### Sampling Methods: From MCMC to Generative Modeling Bayesian learning and Langevin algorithm

#### Anna Korba

CREST, ENSAE, Institut Polytechnique de Paris

Bayesian learning

Bayesian deep learning 00000000 References



#### Bayesian learning

Langevin

Bayesian deep learning

#### Motivation for Sampling (1): Bayesian inference

Goal of Bayesian inference: learn the best distribution over a parameter x to fit observed data.

#### Motivation for Sampling (1): Bayesian inference

# Goal of Bayesian inference: learn the best distribution over a parameter x to fit observed data.

(1) Let  $\mathcal{D} = (w_i, y_i)_{i=1}^p$  a dataset of i.i.d. examples with features w, label y.

(2) Assume an underlying model parametrized by  $x \in \mathbb{R}^d$ , e.g.:

$$y = g(w, x) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathrm{Id}).$$

#### Motivation for Sampling (1): Bayesian inference

# Goal of Bayesian inference: learn the best distribution over a parameter x to fit observed data.

(1) Let  $\mathcal{D} = (w_i, y_i)_{i=1}^p$  a dataset of i.i.d. examples with features w, label y.

(2) Assume an underlying model parametrized by  $x \in \mathbb{R}^d$ , e.g.:

$$y = g(w, x) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathrm{Id}).$$

Step 1. Compute the Likelihood:

$$p(\mathcal{D}|x) \stackrel{(1)}{\propto} \prod_{i=1}^{p} p(y_i|x, w_i) \stackrel{(2)}{\propto} \exp(-\frac{1}{2}\sum_{i=1}^{p} \|y_i - g(w_i, x)\|^2).$$

Bayesian deep learning 00000000 References

#### Step 2. Choose a prior distribution (initial guess) on the parameter:

$$x \sim p_0$$
, e.g.  $p_0(x) \propto \exp(-\frac{\|x\|^2}{2})$ .

References

Step 2. Choose a prior distribution (initial guess) on the parameter:

$$x \sim p_0$$
, e.g.  $p_0(x) \propto \exp(-\frac{\|x\|^2}{2})$ .

Step 3. Bayes' rule yields the formula for the posterior distribution over the parameter *x*:

$$p(x|\mathcal{D}) = rac{p(\mathcal{D}|x)p_0(x)}{Z}$$
 where  $Z = \int_{\mathbb{R}^d} p(\mathcal{D}|x)p_0(x)dx$ 

is called the normalization constant and is intractable.

References

Step 2. Choose a prior distribution (initial guess) on the parameter:

$$x \sim p_0$$
, e.g.  $p_0(x) \propto \exp(-\frac{\|x\|^2}{2})$ .

Step 3. Bayes' rule yields the formula for the posterior distribution over the parameter *x*:

$$p(x|\mathcal{D}) = rac{p(\mathcal{D}|x)p_0(x)}{Z}$$
 where  $Z = \int_{\mathbb{R}^d} p(\mathcal{D}|x)p_0(x)dx$ 

is called the normalization constant and is intractable.

Denoting  $\pi := p(\cdot | \mathcal{D})$  the posterior on parameters  $x \in \mathbb{R}^d$ , we have:

$$\pi(x) \propto \exp(-V(x)), \quad V(x) = rac{1}{2} \sum_{i=1}^{p} \|y_i - g(w_i, x)\|^2 + rac{\|x\|^2}{2}.$$

i.e.  $\pi$ 's density is known "up to a normalization constant".  $\pi$  is a probability distribution over parameters of a model.

The posterior  $\pi$  is interesting for

- measuring uncertainty on prediction through the distribution of  $g(w, \cdot)$ ,  $x \sim \pi$ .
- prediction for a new input w:

$$\hat{y} = \underbrace{\int_{\mathbb{R}^d} g(w, x) d\pi(x)}_{\text{"Bayesian model averaging"}}$$

i.e. predictions of models parametrized by  $x \in \mathbb{R}^d$  are reweighted by  $\pi(x)$ .

References

Here, Sampling methods construct an approximation  $\mu_M = \frac{1}{M} \sum_{m=1}^M \delta_{x_m}$  of  $\pi$ .



References

#### Sampling as Optimization

Actually, in many cases (e.g. it is underlying many algorithms), the sampling problem (approximating  $\pi$ ) can be viewed as optimization over  $\mathcal{P}(\mathbb{R}^d)$ :

 $\min_{\mu\in\mathcal{P}(\mathbb{R}^d)}\mathbb{D}(\mu|\pi)$ 

where D is a divergence or distance, hence that is minimized for  $\mu = \pi$ .

References

#### The Kullback-Leibler divergence

 $\rm D$  could be the (reverse) Kullback-Leibler (KL) divergence:

$$\mathrm{KL}(\mu|\pi) = \begin{cases} \int_{\mathbb{R}^d} \log\left(\frac{\mu}{\pi}(x)\right) d\mu(x) & \text{if } \mu \ll \pi \\ +\infty & \text{otherwise.} \end{cases}$$

We recognize a *f*-divergence  $\int f\left(\frac{\mu}{\pi}\right) d\pi$  where  $f(x) = x \log(x)$ . Taking  $f(x) = -\log(x)$  yields the (forward) KL i.e.  $\text{KL}(\pi|\mu)$ .

References

#### The Kullback-Leibler divergence

 $\rm D$  could be the (reverse) Kullback-Leibler (KL) divergence:

$$\mathrm{KL}(\mu|\pi) = \begin{cases} \int_{\mathbb{R}^d} \log\left(\frac{\mu}{\pi}(x)\right) d\mu(x) & \text{if } \mu \ll \pi \\ +\infty & \text{otherwise.} \end{cases}$$

We recognize a *f*-divergence  $\int f\left(\frac{\mu}{\pi}\right) d\pi$  where  $f(x) = x \log(x)$ . Taking  $f(x) = -\log(x)$  yields the (forward) KL i.e.  $\text{KL}(\pi|\mu)$ .

The (reverse) KL as an objective is convenient when the unnormalized density of  $\pi$  is known since it **does not depend on the normalization constant!** 

Indeed writing  $\pi(x) = e^{-V(x)}/Z$  we have:

$$\mathrm{KL}(\mu|\pi) = \int_{\mathbb{R}^d} \log\left(\frac{\mu}{e^{-V}}(x)\right) d\mu(x) + \log(Z).$$

But, it is not convenient when  $\mu$  or  $\pi$  are discrete, because the KL is  $+\infty$  unless  $supp(\mu) \subset supp(\pi)$ .

