

MiniBatch Monte Carlo simulation

Chernov Andrey

May 9, 2021

Bayesian Inference for black box models

- Compute posterior using Bayes Rule:

$$p(\mathbf{w} \mid \mathcal{D}) = \frac{p(\mathbf{w})p(\mathcal{D} \mid \mathbf{w})}{p(\mathcal{D})}$$

Using MCMC to calculate $p(w|D)$.

- Make predictions using the posterior predictive distribution:

$$\begin{aligned} p(t_{test} \mid x_{test}, \mathcal{D}_{train}) &= \int p(\mathbf{w} \mid \mathcal{D}_{train})p(t_{test} \mid x_{test}, \mathbf{w})d\mathbf{w} \\ &\approx \sum_i w_i p(t|x, w_i) \end{aligned}$$

Where $w_i \sim p(\mathbf{w} \mid \mathcal{D}_{train})$

MCMC

- Markov chain Monte Carlo (MCMC) is a family of methods comprise a class of algorithms for sampling from the desired probability distribution.
- How it works in practice
 - Construct a Markov chain that has the desired distribution as its stationary distribution
 - Sample from a Markov chain while It will not converge to stationary distribution (Warm up the chain)
 - Now, you got the "black box", which samples from the desired distribution

Metropolis-Hastings Algorithm

- **The task** is to generate samples from distribution $p(T)$, known up to a normalization constant.

$$p(T) = \frac{1}{Z} \tilde{p}(T)$$

- **The algorithm**

- The sample θ_{t+1} at time $t + 1$ is generated using a candidate θ' from a (simpler) proposal distribution $q(\theta' | \theta_t)$, filtered by an acceptance test.
- The acceptance test has acceptance probability:

$$\alpha(\theta_t, \theta') = \frac{p(\theta')q(\theta_t|\theta')}{p(\theta_t)q(\theta'|\theta_t)} \wedge 1$$
- if $u < \alpha(\theta_t, \theta')$ $\theta_{t+1} = \theta'$, otherwise set $\theta_{t+1} = \theta_t$. Where $u \sim \mathcal{U}(0, 1)$

We must pass each point from dataset for one acceptance test

For Bayesian inference, the target distribution is $p(\theta \mid x_1, \dots, x_N)$. The acceptance probability is now:

$$\alpha(\theta, \theta') = \frac{p_0(\theta') \prod_{i=1}^N p(x_i \mid \theta') q(\theta \mid \theta')}{p_0(\theta) \prod_{i=1}^N p(x_i \mid \theta) q(\theta' \mid \theta)} \wedge 1$$

It is not scalable for big datasets.

Let us introduce the following notation

$$\begin{aligned}\Lambda_i(\theta, \hat{\theta}) &= \log \frac{p(x_i | \hat{\theta})}{p(x_i | \theta)} & \Lambda(\theta, \hat{\theta}) &= \sum_{i=1}^n \Lambda_i \\ \psi(\theta, \hat{\theta}) &= \log \frac{p(\theta) q(\hat{\theta} | \theta)}{p(\hat{\theta}) q(\theta | \hat{\theta})} & \Delta(\theta, \hat{\theta}) &= \Lambda(\theta, \hat{\theta}) - \psi(\theta, \hat{\theta})\end{aligned}$$

The unbiased estimation for Λ и Δ

$$\begin{aligned}\Lambda^*(\theta, \hat{\theta}) &= \frac{n}{b} \sum_{i=1}^b \log \frac{p(x_i^* | \hat{\theta})}{p(x_i^* | \theta)} \\ \Delta^*(\theta, \hat{\theta}) &= \Lambda^*(\theta, \hat{\theta}) - \psi(\theta, \hat{\theta}) \\ x_i^* &\sim \mathcal{U}[\{x_1, \dots, x_n\}]\end{aligned}$$

Barker Lemma

Lemma 2. If $g(s)$ is any function such that $g(s) = \exp(s)g(-s)$, then the sampling with acceptance function $\alpha(\theta, \theta') \triangleq g(\Delta(\theta, \theta'))$ get desired distribution.

The sigmoid function satisfies this lemma. Indeed,

$$\exp(s)g(-s) = \frac{\exp(s)}{1 + \exp(-s)} = \frac{1}{\exp(-s) + 1} = g(s)$$

MiniBatch Metropolis-Hastings Algorithm

$$\alpha(\theta, \hat{\theta}) \wedge V, \quad V \sim \mathcal{U}[0, 1]$$

Apply Barker's lemma

$$g(\Delta(\theta, \hat{\theta})) \wedge V, \quad V \sim \mathcal{U}[0, 1]$$

Require the monotonicity of the function g .

$$\Delta(\theta, \hat{\theta}) \wedge g^{-1}(V) = X_{\log}$$

Using CLT for Δ

$$\Delta^* \sim \mathcal{N}(\Delta, \sigma^2(\Delta^*)) = \Delta + \mathcal{N}(0, \sigma^2(\Delta^*))$$

