left(-U(\theta) \frac{1}{2}\boldsymbol{r}^{\mathsf{T}}M^{-1}\boldsymbol{r}\right) \propto p(\theta | \mathcal{D})p(\boldsymbol{r})$

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

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

Given a sample (θ, r) from $p(\theta, r)$, ignoring r, θ will be a sample from $p(\theta | D)$

Generating Samples in HMC

 $\frac{\partial \theta}{\partial t} = \frac{\partial H}{\partial r} = \frac{\partial K}{\partial r}$

• Given an initial (θ, r) , Hamiltonian Dynamics defines how (θ, r) changes w.r.t. time t

- $\frac{\partial \mathbf{r}}{\partial t} = -\frac{\partial H}{\partial \theta} = -\frac{\partial U}{\partial \theta}$ • We can use these equations to update $(\theta, \mathbf{r}) \to (\theta^*, \mathbf{r}^*)$ by <u>discretizing time</u>
- For s = 1: S, sample as follows

• For $\ell = 1:L$

Initialize $\theta_0 = \theta^{(s-1)}$, $\mathbf{r}_* \sim \mathcal{N}(0, \mathbf{I})$ and $\mathbf{r}_0 = \mathbf{r}_* - \frac{\rho}{2} \frac{\partial U}{\partial \theta}|_{\theta_0}$

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

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

• Do *L* "leapfrog" steps with learning rates $\rho_\ell = \rho$ for $\ell < L$ and $\rho_L = \rho/2$

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

 $(H(\theta, \mathbf{r}) = U(\theta) + \frac{1}{2}\mathbf{r}^{\top}M^{-1}\mathbf{r} = U(\theta) + K(\mathbf{r}))$

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

• Perform MH accept/reject test on (θ_L, r_L) . If accepted $\theta^{(s)} = \theta_L \stackrel{\text{A single sample generated}}{\longrightarrow}$ by taking L steps

The momentum forces exploring different regions instead of getting driven to regions where the MAP solution is
CS772A: PML

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))

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



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 heta

$$\begin{split} \bar{\Sigma}_{j} &= \text{ sample covariance of } \{\theta_{j,1}, \dots, \theta_{j,T}\} \\ \Sigma &= (\Sigma_{0}^{-1} + \sum_{j=1}^{J} \bar{\Sigma}_{j}^{-1})^{-1} \quad (\Sigma_{0} \text{ is the prior's covariance}) \\ W_{j} &= \Sigma(\Sigma_{0}^{-1}/J + \bar{\Sigma}_{j}^{-1}) \end{split}$$

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

$$\bar{\mu}_{j} = \text{sample mean of } \{\theta_{j,1}, \dots, \theta_{j,T}\}, \quad \bar{\Sigma}_{j} = \text{sample covariance of } \{\theta_{j,1}, \dots, \theta_{j,T}\}$$

$$\hat{\Sigma}_{J} = (\sum_{j=1}^{J} \bar{\Sigma}_{j}^{-1})^{-1}, \quad \hat{\mu}_{J} = \hat{\Sigma}_{J} (\sum_{j=1}^{J} \bar{\Sigma}_{j}^{-1} \bar{\mu}_{j}) \quad (\text{cov and mean of prod. of Gaussians})$$

$$\hat{\theta}_{t} \sim \mathcal{N}(\hat{\mu}_{J}, \hat{\Sigma}_{J}), t = 1, \dots, T \quad (\text{the final consensus samples})$$

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



Approximate Inference: VI vs Sampling

- VI approximates a posterior distribution p(Z|X) by another distribution $q(Z|\phi)$
- Sampling uses S samples $Z^{(1)}, Z^{(2)}, \dots, Z^{(S)}$ to approximate p(Z|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 <u>always</u> requires Monte-Carlo avging for posterior predictive; with VI, <u>sometimes</u> we can get closed form posterior predictive

