

Contents

1	Escaping local minima with Langevin Monte Carlo (Yuchen Zhang)	1
2	Nonasymptotic Analysis of SGLD (Matus Telgarsky)	5
2.1	Basic setting and goal	5
2.2	Analysis overview	7
2.3	Tricks from the analysis	9

1 Escaping local minima with Langevin Monte Carlo (Yuchen Zhang)

Joint work with Moses Charikar and Percy Liang.

The goal is to minimize a nonconvex function $f : K \rightarrow \mathbb{R}$. Nonconvex optimization is universal in ML applications: mixture models, matrix/tensor decomposition, neural networks, etc.

The main challenge is to avoid non-optima local minima. Doing this is NP-hard in general, but we can hope to escape shallow local optima.

If we properly initialize gradient descent we converge to local min. It will get stuck at local min because it is a greedy method, only looking at the neighborhood.

A natural approach is to inject noise,

$$x \leftarrow x - \eta \cdot (\nabla f(x) + \sigma \cdot w), \quad w \sim N(0, I).$$

One thing to pay attention to is how to choose step size. If you choose very small step size, you lose capability of escaping local min. It's 0-mean so this reduces to regular gradient descent. Choose large enough noise to escape before it gets canceled out. But this causes stability problems: overshooting, divergence, especially if f is not very smooth.

Solutoin: Langevin Monte Carlo. Attracting attention because of connection to successful SGD. Instead of scaling with step size, scale with $\sqrt{\eta}$ step size. It imitates Langevin diffusion in physics. $1/\xi$ controls temperature.

$$x \leftarrow x - \eta \cdot \nabla f(x) + \sqrt{\frac{2\eta}{\xi}} w.$$

Note when $\eta \rightarrow 0$, $\sqrt{\frac{\eta}{\xi}} \gg \eta$, still able to escape.

Why scale with $\sqrt{\eta}$? Only by this choice does it converge to the stationary distribution, $\mu(x) \propto e^{-\xi f(x)}$. When ξ is large and $x \sim \mu(x)$, x minimizes f . This is good for nonconvex optimization. You need f to be smooth and $e^{-\xi f}$ to be integrable. ¹

I assume f is smooth and bounded on K , and K is compact. It doesn't need to be convex, but I make assumptions on the corner (no sharp corners). I think this is not necessary but is convenient.

¹There is a noncompact analysis in 80's, purely asymptotic. For a compact set you can do reflected Brownian motion, or reject steps outside.

For smooth functions and small enough stepsize, LMC asymptotically converges (Roberts, Tweedie 96). For convex f , LMC converges to μ in polytime (Bubeck 2015, Dalalyan 2016). For nonconvex functions LMC converges in exponential time, for simple 1-D nonconvex functions. This is an area that is underexplored.

On the practical side, LMC is successful in practice.

- Prevents overfitting: logistic regression, ICA (Welling Teh 2011)
- Learn deep nets: neural programmer (Neelakantan 2015), neural RAM (Kurach 2015), GPUs (Kaiser, Sutskever 2015), bidirectional LSTM (Zeyer 2016).

How to explain good performance?

LMC can hit a good solution much earlier than it converges to the stationary distribution. Ex. W-shaped function. Sampling means finding all good solutions; optimization means finding one.

The reason we don't have a polytime guarantee is that mixing time is too pessimistic for optimization. We look at hitting times to a set U (target set) instead. If we have an arbitrary defined set, we get into NP-hard problems. Whether we get a polytime bound depends on how we define this set.

Empirical risk can have many local min, but can have shallow local min not present in the global min. We want to get to population local min.

An ε -approximate local min is $\{x \in K : \|\nabla f(x)\| \leq \varepsilon, \nabla f(x) \succeq -\sqrt{\varepsilon}I\}$ Can also get by cubic regularized Newton, noisy gradient descent. *Here we can achieve this for population risk even if we only get access to empirical risk!*

We have bad complexity in d in number of iterations.

We only need the Lipschitz constant for population loss function to be small.