If $\sigma^2(\Delta^*) < \text{Var}(X_{\log}) = \frac{\pi^2}{3}$, we can decompose X_{\log} as

$$X_{\log} = \mathcal{N}(0, \sigma^2(\Delta^*)) + X_{\text{corr}}$$

MiniBatch Metropolis-Hastings Algorithm

Thus, we accept the new point, if:

$$\Delta + X_{\log} > 0$$

$$\Delta + X_{\log} = (\Delta + X_{\text{norm}}) + X_{\text{corr}} = \Delta^* + X_{\text{corr}} > 0$$

where $X_{\text{norm}} \sim \mathcal{N}(\Delta, \sigma^2(\Delta^*))$

MiniBatch Metropolis-Hastings Algorithm

Algorithm 2 Minibatch Metropolis Hastings acceptance test

Input: $\hat{\theta}$, sampled from $q(\hat{\theta} | \theta^t)$. and previous point θ^t

Output: New point θ^{t+1} from the desired distribution.

$$\Delta^* = 0, \sigma^2(\Delta^*) = \infty$$

while $\sigma^2(\Delta^*) > \sigma^2$ do

Sample new batch $\{x_i^*\}_1^b$ from $\mathcal{D} : x_i^* \sim \mathcal{U}[\{x_1, \dots, x_n\}]$. Recalculated $\Delta^*, \sigma^2(\Delta^*)$, using $\{x_i^*\}_1^b$ end while

$$X_{\text{corr}} \sim C_{\sigma}^*, X_{\text{norm}} \sim \mathcal{N}(0, \sigma^2 - \sigma^2(\Delta^*))$$

$$\theta^{t+1} = \begin{cases} \hat{\theta}, & \text{if } \Delta^* + X_{\text{norm}} \leq X_{\text{corr}} \\ \theta^t, & \text{if } \Delta^* + X_{\text{norm}} > X_{\text{corr}} \end{cases}$$

Approximate logit distribution using normal distribution

logit PDF:

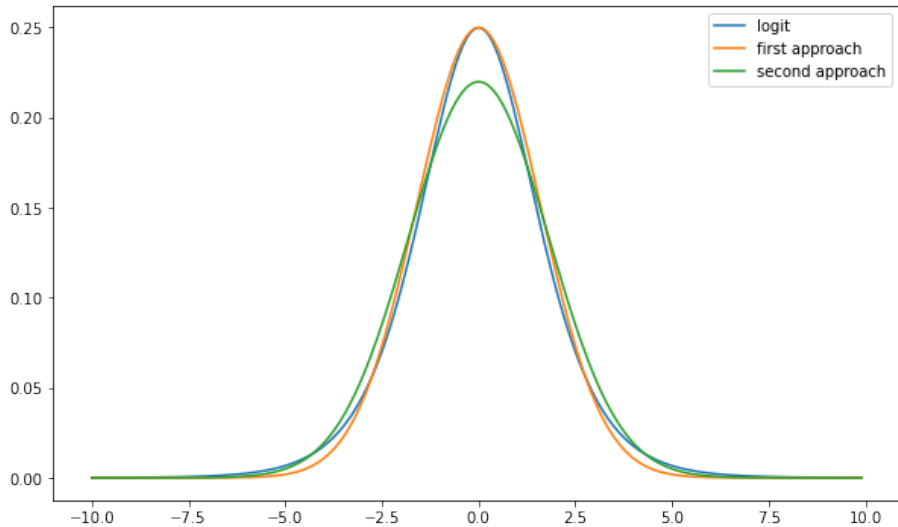
$$f(x) = \frac{\exp(-x)}{(1 + \exp(-x))^2}$$

$$f(0) = \frac{1}{4}$$

$$\phi(0) = \frac{1}{\sqrt{2\pi}\sigma_{prox}}$$

First Idea: Find σ_{prox} from equation $\phi(0) = f(0)$, $\sigma_{prox} = \frac{4}{\sqrt{2\pi}}$

Second Idea: $\operatorname{argmin}_{\sigma_{prox}} KL(f(x) || \phi(x))$; $\sigma_{prox} = \frac{\pi}{\sqrt{3}}$



OK, It does not work anyway

```
// (Part 3.2) Abnormally good or bad minibatches.
else if (math.abs(numStd) > 5.0) {
    if (opts.verboseMH) {
        println("\tCASE 1: math.abs(numStd) = " + math.abs(numStd))
    }
    newMinibatch = true
    if (numStd > 0) {
        accept = true
    }
}

// (Part 3.3) If sample variance is too large, don't do anything.
else if (sampleVariance >= targetVariance) {
    if (opts.verboseMH) {
        println("\tCASE 2: sample >= target = "+targetVariance)
    }
}
```

For MNIST problem with one fully connected layer $\sigma(\Delta^*) > 500$ on **full** dataset

The Idea to fix problem with variance

Let's $g(s) = \frac{K}{1 + \exp(-s)}$. It still satisfy Barker Lemma.

But It is incorrect $g(s)$, because $g(+\infty) = K$

$$\Delta + K * X_{\log} > 0$$

$$\frac{\Delta}{K} + X_{\log} = \left(\frac{\Delta}{K} + X_{\text{norm}}^* \right) + X_{\text{corr}} = \Delta^* + X_{\text{corr}} > 0$$

where $X_{\text{norm}}^* \sim \mathcal{N}\left(\Delta, \frac{\sigma^2(\Delta^*)}{K^2}\right)$.

If $\sigma^2 \approx K^2$ and use Langevin Sampler it works fine.