#### Examples

(Parametric methods) Variational Inference : Restrict the search space to a parametric families {μ<sub>θ</sub>, θ ∈ ℝ<sup>ρ</sup>}. The problem rewrites as a finite-dimensional optimization problem (i.e. over ℝ<sup>ρ</sup>):

## $\min_{\theta \in \mathbb{R}^p} \mathrm{D}(\mu_{\theta} | \pi)$

- Example: Gaussians with diagonal covariance matrices can be parametrized by θ = (m, σ) ∈ ℝ<sup>2d</sup> (see Bayes by Backprop in the last section)
- Example: use normalizing flows to construct a family  $\mu_{\theta} = f_{\theta \#} p$  and optimize the previous objective<sup>1</sup>. <sup>1</sup>Rezende, D., Mohamed, S. (2015, June). Variational inference with normalizing flows. In International conference on machine learning.

#### Examples

(Parametric methods) Variational Inference : Restrict the search space to a parametric families {μ<sub>θ</sub>, θ ∈ ℝ<sup>ρ</sup>}. The problem rewrites as a finite-dimensional optimization problem (i.e. over ℝ<sup>ρ</sup>):

## $\min_{\theta \in \mathbb{R}^p} \mathrm{D}(\mu_{\theta} | \pi)$

- Example: Gaussians with diagonal covariance matrices can be parametrized by θ = (m, σ) ∈ ℝ<sup>2d</sup> (see Bayes by Backprop in the last section)
- Example: use normalizing flows to construct a family  $\mu_{\theta} = f_{\theta \#} p$  and optimize the previous objective<sup>1</sup>. <sup>1</sup>Rezende, D., Mohamed, S. (2015, June). Variational inference with normalizing flows. In International conference on machine learning.
- (Non parametric methods) Markov Chain Monte Carlo (MCMC) methods, Sequential Monte Carlo (SMC)...: generate a Markov chain in  $\mathbb{R}^d$  whose law converges to  $\pi \propto \exp(-V)$
- Example: Langevin (next section)

References

#### Langevin Monte Carlo

Langevin Monte Carlo (LMC) [Roberts and Tweedie (1996)]

$$x_{m+1} = x_m + \gamma \nabla \log \pi(x_m) + \sqrt{2\gamma} \eta_m, \quad \eta_m \sim \mathcal{N}(0, \mathrm{Id}).$$



Picture from https://chi-feng.github.io/mcmc-demo/app.html.

Note that in the Bayesian inference setting, where  $\pi = \frac{\exp(-V)}{Z}$ , it is easily implementable since the score  $\nabla_x \log \pi(x) = -\nabla_x (V(x) + \log(Z)) = -\nabla V(x)$ since  $\nabla_x \log(Z) = 0$ . Bayesian learning

Bayesian deep learning 00000000 References

#### Outline

#### Bayesian learning

Langevin

Bayesian deep learning

References

#### Langevin diffusion

Langevin diffusion is the Stochastic Differential Equation (SDE):

$$\mathrm{d}x_t = -\nabla V(x_t)dt + \sqrt{2}\mathrm{d}B_t, \quad x_t \sim p_t$$

where  $B_t$  denotes the standard Brownian motion in  $\mathbb{R}^d$ , defined as:

- $B_0 = 0$  almost surely;
- For any  $t_0 < t_1 < \cdots < t_N$ , the increments  $B_{t_n} B_{t_{n-1}}$  are independent,  $n = 1, 2, \dots, N$ ;
- The difference  $B_t B_s$  and  $B_{t-s}$  have the same distribution:  $\mathcal{N}(0, (t-s) \operatorname{Id})$  for s < t;
- *B<sub>t</sub>* is continuous almost surely.

Langevin diffusion defines a Markov process as follows:

$$x_t = x_0 - \int_0^t \nabla V(x_s) ds + \sqrt{2}B_t,$$

where  $x_0$  is some initialization.

Bayesian deep learning 00000000 References

#### Time-discretization

An Euler-Maruyama time-discretization of Langevin diffusion yields:

$$x_{t+1} = x_t - \gamma \nabla V(x_t) + \sqrt{2\gamma} \eta_t, \quad \eta_t \sim \mathcal{N}(0, \mathrm{Id}). \tag{1}$$

Bayesian deep learning 00000000 References

#### Time-discretization

An Euler-Maruyama time-discretization of Langevin diffusion yields:

$$x_{t+1} = x_t - \gamma \nabla V(x_t) + \sqrt{2\gamma} \eta_t, \quad \eta_t \sim \mathcal{N}(0, \mathrm{Id}).$$
 (1)

Proof:

$$egin{aligned} & x_\gamma pprox x_0 - \int_0^\gamma 
abla V(x_0) \, dt + \sqrt{2\gamma} \, \eta \ & = x_0 - \left(\int_0^\gamma dt 
ight) 
abla V(x_0) + \sqrt{2\gamma} \, \eta \ & = x_0 - \gamma 
abla V(x_0) + \sqrt{2\gamma} \, \eta. \end{aligned}$$

Bayesian deep learning 00000000 References

#### Time-discretization

An Euler-Maruyama time-discretization of Langevin diffusion yields:

$$x_{t+1} = x_t - \gamma \nabla V(x_t) + \sqrt{2\gamma} \eta_t, \quad \eta_t \sim \mathcal{N}(0, \mathrm{Id}).$$
 (1)

Proof:

$$egin{aligned} & x_\gamma pprox x_0 - \int_0^\gamma 
abla \mathcal{V}(x_0) \, dt + \sqrt{2\gamma} \, \eta \ & = x_0 - \left(\int_0^\gamma dt 
ight) 
abla \mathcal{V}(x_0) + \sqrt{2\gamma} \, \eta \ & = x_0 - \gamma 
abla \mathcal{V}(x_0) + \sqrt{2\gamma} \, \eta. \end{aligned}$$

We can now iterate this approach k times, which gives us a recursion, which can be easily implementable on a computer:

$$x_{k\gamma} pprox x_{(k-1)\gamma} - \gamma 
abla V(x_{(k-1)\gamma}) + \sqrt{2\gamma} \eta_k,$$

where  $\eta_k \sim \mathcal{N}(0, \mathrm{Id})$  for all k. Dropping the dependency on  $\gamma$  in the indices yields the scheme (1).

Bayesian deep learning 00000000 References

#### **Ornstein-Uhlenbeck**

Example:  $\pi \propto \exp(-\frac{\|x\|^2}{2})$ ,

References

#### **Ornstein-Uhlenbeck**

Example: 
$$\pi \propto \exp(-\frac{\|x\|^2}{2})$$
,  $\log \pi(x) = -V(x) = -\frac{\|x\|^2}{2}$ ,  $\nabla \log \pi(x) = -x$ .

References

#### **Ornstein-Uhlenbeck**

Example:  $\pi \propto \exp(-\frac{\|x\|^2}{2})$ ,  $\log \pi(x) = -V(x) = -\frac{\|x\|^2}{2}$ ,  $\nabla \log \pi(x) = -x$ .