We present non-asymptotic analysis for the LMC algorithm (and stochastic version SGLD) on general nonconvex functions, polynomial-time guarantee for hitting certain optimality set, and apply to ERM finding local minima of the population risk, and to learning halfspaces under 0-1 loss (note there is no guarantee on concentration of gradient and Hessian matrices, which are 0 almost everywhere).

First, some definitions and notations.

For $f : K \rightarrow \mathbb{R}$, define

$$\mu_f(x) := \frac{e^{-f(x)}}{\int_K e^{-f(x)}} \quad (1)$$

$$\mu_f(\partial A) := \lim_{\varepsilon \rightarrow 0} \frac{\mu_f(\{x \in K : d(x, A) \leq \varepsilon\}) - \mu_f(A)}{\varepsilon} \quad (2)$$

$$C_f(V) := \inf_{A \subseteq V} \frac{\mu_f(\partial A)}{\mu_f(A)} \quad (3)$$

(restricted Cheeger constant), min ratio betw

If $C_f(V)$ is very small, then there exists some set $A \subseteq V$ with a very narrow exit. If there is a Markov process with stationary distribution μ_f and it was initialized in A , then it takes a long time to move out. This defines geometric property of function, cf. conductance which is property on Markov process.

Restricted Cheeger constant $C_{\xi f}(V)$ measures how fast LMC can escape V . Ex. for the W function, the restricted Cheeger constant for neighborhood of one local min is small.

If you run too long, won't the bias accumulate? I can choose step large enough so we hit before bias accumulates too much. With small Cheeger constant I have to choose small step size.

Theorem 1.1. *For arbitrarily smooth f and $U \subseteq K$, LMC (with some set size and fixed temperature) hits U whp in T iterations, where*

$$T \leq \frac{\text{poly}(d)}{C_{(\xi f)}(K \setminus U)^4}.$$

Initialization is arbitrary. If I choose small step size, I need more iterations. Note that Cheeger constant depends on temperature. Temperature is chosen according to effective dimension of problem, volume of set you want to hit. (Choose proportional to $\frac{1}{d}$ in worst-case. If region around manifold, then you can choose temperature larger.)

Given the function and the set, there exists a step size. Step size: (28) on p. 17.

Lower bound: inversely proportional to conductance. Upper bound is at most square of reciprocal of conductance. This is a mixing time lower bound. If you look at a convex function it is the same thing. The bound we have is square of mixing time for convex case.

Is this an algorithm that can be implemented? There is a theoretical step size, but in practice we choose the step-size by cross-validation.

MT: We go through mixing time lemma, to Gibbs measure, as corollary get suboptimality. Ours is global optimal guarantee. We pay exponentially in the saddle height. We also have $\frac{1}{\varepsilon^4}$. If you choose the step size and temperature, then things become bad. There is a lower bound.

You can make function grow quadratically; choose temperature schedule.

Proof. • Construct a time-reversible Markov chain; prove hitting time on par with LMC. (LMC is not reversible.)

- Prove hitting time inversely depends on reverse conductance.

(Look at specification of algorithm, specifics of Markov chain.)

- Restricted conductance is lower bounded by $(C_{(\xi f)}(K \setminus U))^2$.

□

This is a general theorem that doesn't give a polytime guarantee.

Now we look at certain cases where we can show Cheeger constant is not exponentially small.

A simple class of functions is convex functions. Choosing temperature to be low enough,

Proposition 1.2: If $\xi \geq O\left(\frac{d}{\varepsilon^2}\right)$, $C_{(\xi f)}(K \setminus U)$ is lower bounded by $\Omega(1)$ for U set of ε -approximate global min,

$$U = \left\{x : f(x) \leq \inf_{x \in K} f(x) + \varepsilon\right\}.$$

More interesting case is nonconvex.

Proposition 1.3: If $\xi \geq O\left(\frac{\text{poly}(\text{params})}{\varepsilon^2}\right)$, then C lower bounded by $\Omega(\sqrt{\varepsilon})$, where

$$U := \{x : \|\nabla f(x)\|_2 \leq \varepsilon, \nabla^2 f(x) \succeq -\sqrt{\varepsilon}I\}.$$

There is implicit dependence on dimension, trace of hessian. Local minimum has 0 gradient, saddle point also has 0 gradient but H is non-PSD. U avoids strict saddle points. d^4 is impractical but there has to be $\text{poly}(d)$ dependence.