PPD if using sampling:
$$p(x_*|X) = \int p(x_*|Z)p(Z|X)dZ$$

 $a \approx \frac{1}{S} \sum_{s=1}^{S} p(x_* | Z^{(s)})$ Closed form if integral is tractable (otherwise Monte Carlo avg still needed for PPD)

CS772A: PML

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

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

Latent Variable Models for Generation Tasks

• Assume a K-dim latent variable \boldsymbol{z}_n is transformed to generate to D-dim observation \boldsymbol{x}_n



- It is common to use a Gaussian prior for \boldsymbol{z}_n (though other priors can be used)
- If we use a neural net or GP, such models can generate very high-quality data
 - Take the trained network, generate a random $m{z}$ from prior, pass it through the model to generate $m{x}$



Some sample images generated by Vector Quantized Variational Auto-Encoder (VQ-VAE), a state-of-the-art latent variable model for generation



Factor Analysis and Probabilistic PCA



CS772A: PML

- FA and PPCA assume f to be a linear model
- \blacksquare In FA/PPCA, latent variables $\pmb{z}_n \in \mathbb{R}^K$ typically assumed to have a Gaussian prior
 - If we want sparse latent variabled, can use Laplace or spike-and-slab prior on z_n
 - More complex extensions of FA/PPCA use a mixture of Gaussians prior on z_n
- Assumption: Observations $x_n \in \mathbb{R}^D$ typically assumed to have a Gaussian likelihood
 - Other likelihood models (e.g., exp-family) can also be used if data not real-valued
- Relationship between z_n and x_n modeled by a noisy linear mapping

$$\begin{aligned} \mathbf{x}_n &= \mathbf{W} \mathbf{z}_n + \epsilon_n \\ \text{Zero-mean and diagonal or spherical Gaussian noise} \\ \end{aligned} \\ \mathbf{x}_n &= \mathbf{W} \mathbf{z}_n + \epsilon_n \\ \text{Linear combination} \\ \text{of the columns of } \mathbf{W} \end{aligned} \\ \begin{aligned} \mathbf{y}_{k=1} &= \mathbf{y}_{k=1} \\ \mathbf{y}_{k=1} &= \mathbf{y}_{k=1} \\ \text{Linear combination} \\ \text{Linear combination} \\ \text{of the columns of } \mathbf{W} \end{aligned} \\ \begin{aligned} \mathbf{y}_{k=1} &= \mathbf{y}_{k=1} \\ \text{Linear combination} \\ \mathbf{y}_{k=1} &= \mathbf{y}_{k=1} \\ \mathbf{y}_{k=1} \\ \mathbf{y}_{k=1} &= \mathbf{y}_{k=1} \\ \mathbf{y}_{k=1} &= \mathbf{y}_{k=1} \\ \mathbf{y}_{k=1} &= \mathbf{y}_{k=1} \\ \mathbf{y}_{k=1} \\$$

- Linear Gaussian Model. W, z_n 's, and Ψ can be learned (e.g., using EM, VI, MCMC)

Some Variants of FA/PPCA

Gamma-Poisson latent factor model.

Popular for modeling countvalued data (in text analysis, recommender systems, etc) Non-negative priors often give a nice interpretability to such latent variable models (will see some more examples of such models shortly)

- Assumes K-dim non-negative latent variable \mathbf{z}_n and D-dim count-valued observations \mathbf{x}_n
- \hfill An example: Each x_n is the word-count vector representing a document

 $p(\mathbf{z}_{n}) = \prod_{k=1}^{K} \text{Gamma}(\mathbf{z}_{nk}|\mathbf{a}_{k}, \mathbf{b}_{k}))$ $p(\mathbf{x}_{n}|\mathbf{z}_{n}) = \prod_{d=1}^{D} \text{Poisson}(\mathbf{x}_{nd}|f(\mathbf{w}_{d}, \mathbf{z}_{n}))^{\checkmark}$