(continuous time) Langevin diffusion = Ornstein-Uhlenbeck process:

 $\mathrm{d}x_t = -x_t + \mathrm{d}B_t.$ 

#### References

#### **Ornstein-Uhlenbeck**

Example: 
$$\pi \propto \exp(-\frac{\|x\|^2}{2})$$
,  $\log \pi(x) = -V(x) = -\frac{\|x\|^2}{2}$ ,  $\nabla \log \pi(x) = -x$ .

(continuous time) Langevin diffusion = Ornstein-Uhlenbeck process:

$$\mathrm{d}x_t = -x_t + \mathrm{d}B_t.$$

(discrete time)  $x_{t+1} = x_t - \gamma x_t + \sqrt{2\gamma} \eta_t$ ,  $\eta_t \sim \mathcal{N}(0, \mathrm{Id})$ .

#### **Ornstein-Uhlenbeck**

Example:  $\pi \propto \exp(-\frac{\|x\|^2}{2})$ ,  $\log \pi(x) = -V(x) = -\frac{\|x\|^2}{2}$ ,  $\nabla \log \pi(x) = -x$ .

(continuous time) Langevin diffusion = Ornstein-Uhlenbeck process:

$$\mathrm{d}x_t = -x_t + \mathrm{d}B_t.$$

(discrete time)  $x_{t+1} = x_t - \gamma x_t + \sqrt{2\gamma}\eta_t, \quad \eta_t \sim \mathcal{N}(0, \mathrm{Id}).$ 



Recall above we plot  $x_{t+1} = x_t + \gamma \nabla \log \pi(x_t) + \sqrt{2\gamma} \eta_t$  for  $\pi \propto \exp(-\frac{||x||^2}{2})$ .

Bayesian deep learning 00000000 References

#### The Fokker-Planck equation

**Question:** how does the law  $p_t$  of  $x_t$  evolve? does it converge to  $\pi$ ?

Bayesian deep learning 00000000 References

#### The Fokker-Planck equation

**Question:** how does the law  $p_t$  of  $x_t$  evolve? does it converge to  $\pi$ ?

For simplicity, let us assume d = 1, so that Langevin diffusion becomes:

 $\mathrm{d}x_t = -\partial_x V(x_t) \,\mathrm{d}t + \sqrt{2} \,\mathrm{d}B_t,$ 

Bayesian deep learning 00000000

References

#### The Fokker-Planck equation

**Question:** how does the law  $p_t$  of  $x_t$  evolve? does it converge to  $\pi$ ?

For simplicity, let us assume d = 1, so that Langevin diffusion becomes:

$$\mathrm{d} x_t = -\partial_x V(x_t) \,\mathrm{d} t + \sqrt{2} \,\mathrm{d} B_t,$$

To understand how p(x, t) evolves, we will use the Fokker–Planck equation, which governs the evolution of p(x, t) through the following partial differential equation (PDE):

$$\partial_t p(x,t) = \partial_x [\partial_x V(x)p(x,t)] + \partial_x^2 p(x,t).$$

This equation characterizes how the "change" in  $p(\cdot, t)$  behaves, i.e.,  $\partial_t p(x, t)$ .

Bayesian deep learning 00000000

References

#### The Fokker-Planck equation

**Question:** how does the law  $p_t$  of  $x_t$  evolve? does it converge to  $\pi$ ?

For simplicity, let us assume d = 1, so that Langevin diffusion becomes:

$$\mathrm{d} x_t = -\partial_x V(x_t) \,\mathrm{d} t + \sqrt{2} \,\mathrm{d} B_t,$$

To understand how p(x, t) evolves, we will use the Fokker–Planck equation, which governs the evolution of p(x, t) through the following partial differential equation (PDE):

$$\partial_t p(x,t) = \partial_x [\partial_x V(x)p(x,t)] + \partial_x^2 p(x,t).$$

This equation characterizes how the "change" in  $p(\cdot, t)$  behaves, i.e.,  $\partial_t p(x, t)$ .

**Remark:** for d > 1, the Fokker-Planck equation writes:

$$\partial_t p(x,t) = \nabla \cdot (\nabla V(x)p(x,t)) + \Delta(p(x,t)).$$

(where  $\nabla \cdot$  and  $\Delta$  are the divergence and Laplacian operators: analog to above but summing all partial derivatives for  $x_1, \ldots, x_d$ ).

Bayesian deep learning 00000000 References

#### The Fokker-Planck equation

Now, the idea is: if  $p(\cdot, t)$  converges to a distribution as  $t \to \infty$ , then whenever this limit is reached, there should not be any more changes in p. In other words, whenever  $p(\cdot, t)$  hits its limit,  $\partial_t p(x, t)$  has to be equal to 0.

Bayesian deep learning 00000000

#### The Fokker-Planck equation

Now, the idea is: if  $p(\cdot, t)$  converges to a distribution as  $t \to \infty$ , then whenever this limit is reached, there should not be any more changes in p. In other words, whenever  $p(\cdot, t)$  hits its limit,  $\partial_t p(x, t)$  has to be equal to 0.

Therefore, we can simply "check" if  $\pi \propto \exp(-V)$  is a limit of  $p(\cdot, t)$  by replacing p(x, t) with  $\pi(x)$  in the Fokker–Planck equation and observing whether the right-hand side is equal to 0 or not. Let us apply this procedure:

Bayesian deep learning 00000000 References

#### The Fokker-Planck equation

Now, the idea is: if  $p(\cdot, t)$  converges to a distribution as  $t \to \infty$ , then whenever this limit is reached, there should not be any more changes in p. In other words, whenever  $p(\cdot, t)$  hits its limit,  $\partial_t p(x, t)$  has to be equal to 0.

Therefore, we can simply "check" if  $\pi \propto \exp(-V)$  is a limit of  $p(\cdot, t)$  by replacing p(x, t) with  $\pi(x)$  in the Fokker–Planck equation and observing whether the right-hand side is equal to 0 or not. Let us apply this procedure:

$$\partial_x \left[\partial_x V(x)\pi(x)
ight] + \partial_x^2 \pi(x) = \partial_x \left[\partial_x V(x)\pi(x) + \partial_x \pi(x)
ight]$$
  
=  $\partial_x \left[\partial_x V(x)\pi(x) - \partial_x V(x)\pi(x)
ight]$   
= 0,

where we used the fact that

$$\partial_x V(x) = -\partial_x \log \pi(x) = -rac{1}{\pi(x)} \partial_x \pi(x),$$

hence

$$\partial_x \pi(x) = -\pi(x)\partial_x V(x).$$

Bayesian deep learning 00000000

#### The Fokker-Planck equation

Now, the idea is: if  $p(\cdot, t)$  converges to a distribution as  $t \to \infty$ , then whenever this limit is reached, there should not be any more changes in p. In other words, whenever  $p(\cdot, t)$  hits its limit,  $\partial_t p(x, t)$  has to be equal to 0.