If f is nonconvex and uniformly close to F :

Proposition 1.4: $\|f - F\|_\infty := \sup_{x \in K} |f(x) - F(x)| \leq \nu$,

$$C_{(\xi f)}(K \setminus U) \geq e^{-2\xi\nu} C_{(\xi F)}(K \setminus U).$$

If $\|f - F\|_\infty = O\left(\frac{1}{\xi}\right)$ then $C_{(\xi f)}(K \setminus U) = C_{(\xi F)}(K \setminus U)$.

Ex. 0-1 loss. First define smoothed version on 0-1 loss. Guarantee smoothed version of function is uniformly close. But we cannot guarantee all local minima are close to population global minimum.

Corollary 1.5. *Run LMC on smooth f .*

1. If F convex, $\|f - F\|_\infty = O\left(\frac{\varepsilon^2}{d}\right)$, then LMC hits ε -approximate global min of F in poly time.
2. If F is smooth and $\|f - F\|_\infty = P\left(\frac{\varepsilon^2}{(\text{params})}\right)$, then LMC hits ε -approx local min of F in poly time.

If f is nonconvex, nonsmooth,

1. define $\tilde{f}_\sigma(x) = \mathbb{E}_z f(x + z)$ for $z \sim N(0, \sigma^2 I)$
2. Run LMC on \tilde{f}_σ ,

$$\nabla \tilde{f}_\sigma(x) = \mathbb{E}_z \left[\frac{z}{\sigma^2} (f(x + z) - f(x)) \right]$$

(only depends on values of f).

For any $\sigma > 0$, \tilde{f}_σ smooth. We can choose σ small enough so $\|f - \tilde{f}_\sigma\|_\infty = O\left(\frac{1}{\xi}\right)$, which implies $C_{(\xi f)}(K \setminus U) \approx C_{(\xi \tilde{f}_\sigma)}(K \setminus U)$.

I care less about computational, more about sample complexity. Choice of σ only affects computational complexity.

Do I need to choose large enough σ so that smoothness is similar to original function? No, only need function value to be close.

In small neighborhood of discontinuity, there constant probability of choosing z that makes it jump to the opposite side. Assume jump is $\leq \frac{1}{\xi}$. (Use upper bound on VC dimension.)

Apply to $f(x) = \frac{1}{n} \sum_{i=1}^n l(x; a_i)$ and $F(x) = \mathbb{E}_{a \sim \mathbb{P}}[l(x; a)]$. Under mild conditions $\|f - F\|_\infty \rightarrow 0$ as $n \rightarrow \infty$. However, ∇, ∇^2 don't converge unless l is smooth, so GD is unreliable.

Can I just do analysis on the smoothed function?

There are 2 notions of samples. I assume I have a set of n samples, and I can take many passes over it, which is relevant to computation cost.

LMC is not necessarily better than smoothed version of gradient descent—you can also eliminate local minimum.

Add noise in heuristic scaling to gradient in neural networks helps. Would be interesting to compare. People claim LMC is better because they didn't push hard enough in other methods. LMC has enough power to escape, but batch normalization can also avoid getting into bad situation. It's hard to say in practice what is the best procedure.

Learning halfspace with 0-1 loss. LMC learns halfspace in polytime for arbitrary noise level. Better than best previous result (Awasthi 2015, learn halfspace in polytime if noise $\leq 10^{-6}$). Here $l(x; (a, y)) = \mathbb{1}(\text{sign}(\langle x, a \rangle) \neq y)$ for $(a, y) \in \mathbb{R}^d \times \{\pm 1\}$. See Theorem 3 in paper, p. 11. Choose $n \geq \tilde{O}(1) \frac{d^4}{\delta_0^2 \epsilon^4}$.

(Here the population risk is not smooth. This is a consequence of a more general theorem that doesn't need smoothness on population risk.)