This is the rate of the Poisson. It should be non-negative, $\exp(\mathbf{w}_d^{\mathsf{T}} \mathbf{z}_n)$, or simply $\mathbf{w}_d^{\mathsf{T}} \mathbf{z}_n$ if \mathbf{w}_d is also non-negative (e.g., using a gamma/Dirichlet prior on it)

- This can be thought of as a probabilistic non-negative matrix factorization model
- Dirichlet-Multinomial/Multinoulli PCA
 - Assumes K-dim non-negative latent variable \mathbf{z}_n and D categorical obs $\mathbf{x}_n = \{\mathbf{x}_{nd}\}_{d=1}^D$
 - An example: Each \mathbf{x}_n is a document with D words in it (each word is a categorical value)

Also sums to 1

 $p(\mathbf{z}_n) = \text{Dirichlet}(\mathbf{z}_n | \boldsymbol{\alpha})$

 $p(\mathbf{x}_{n}|\mathbf{z}_{n}) = \prod_{d=1}^{D} \text{Multinoulli}(\mathbf{x}_{nd}|f(\mathbf{w}_{d},\mathbf{z}_{n}))$

This should give the probability vector of the multinoulli over x_{nd} . It should be non-negative and should sums to 1

A Deep Generative Model: Variational Auto-encoder (VAE)

VAE* is a probabilistic extension of autoencoders (AE). An AE is shown below



- The basic difference is that VAE assumes a prior p(z) on the latent code z
 - This enables it to not just compress the data but also generate synthetic data
 - How: Sample \boldsymbol{z} from a prior and pass it through the decoder
- Thus VAE can learn good latent representation + generate novel synthetic data
- The name has "Variational" in it since it is learned using VI principles

12

Variational Autoencoder (VAE)



The Reparametrization Trick is commonly used to optimize the ELBO

Posterior is inferred only over z, and usually only point estimate on θ

Amortized Inference

- Latent variable models need to infer the posterior $p(\mathbf{z}_n | \mathbf{x}_n)$ for each observation \mathbf{x}_n
- This can be slow if we have lots of observations because
 - 1. We need to iterate over each $p(\boldsymbol{z}_n | \boldsymbol{x}_n)$
 - 2. Learning the global parameters needs wait for step 1 to finish for all observations
- One way to address this is via Stochastic VI
- Amortized inference is another appealing alternative (used in VAE and other LVMs too)

 $p(\mathbf{z}_n | \mathbf{x}_n) \approx q(\mathbf{z}_n | \phi_n) = q(\mathbf{z}_n | NN(\mathbf{x}_n; \mathbf{W}))$ If q is Gaussian then the NN will output a mean and a variance

- Thus no need to learn ϕ_n 's (one per data point) but just a single NN with params W
 - This will be our "encoder network" for learning \boldsymbol{z}_n
 - Also very efficient to get $p(\pmb{z}_*|\pmb{x}_*)$ for a new data point \pmb{x}_*

Variational Autoencoder: The Complete Pipeline

15

CS772A: PML

Both probabilistic encoder and decoder learned jointly by maximizing the ELBO





Pic source: https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae.html

Generative Adversarial Network (GAN)

- GAN is an implicit generative latent variable model
- Can generate from it but can't compute p(x) the model doesn't define it explicitly
- GAN is trained using an adversarial way (Goodfellow et al, 2013)



Thus can't train

Unlike VAE, no explicit parametric

likelihood model p(x|z)

Generative Adversarial Network (GAN)

The GAN training criterion was

 $\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{data}(\boldsymbol{x})}[\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z})))]$