Therefore, we can simply "check" if  $\pi \propto \exp(-V)$  is a limit of  $p(\cdot, t)$  by replacing p(x, t) with  $\pi(x)$  in the Fokker–Planck equation and observing whether the right-hand side is equal to 0 or not. Let us apply this procedure:

$$\partial_x \left[\partial_x V(x)\pi(x)
ight] + \partial_x^2 \pi(x) = \partial_x \left[\partial_x V(x)\pi(x) + \partial_x \pi(x)
ight]$$
  
=  $\partial_x \left[\partial_x V(x)\pi(x) - \partial_x V(x)\pi(x)
ight]$   
= 0,

where we used the fact that

$$\partial_x V(x) = -\partial_x \log \pi(x) = -\frac{1}{\pi(x)} \partial_x \pi(x),$$

hence

$$\partial_x \pi(x) = -\pi(x)\partial_x V(x).$$

Conclusion:  $\pi$  is an equilibrium for the FP equation !

Bayesian deep learning 00000000 References

#### **Ornstein–Uhlenbeck Process**

We now focus on a specific case of a Langevin diffusion and we will prove that: For the SDE:

$$dX_t = -\beta X_t \, dt + \sigma \, dB_t$$

The solution is:

$$X_t = e^{-\beta t} X_0 + \sigma e^{-\beta t} \int_0^t e^{\beta s} \, dB_s$$

with stationary/limiting distribution  $\pi = \mathcal{N}(0, \frac{\sigma^2}{2\beta})$ and we have

$$X_t \mid X_0 \sim \mathcal{N}\left(e^{-eta t}X_0, \frac{\sigma^2}{2eta}(1-e^{-2eta t})
ight)$$

#### Observe that:

• The farther into the future, the more the initial value gets "forgotten"

Bayesian deep learning

References

#### Proof

#### Step 1 (Multiply by the integrating factor) Multiply both sides of the SDE by $\mu(t) = e^{\beta t}$ :

$$e^{\beta t}dX_t = -\beta e^{\beta t}X_t\,dt + \sigma e^{\beta t}dB_t$$
Bayesian deep learning

References

#### Proof

#### Step 1 (Multiply by the integrating factor) Multiply both sides of the SDE by $\mu(t) = e^{\beta t}$ :

$$e^{\beta t}dX_t = -\beta e^{\beta t}X_t dt + \sigma e^{\beta t}dB_t$$

But using (Itô's) product rule:

$$d\left(e^{\beta t}X_{t}\right)=e^{\beta t}dX_{t}+\beta e^{\beta t}X_{t}\,dt$$

Bayesian deep learning

References

#### Proof

#### Step 1 (Multiply by the integrating factor) Multiply both sides of the SDE by $\mu(t) = e^{\beta t}$ :

$$e^{\beta t}dX_t = -\beta e^{\beta t}X_t dt + \sigma e^{\beta t}dB_t$$

But using (Itô's) product rule:

$$d\left(e^{\beta t}X_{t}\right)=e^{\beta t}dX_{t}+\beta e^{\beta t}X_{t}\,dt$$

So we get:

$$d\left(e^{\beta t}X_{t}\right)=\sigma e^{\beta t}dB_{t}$$

Bayesian deep learning

References

#### Proof

#### Step 1 (Multiply by the integrating factor) Multiply both sides of the SDE by $\mu(t) = e^{\beta t}$ :

$$e^{\beta t}dX_t = -\beta e^{\beta t}X_t dt + \sigma e^{\beta t}dB_t$$

But using (Itô's) product rule:

$$d\left(e^{\beta t}X_{t}\right)=e^{\beta t}dX_{t}+\beta e^{\beta t}X_{t}\,dt$$

So we get:

$$d\left(e^{\beta t}X_{t}\right)=\sigma e^{\beta t}dB_{t}$$

**Step 2 (Integrate both sides)** Now integrate from 0 to *t*:

$$e^{\beta t}X_t - X_0 = \sigma \int_0^t e^{\beta s} \, dB_s$$

Bayesian deep learning

References

#### Proof

#### Step 1 (Multiply by the integrating factor) Multiply both sides of the SDE by $\mu(t) = e^{\beta t}$ :

$$e^{\beta t}dX_t = -\beta e^{\beta t}X_t dt + \sigma e^{\beta t}dB_t$$

But using (Itô's) product rule:

$$d\left(e^{\beta t}X_{t}\right)=e^{\beta t}dX_{t}+\beta e^{\beta t}X_{t}\,dt$$

So we get:

$$d\left(e^{\beta t}X_{t}\right)=\sigma e^{\beta t}dB_{t}$$

**Step 2 (Integrate both sides)** Now integrate from 0 to *t*:

$$e^{\beta t}X_t - X_0 = \sigma \int_0^t e^{\beta s} \, dB_s$$

Rewriting:

$$X_t = e^{-\beta t} X_0 + \sigma e^{-\beta t} \int_0^t e^{\beta s} \, dB_s$$

Bayesian deep learning

References

### Proof (continued)

#### Step 3 (Distribution of the integral term )

Let:  $I_t := \int_0^t e^{\beta s} dB_s$ . This is an Itô integral of a deterministic function  $\Rightarrow$  it's a **Gaussian random variable** with:

Bayesian deep learning 00000000 References

### Proof (continued)

## Step 3 (Distribution of the integral term ) Let: $I_t := \int_0^t e^{\beta s} dB_s$ . This is an Itô integral of a deterministic function $\Rightarrow$ it's a Gaussian random variable with:

• Mean:  $\mathbb{E}[I_t] = 0$ 

Bayesian deep learning 00000000 References

## Proof (continued)

# Step 3 (Distribution of the integral term ) Let: $I_t := \int_0^t e^{\beta s} dB_s$ . This is an Itô integral of a deterministic function $\Rightarrow$ it's a Gaussian random variable with:

- Mean:  $\mathbb{E}[I_t] = 0$
- Variance :

$$\begin{aligned} \mathsf{Var}(I_t) &= \mathbb{E}\left[\left(\int_0^t e^{\beta s} \, dB_s\right)^2\right] = \int_0^t \left(e^{\beta s}\right)^2 ds \quad (\text{using Itô isometry}) \\ &= \int_0^t e^{2\beta s} \, ds = \left[\frac{1}{2\beta}e^{2\beta s}\right]_0^t = \frac{1}{2\beta}(e^{2\beta t} - 1). \end{aligned}$$

Bayesian deep learning 00000000 References

### Proof (continued)

# Step 3 (Distribution of the integral term ) Let: $I_t := \int_0^t e^{\beta s} dB_s$ . This is an Itô integral of a deterministic function $\Rightarrow$ it's a Gaussian random variable with:

- Mean:  $\mathbb{E}[I_t] = 0$
- Variance :

$$\begin{aligned} \mathsf{Var}(I_t) &= \mathbb{E}\left[\left(\int_0^t e^{\beta s} \, dB_s\right)^2\right] = \int_0^t \left(e^{\beta s}\right)^2 ds \quad (\text{using Itô isometry}) \\ &= \int_0^t e^{2\beta s} \, ds = \left[\frac{1}{2\beta}e^{2\beta s}\right]_0^t = \frac{1}{2\beta}(e^{2\beta t} - 1). \end{aligned}$$

Therefore:

$$\sigma e^{-\beta t} I_t \sim \mathcal{N}\left(0, \ \sigma^2 e^{-2\beta t} \cdot \frac{1}{2\beta} (e^{2\beta t} - 1)\right) = \mathcal{N}\left(0, \ \frac{\sigma^2}{2\beta} (1 - e^{-2\beta t})\right).$$

Bayesian deep learning 00000000 References

### Proof (continued)

## Step 3 (Distribution of the integral term ) Let: $I_t := \int_0^t e^{\beta s} dB_s$ . This is an Itô integral of a deterministic function $\Rightarrow$ it's a Gaussian random variable with:

- Mean:  $\mathbb{E}[I_t] = 0$
- Variance :

$$\begin{aligned} \mathsf{Var}(I_t) &= \mathbb{E}\left[\left(\int_0^t e^{\beta s} \, dB_s\right)^2\right] = \int_0^t \left(e^{\beta s}\right)^2 ds \quad (\text{using Itô isometry}) \\ &= \int_0^t e^{2\beta s} \, ds = \left[\frac{1}{2\beta}e^{2\beta s}\right]_0^t = \frac{1}{2\beta}(e^{2\beta t} - 1). \end{aligned}$$

Therefore:

$$\sigma e^{-\beta t} I_t \sim \mathcal{N}\left(0, \ \sigma^2 e^{-2\beta t} \cdot \frac{1}{2\beta} (e^{2\beta t} - 1)\right) = \mathcal{N}\left(0, \ \frac{\sigma^2}{2\beta} (1 - e^{-2\beta t})\right).$$

So the full solution is :  $X_t = e^{-\beta t} X_0 + \sigma e^{-\beta t} I_t$ , where  $X_t \mid X_0 \sim \mathcal{N}\left(e^{-\beta t} X_0, \frac{\sigma^2}{2\beta}(1-e^{-2\beta t})\right)$ . Done!

 Bayesian deep learning 00000000 References

#### (Very) Important remarks



Figure: Representing  $X_t$  an OU process (with  $\beta = \sigma = 1$ ), and  $p_t$  its (time) marginals

• We know that the full solution :

$$X_t = e^{-\beta t} X_0 + \text{Gaussian noise}$$
(2)

where Gaussian noise  $\sim \mathcal{N}\left(0, \ \frac{\sigma^2}{2\beta}(1-e^{-2\beta t})\right)$  and that conditionally on  $X_0$ :

$$X_t \mid X_0 \sim \mathcal{N}\left(e^{-\beta t}X_0, \frac{\sigma^2}{2\beta}(1 - e^{-2\beta t})\right)$$
(3)

 Bayesian deep learning 00000000 References

#### (Very) Important remarks



Figure: Representing  $X_t$  an OU process (with  $\beta = \sigma = 1$ ), and  $p_t$  its (time) marginals

• We know that the full solution :

$$X_t = e^{-\beta t} X_0 + \text{Gaussian noise}$$
(2)

where Gaussian noise  $\sim \mathcal{N}\left(0, \; rac{\sigma^2}{2eta}(1-e^{-2eta t})
ight)$  and that conditionally on  $X_0$ :

$$X_t \mid X_0 \sim \mathcal{N}\left(e^{-\beta t}X_0, \frac{\sigma^2}{2\beta}(1 - e^{-2\beta t})\right)$$
(3)

 The marginals (p<sub>t</sub>)<sub>t≥0</sub>, where p<sub>t</sub> the law of X<sub>t</sub> in (2) are not Gaussian in general !! (see gray density in the figure above)

 Bayesian deep learning 00000000 References

#### (Very) Important remarks



Figure: Representing  $X_t$  an OU process (with  $\beta = \sigma = 1$ ), and  $p_t$  its (time) marginals

We know that the full solution :

$$X_t = e^{-\beta t} X_0 + \text{Gaussian noise}$$
(2)

where Gaussian noise  $\sim \mathcal{N}\left(0, \; rac{\sigma^2}{2eta}(1-e^{-2eta t})
ight)$  and that conditionally on  $X_0$ :

$$X_t \mid X_0 \sim \mathcal{N}\left(e^{-\beta t}X_0, \frac{\sigma^2}{2\beta}(1 - e^{-2\beta t})\right)$$
(3)

- The marginals (p<sub>t</sub>)<sub>t≥0</sub>, where p<sub>t</sub> the law of X<sub>t</sub> in (2) are not Gaussian in general !! (see gray density in the figure above)
- but the conditional laws in (3) are Gaussian

Bayesian deep learning 00000000

#### Introducing some initial Condition

#### When are the marginals $p_t$ Gaussian? Answer: when $p_0$ is Gaussian.

Bayesian deep learning 00000000 References

### Introducing some initial Condition

#### When are the marginals $p_t$ Gaussian? Answer: when $p_0$ is Gaussian.

Assume 
$$X_0 \sim \mathcal{N}\left(0, rac{\sigma^2}{2eta}
ight).$$
  
Then we have  $\Rightarrow X_t \sim \mathcal{N}\left(0, rac{\sigma^2}{2eta}
ight).$ 

Proof: Recall  $X_t = A + B$  where  $A = e^{-\beta t} X_0$ ,  $B = \sigma e^{-\beta t} \int_0^t e^{\beta s} dW_s$ .

• 
$$A \sim \mathcal{N}(0, e^{-2\beta t} \cdot \frac{\sigma^2}{\beta})$$

• 
$$B \sim \mathcal{N}(0, \frac{\sigma^2}{2\beta}(1 - e^{-2\beta t}))$$

Bayesian deep learning 00000000 References

#### Introducing some initial Condition

When are the marginals  $p_t$  Gaussian? Answer: when  $p_0$  is Gaussian.

Assume 
$$X_0 \sim \mathcal{N}\left(0, rac{\sigma^2}{2eta}
ight)$$
.  
Then we have  $\Rightarrow X_t \sim \mathcal{N}\left(0, rac{\sigma^2}{2eta}
ight)$ .

Proof: Recall  $X_t = A + B$  where  $A = e^{-\beta t} X_0$ ,  $B = \sigma e^{-\beta t} \int_0^t e^{\beta s} dW_s$ .