- With G fixed, the optimal D (exercise) $D_{G}^{*}(x) = \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)}$ Distribution of synthetic data
- Given the optimal D, The optimal generator G is found by minimizing

$$V(D_{g}^{*},G) = \mathbb{E}_{x \sim p_{data}} \left[\log \frac{p_{data}(x)}{p_{data}(x) + p_{g}(x)} \right] + \mathbb{E}_{x \sim p_{g}} \left[\log \frac{p_{g}(x)}{p_{data}(x) + p_{g}(x)} \right]$$

$$Jensen-Shannon divergence between p_{data} and p_{g}.$$

$$Minimized when p_{g} = p_{data}$$

$$P_{g} = p_{data}$$

$$Full SGAN can learn the true data distribution if the generator and discriminator have enough modeling power
$$CS772A: PM$$$$

GAN Optimization $\min_{G} \max_{D} V(D, G) = \mathbb{E}_{x \sim p_{data}(x)}[\log D(x)] + \mathbb{E}_{z \sim p_{z}(z)}[\log(1 - D(G(z))]]$ **•** The GAN training procedure can be summarized as **•** Initialize $\theta_{g}, \theta_{d};$ **•** for each training iteration do **•** In practice, for stable training, we run K > 1 steps of **•** optimizing w.rt. D and 1 step of optimizing w.rt. G

- 4 Sample minibatch of M noise vectors $\mathbf{z}_m \sim q_z(\mathbf{z})$;
- 5 Sample minibatch of M examples $\mathbf{x}_m \sim p_D$;
- 6 Update the discriminator by performing stochastic gradient *ascent* using this gradient: $\nabla_{\boldsymbol{\theta}_d} \frac{1}{M} \sum_{m=1}^{M} \left[\log D(\mathbf{x}_m) + \log(1 - D(G(\mathbf{z}_m))) \right].$;
- 7 Sample minibatch of M noise vectors $\mathbf{z}_m \sim q_z(\mathbf{z})$;
- 8 Update the generator by performing stochastic gradient *descent* using this gradient: $\nabla \theta_g \frac{1}{M} \sum_{m=1}^{M} \log(1 - D(G(\mathbf{z}_m))). ;$

9 Return $\boldsymbol{\theta}_g, \, \boldsymbol{\theta}_d$

In practice, in this step, instead of minimizing

 $\log(1 - D(G(z)))$, we maximize $\log(D(G(z)))$

GANs that also learn latent representations

- The standard GAN can only generate data. Can't learn the latent z from x
- Bidirectional GAN* (BiGAN) is a GAN variant that allows this



19

Evaluating GANs

- Two measures that are commonly used to evaluate GANs
 - Inception score (IS): Evaluates the distribution of generated data
 - Frechet inception distance (FID): Compared the distribution of real data and generated data
- Inception Score defined as $\exp(\mathbb{E}_{x \sim p_q}[\mathrm{KL}(p(y|x)||p(y))])$ will be high if
 - Very few high-probability classes in each sample x: Low entropy for p(y|x)
 - We have diverse classes across samples: Marginal p(y) is close to uniform (high entropy)
- FID uses extracted features (using a deep neural net) of real and generated data
 - Usually from the layers closer to the output layer
- These features are used to estimate two Gaussian distributions

Using real data $\mathcal{N}(\mu_R, \Sigma_R)$ $\mathcal{N}(\mu_G, \Sigma_G)$ Using generated data

- FID is then defined as FID = $|\mu_G \mu_R|^2 + \text{trace}(\Sigma_G + \Sigma_R (\Sigma_G \Sigma_R)^{1/2})$
- These measures can also be used for evaluating other deep gen models like VAEML

High IS and low FID is desirable

Both IS and FID measure how

realistic the generated data is

GAN: Some Issues/Comments

- GAN training can be hard and the basic GAN suffers from several issues
- Instability of training procedure
- Mode Collapse problem: Lack of diversity in generated samples
 - Generator may find some data that can easily fool the discriminator
 - It will stuck at that mode of the data distribution and keep generating data like that



GAN 1: No mode collapse (all 10 modes captured in generation)

GAN 2: Mode collapse (stuck on one of the modes)

Some work on addressing these issues (e.g., Wasserstein GAN, Least Squares GAN, etc)