•  $A \sim \mathcal{N}(0, e^{-2\beta t} \cdot rac{\sigma^2}{\beta})$ 

• 
$$B \sim \mathcal{N}(0, \frac{\sigma^2}{2\beta}(1 - e^{-2\beta t}))$$

•  $A \perp B \Rightarrow A + B \sim \mathcal{N}(0, \text{sum of variances})$ 

Above, the law of  $X_t$  does not depend on time, because we have started the process at the stationary distribution  $\pi(x) = \mathcal{N}\left(0, \frac{\sigma^2}{2\beta}\right)$ :

$$\mathsf{lf:} \ X_0 \sim \pi(x) \Rightarrow X_t \sim \pi(x) \quad \mathsf{for all} \ t$$

Bayesian deep learning 00000000 References

#### Introducing some initial Condition

#### When are the marginals $p_t$ Gaussian? Answer: when $p_0$ is Gaussian.

Assume 
$$X_0 \sim \mathcal{N}\left(0, rac{\sigma^2}{2eta}
ight)$$
.  
Then we have  $\Rightarrow X_t \sim \mathcal{N}\left(0, rac{\sigma^2}{2eta}
ight)$ .

Proof: Recall  $X_t = A + B$  where  $A = e^{-\beta t} X_0$ ,  $B = \sigma e^{-\beta t} \int_0^t e^{\beta s} dW_s$ .

• 
$$A \sim \mathcal{N}(0, e^{-2\beta t} \cdot \frac{\sigma^2}{\beta})$$

• 
$$B \sim \mathcal{N}(0, \frac{\sigma^2}{2\beta}(1 - e^{-2\beta t}))$$

•  $A \perp B \Rightarrow A + B \sim \mathcal{N}(0, \text{sum of variances})$ 

Above, the law of  $X_t$  does not depend on time, because we have started the process at the stationary distribution  $\pi(x) = \mathcal{N}\left(0, \frac{\sigma^2}{2\beta}\right)$ :

$$\text{If: } X_0 \sim \pi(x) \Rightarrow X_t \sim \pi(x) \quad \text{for all } t \\$$

In general, for a  $X_0 \sim \mathcal{N}(0, \sigma_0^2)$ , we would have

$$X_t \sim \mathcal{N}\left(0, \ e^{-2eta t}\sigma_0^2 + rac{\sigma^2}{2eta}(1-e^{-2eta t})
ight).$$

Bayesian deep learning 00000000

### Back to general Langevin diffusion

• We have spent quite a lot of time on Ornstein-Uhlenbeck (OU):

$$dx_t = -\beta x_t \, dt + \sigma \, dB_t$$

Solution:

$$x_t = e^{-\beta t} x_0 + \sigma e^{-\beta t} \int_0^t e^{\beta s} \, dB_s$$

Distribution:

$$X_t \mid X_0 \sim \mathcal{N}\left(e^{-eta t}X_0, rac{\sigma^2}{2eta}(1-e^{-2eta t})
ight)$$

......

Bayesian deep learning 00000000

### Back to general Langevin diffusion

• We have spent quite a lot of time on Ornstein-Uhlenbeck (OU):

$$dx_t = -\beta x_t \, dt + \sigma \, dB_t$$

Solution:

$$x_t = e^{-\beta t} x_0 + \sigma e^{-\beta t} \int_0^t e^{\beta s} \, dB_s$$

Distribution:

$$X_t \mid X_0 \sim \mathcal{N}\left(e^{-eta t}X_0, rac{\sigma^2}{2eta}(1-e^{-2eta t})
ight)$$

• Let's go back to a general Langevin diffusion :

$$\mathrm{d}x_t = -\nabla V(x_t) dt + \sqrt{2} \mathrm{d}B_t, \quad x_t \sim p_t$$

Solution:

$$x_t = x_0 - \int_0^t \nabla V(x_s) ds + \sqrt{2}B_t,$$

- Remember that OU is a specific case of Langevin, where the target/stationary distribution is:  $\pi = \mathcal{N}(0, \frac{\sigma^2}{2\beta})$ , where  $\pi(x) \propto \exp(-\frac{\beta \|x\|^2}{\sigma^2})$
- for general Langevin, the stationary distribution is  $\pi \propto \exp(-V)$ .

Langevin diffusion (and its discretized versions) is an example of a non-parametric method: we built a process  $x_t \in \mathbb{R}^d$ , whose distribution  $p_t$  converges to  $\pi$  as  $t \to \infty$ 

The law (p<sub>t</sub>)<sub>t≥0</sub> of Langevin diffusion (x<sub>t</sub>)<sub>t≥0</sub> is known to follow a gradient flow to minimize D(p|π) = KL(p|π): dp<sub>t</sub> = −∇<sub>W2</sub> KL(p<sub>t</sub>|π)dt (see <sup>1</sup>)



Recall above we plot  $x_{t+1} = x_t + \gamma \nabla \log \pi(x_t) + \sqrt{2\gamma} \eta_t$  for  $\pi \propto \exp(-\frac{\|x\|^2}{2})$ ,  $x_0 \sim p_0$ .

 $^1$  Jordan, R., Kinderlehrer, D., & Otto, F. (1998). The variational formulation of the Fokker–Planck equation. SIAM journal on mathematical analysis.

### When does Langevin diffusion's law converges (fast) to $\pi$ ?

• Consider a standard Gaussian distribution  $\pi(x) \propto \exp(-\frac{\|x\|^2}{2})$ , i.e.  $\pi \propto \exp(-V)$  with V 1-strongly convex, i.e.  $\pi$  is (1-)strongly log-concave.



Then  $\operatorname{KL}(p_t|\pi) = \exp(-2t) \operatorname{KL}(p_0|\pi)$ .

## When does Langevin diffusion's law converges (fast) to $\pi$ ?

• Consider a standard Gaussian distribution  $\pi(x) \propto \exp(-\frac{||x||^2}{2})$ , i.e.  $\pi \propto \exp(-V)$  with V 1-strongly convex, i.e.  $\pi$  is (1-)strongly log-concave.



Then  $\operatorname{KL}(p_t|\pi) = \exp(-2t) \operatorname{KL}(p_0|\pi)$ .

• If *π* is a perturbation of a strongly-log-concave distribution, then the rate degrades with the size of the perturbation.



(see Holley–Stroock theorem and log-Sobolev inequalities, (Bakry et al., 2014)).

Langevin 00000000000000 Bayesian deep learning 00000000 References

### Langevin in the multimodal case



Mixture of equally weighted 16 Gaussians with unit variance and uniformly chosen centers in  $[-40, 40]^2$ , a standard sampling benchmark. ULA was initialized with  $\mathcal{N}(0, I_2)$ , step-size h = 0.01. ULA was run with  $5.10^4$  steps (one minute run).

Langevin 00000000000000 Bayesian deep learning 00000000 References

### Langevin in the multimodal case



Mixture of equally weighted 16 Gaussians with unit variance and uniformly chosen centers in  $[-40, 40]^2$ , a standard sampling benchmark. ULA was initialized with  $\mathcal{N}(0, I_2)$ , step-size h = 0.01. ULA was run with  $5.10^4$  steps (one minute run).

The theoretical convergence is so slow, that in practice Langevin gets stuck for infinite time the modes close to its initialization !

Bayesian deep learning •0000000 References

#### Outline

Bayesian learning

Langevin

Bayesian deep learning

#### References

#### Recall Bayesian inference

Given labelled data  $(w_i, y_i)_{i=1}^p$ , we want to sample from the posterior distribution over the parameters of a model  $g(\cdot, x)$ 

$$\pi(x) \propto \exp\left(-V(x)
ight), \quad V(x) = \sum_{\substack{i=1 \ \text{loss on labeled data } (w_i, y_i)_{i=1}^p}^p + rac{\|x\|^2}{2},$$

I.e.,  $\pi(x) = \frac{\exp(-V(x))}{Z}$ ,  $V(x) = -\log p(\mathcal{D}|x) - \log p_0(x)$  with Z intractable.

### Recall Bayesian inference

Given labelled data  $(w_i, y_i)_{i=1}^p$ , we want to sample from the posterior distribution over the parameters of a model  $g(\cdot, x)$ 

$$\pi(x) \propto \exp\left(-V(x)
ight), \quad V(x) = \sum_{\substack{i=1 \ \text{loss on labeled data } (w_i, y_i)_{i=1}^p}}^p + \frac{\|x\|^2}{2},$$

I.e.,  $\pi(x) = \frac{\exp(-V(x))}{Z}$ ,  $V(x) = -\log p(\mathcal{D}|x) - \log p_0(x)$  with Z intractable.

Ensemble prediction for an input w:



Predictions of models parametrized by  $x \in \mathbb{R}^d$  are reweighted by  $\pi(x)$ .



### Recall Bayesian inference

Given labelled data  $(w_i, y_i)_{i=1}^p$ , we want to sample from the posterior distribution over the parameters of a model  $g(\cdot, x)$ 

$$\pi(x) \propto \exp\left(-V(x)
ight), \quad V(x) = \sum_{\substack{i=1 \ \text{loss on labeled data } (w_i, y_i)_{i=1}^p}}^p + \frac{\|x\|^2}{2},$$

I.e.,  $\pi(x) = \frac{\exp(-V(x))}{Z}$ ,  $V(x) = -\log p(\mathcal{D}|x) - \log p_0(x)$  with Z intractable.

Ensemble prediction for an input w:



Predictions of models parametrized by  $x \in \mathbb{R}^d$ are reweighted by  $\pi(x)$ .



angevin

Bayesian deep learning 0000000

### Langevin for (Bayesian) deep NN?

Given labelled data  $\mathcal{D} = (w_i, y_i)_{i=1}^p$ , we want to sample from the posterior distribution over the parameters of a model  $g(\cdot, x)$ 

$$\pi(x) \propto \exp\left(-V(x)\right), \quad V(x) = \underbrace{\sum_{i=1}^{p} \left\|y_i - g(w_i, x)\right\|^2}_{\text{loss on labeled data } (w_i, y_i)_{i=1}^{p}} + \underbrace{\frac{\left\|x\right\|^2}{2}}_{\text{prior reg.}}.$$

angevin

Bayesian deep learning 0000000

### Langevin for (Bayesian) deep NN?

Given labelled data  $\mathcal{D} = (w_i, y_i)_{i=1}^p$ , we want to sample from the posterior distribution over the parameters of a model  $g(\cdot, x)$ 

$$\pi(x) \propto \exp\left(-V(x)
ight), \quad V(x) = \underbrace{\sum_{i=1}^{p} \left\|y_i - g(w_i, x)
ight\|^2}_{\text{loss on labeled data } (w_i, y_i)_{i=1}^{p}} + \underbrace{\frac{\left\|x
ight\|^2}{2}}_{\text{prior reg.}}.$$

 Recall that we know that the convergence speed of Langevin diffusion depends on how much "V is convex" and if it has few local minimas angevin

Bayesian deep learning

### Langevin for (Bayesian) deep NN?

Given labelled data  $\mathcal{D} = (w_i, y_i)_{i=1}^p$ , we want to sample from the posterior distribution over the parameters of a model  $g(\cdot, x)$ 

$$\pi(x) \propto \exp\left(-V(x)
ight), \quad V(x) = \underbrace{\sum_{i=1}^{p} \left\|y_i - g(w_i, x)
ight\|^2}_{\text{loss on labeled data } (w_i, y_i)_{i=1}^{p}} + \underbrace{\frac{\left\|x
ight\|^2}{2}}_{\text{prior reg.}}.$$

- Recall that we know that the convergence speed of Langevin diffusion depends on how much "V is convex" and if it has few local minimas
- is  $x \mapsto V(x)$  convex for g(.,x) a neural network parametrized by x?

### Langevin for (Bayesian) deep NN?

Given labelled data  $\mathcal{D} = (w_i, y_i)_{i=1}^p$ , we want to sample from the posterior distribution over the parameters of a model  $g(\cdot, x)$ 

$$\pi(x) \propto \exp\left(-V(x)
ight), \quad V(x) = \underbrace{\sum_{i=1}^{p} \left\|y_i - g(w_i, x)
ight\|_{i=1}^2}_{\text{loss on labeled data } (w_i, y_i)_{i=1}^p} + \underbrace{\frac{\left\|x
ight\|_{i=1}^2}{2}}_{\text{prior reg.}}.$$

- Recall that we know that the convergence speed of Langevin diffusion depends on how much "V is convex" and if it has few local minimas
- is x → V(x) convex for g(.,x) a neural network parametrized by x?



A highly nonconvex loss surface, as is common in deep neural nets. From https://www.telesens.co/2019/01/16/neural-network-loss-visualization.

#### Different strategies in practice/in the literature

Close to what we've seen previously:

- Stochastic Langevin dynamics: approximate  $\nabla V(x) = \nabla \left( \sum_{i=1}^{p} \|y_i - g(w_i, x)\|^2 + \frac{\|x\|^2}{2} \right) \text{ by a batch of data samples}$   $(w_i, y_i)_{i=1}^m \text{ with } m \ll p$
- Variational Inference

find 
$$q_{\theta} = \operatorname*{arg\,min}_{p \in P_{\theta}} \operatorname{KL}(p|\pi)$$

where  $P_{\theta}$  is a family of parametric distributions (upcoming in few slides).

### Different strategies in practice/in the literature

#### More heuristic:

#### Monte Carlo Dropout

*Gal*, Y., & Ghahramani, Z. (2016). Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In international conference on machine learning.

#### Deep ensembles

Lakshminarayanan, B., Pritzel, A., & Blundell, C. (2017). Simple and scalable predictive uncertainty estimation using deep ensembles. Advances in neural information processing systems.



#### Variational Inference for BNN - Bayes by Backprop example

Variational Inference

find 
$$q_{ heta} = rgmin_{p \in P_{ heta}} \operatorname{KL}(p|\pi)$$

where  $P_{\theta}$  is a family of parametric distributions.

### Variational Inference for BNN - Bayes by Backprop example

Variational Inference

find 
$$q_{ heta} = \operatorname*{arg\,min}_{p \in P_{ heta}} \operatorname{KL}(p|\pi)$$

where  $P_{\theta}$  is a family of parametric distributions.

A typical neural network of depth L (with non-linearity  $h(\cdot)$ ) for input w and parameter x writes:

$$g(w,x) = A^{L}h\left(A^{L-1}h\left(\ldots h\left(A^{1}w+b^{1}\right)\right)+b^{L-1}\right)+b^{L},$$

$$h' = h(A'h'^{-1} + b'), \quad h^1 = h(A^1w + b^1).$$

Neural network parameters:  $x = \{A', b'\}_{l=1}^{L}$ .

#### We will describe the approach of "Bayes by Backprop"<sup>1</sup>.

Blundell, C., Cornebise, J., Kavukcuoglu, K., & Wierstra, D. (2015). Weight uncertainty in neural network. In International conference on machine learning.

Bayesian deep learning 000000000

#### Step 1: Construct the $q_{\theta}(x) \approx p(x \mid D) = \pi(x)$ Distribution Example: Mean-field (="factorized") Gaussian distribution:

$$q_{\theta} = \prod_{l=1}^{L} q(A') q(b')$$

.

$$q(A_{i}) = \prod_{ij} q(A_{ij}^{l}), \quad q(A_{ij}^{l}) = \mathcal{N}(A_{ij}^{l}; M_{ij}^{l}, V_{ij}^{l})$$
$$q(A_{ij}^{l}) = \prod_{ij} q(A_{ij}^{l}), \quad q(A_{ij}^{l}) = \mathcal{N}(A_{ij}^{l}; M_{ij}^{l}, V_{ij}^{l})$$

$$q(b') = \prod_{i} q(b'_{i}), \quad q(b'_{i}) = \mathcal{N}(b'_{i}; m'_{i}, v'_{i})$$

Variational parameters:  $\theta = \{M_{ij}^{l}, V_{ij}^{l}, m_{i}^{l}, v_{i}^{l}\}_{l=1}^{L}$ 





In dimension two, a simple example of  $q_{\theta}$  is a factorized Gaussian:

$$q_{\theta}(A_{11}^{1}, A_{12}^{1}) = \mathcal{N}(A_{11}^{1}; 0, 1) \cdot \mathcal{N}(A_{12}^{1}; 0, 1),$$

where  $q_{\theta}$  is the product of two independent standard normal distributions over the parameters  $A_{11}^1$  and  $A_{12}^1.$ 

Note that the "factor" assumption in mean-field decorrelates variables.
## Step 2: Fit the $q_{\theta}$ Distribution

**Variational inference:**  $\theta^* = \arg \max L(\theta)$  where L is the ELBO

$$L(\theta) = \mathbb{E}_{q_{\theta}}[\log p(D \mid x)] - \mathrm{KL}[q_{\theta} \parallel p_0(x)]$$

First scalable technique: Stochastic optimization

- i.i.d. assumption:  $\log p(D \mid x) = \sum_{i=1}^{N} \log p(y_i \mid w_i, x)$
- Mini-batch training:  $\{(w_m, y_m)\}_{m=1}^M \sim D^M$

$$L(\theta) pprox rac{N}{M} \sum_{i=1}^{M} \mathbb{E}_{q_{\theta}}[\log p(y_i \mid w_i, x)] - \mathrm{KL}[q_{\theta} \parallel p_0(x)]$$

Reweighting to ensure calibrated posterior concentration.

## Step 2: Fit the $q_{\theta}$ Distribution Variational inference: $\theta^* = \arg \max L(\theta)$ where L is the ELBO

$$L(\theta) = \mathbb{E}_{q_{\theta}}[\log p(D \mid x)] - \mathrm{KL}[q_{\theta} \parallel p_0(x)]$$

2nd Scalable Technique: Monte Carlo Sampling

- $\mathbb{E}_{q_{\theta}}[\log p(y \mid w, x)]$  is intractable even with Gaussian  $q_{\theta}$
- Solution: Monte Carlo estimate:

$$\mathbb{E}_{q_{\theta}}[\log p(y \mid w, x)] \approx \frac{1}{K} \sum_{k=1}^{K} \log p(y \mid w, x_k), \quad x_k \sim q_{\theta}$$

• Reparameterization trick to sample from mean-field Gaussians:

$$x_k = m_{\theta} + \sigma_{\theta} \odot \epsilon_k, \quad \epsilon_k \sim \mathcal{N}(0, I)$$



• Therefore:

$$\mathbb{E}_{q_{\theta}}[\log p(y \mid w, x)] \approx \frac{1}{K} \sum_{k=1}^{K} \log p(y \mid w, x_k), \ x_k = m_{\theta} + \sigma_{\theta} \epsilon_k$$

## Combining both steps and final prediction Full ELBO approximation:

$$L(\theta) \approx \frac{N}{M} \sum_{m=1}^{M} \frac{1}{K} \sum_{k=1}^{K} \log p(y_m \mid w_m, x_k) - \mathrm{KL}[q_\theta \parallel p(x)], \quad x_k \sim q_\theta$$

analytic between two Gaussians (if not, can also be estimated with Monte Carlo)

In regression:  $p(y \mid w, x) = \mathcal{N}(f_x(w), \sigma^2)$ , In classification:  $p(y \mid w, x) = \text{Categorical}(\text{logit} = f_x(w))$ 

Step 3: Compute Prediction with Monte Carlo Approximations

$$p(y^* \mid w^*, D) pprox rac{1}{\mathcal{K}} \sum_{k=1}^{\mathcal{K}} p(y^* \mid w^*, x_k), \quad x_k \sim q_k$$

<u>Mean-field Gaussian case</u>:  $x_k = m_\theta + \sigma_\theta \odot \epsilon_k$ ,  $\epsilon_k \sim \mathcal{N}(0, I)$ 



## References I

- Bakry, D., Gentil, I., Ledoux, M., et al. (2014). *Analysis and geometry of Markov diffusion operators*, volume 103. Springer.
- Roberts, G. O. and Tweedie, R. L. (1996). Exponential convergence of langevin distributions and their discrete approximations. *Bernoulli*, pages 341–363.