Summary: LMC asymptotically optimal for nonconvex optimization but convergence rates not well understood. We prove hitting inversely depends on restricted Cheeger constant. We lower bound restricted Cheeger constant for convex functions and smooth nonconvex functions. If $\|f - F\|_\infty$ small, then running LMC on f achieves optimal points of F (stability property). LMC is more reliable in GD in empirical risk minimization because it escapes shallow/tiny local min. It is asymptotically consistent too.

2 Nonasymptotic Analysis of SGLD (Matus Telgarsky)

Joint work with Maxim Raginsky and Alexander Rakhlin.

I tried to use the Lovasz-Vempala papers to do nonconvex optimization. You don't need convexity in the stationary distribution proof. What could go time if I try to prove a mixing time? I pushed through the proofs and got a nonconvex analysis. All the mixing time bounds were uninterpretable; I don't know when it's polynomial. I couldn't make them work for tensor problems.

This project has been eye-opening; these techniques (stochastic differential equations) are powerful. The Lovasz-Vempala papers are nasty; you shouldn't have to be so clever...

2.1 Basic setting and goal

$$w_{k+1} = w_k - \eta g_{k+1} - \sqrt{\frac{2\pi}{\mu}} \xi_k$$

where $\xi_k \sim N(0, I)$. Assume conditional unbiased estimate of true gradient $\mathbb{E}[g_{k+1}|w_k] = \nabla f(w_k)$. It's more difficult than I expected to convert this into a population gradient. Suppose $w_0 \sim p_0$ and let p_k be the law of w_k .

The stationary distribution (using the true gradient) is $d\pi(w) \propto \exp(-\beta f(w))$. The

Wasserstein distance is

$$W_2(\mu, \nu) = \inf \left\{ \sqrt{\mathbb{E} \|w - w'\|_2^2} : \mu = \text{Law}(w), \nu = \text{Law}(w') \right\}.$$

I'll highlight why we're using W_2 distance; this is a key part of this paper.

There are 2 tensions in the analysis: the more time that passes, the better the walk mixes; but the more the discretized and continuous version drift from each other.

Informal theorem is as follows:

Theorem 2.1. *Suppose f satisfies “some regularity conditions”. Using $K = \Theta\left(\frac{\beta(d+\beta)}{\lambda\epsilon^4}\right)$, $\eta = \Theta(\epsilon^4)$. Then the Wasserstein distance between the discretized walk and the stationary distribution is*

$$W_2(\mu_k, \pi) = \tilde{O}\left(\frac{\beta(d+\beta)^2}{\lambda}(\dots + \epsilon)\right)$$

2

As a corollary we get an optimization guarantee.

Corollary 2.2.

$$\mathbb{E}f(w_k) - \inf_w f(w) = O(W_2(\mu_k, \pi) + \tilde{O}\left(\frac{d}{R}\right))$$

Remark 2.3: 1. In general the only bound we have is exponential, $\frac{1}{\lambda} = \exp(O(\beta + d))$.

2. The arXiv version has a fixed training set. (In general you're running so many iterations that you won't have enough samples.)

I'll give analysis for true stochastic gradients so we can use the population spectral gap.

3. Dalalyan uses a TV bound. It's not clear how to get the corollary with only a TV bound. You can't control where the difference in the mass between the probability distribution goes—it could be moved arbitrarily far away. For W_1 you would not be able to do it with smooth f .

I will return to Dalalyan's proof. I'll step through where we use the same vs. different proofs.

Full assumptions are

- $|f(0)| \leq C$, $\|\nabla f(0)\| \leq C$.
- $\|\nabla f(v) - \nabla f(w)\| \leq C\|v - w\|$.
- (from control theory literature) Dissipativity: for all w ,

$$\langle w, \nabla f(w) \rangle \geq \frac{1}{C} \|w\|^2 - C.$$

It's almost saying you're lower bounded by a quadratic.

²hiding numerical constants, $\ln\left(\frac{1}{\epsilon}\right)$, $\ln\beta$, smoothness constants not depending on dimension.

•

$$\mathbb{E} \|\hat{g}_{k+1} - \nabla f(w_k)\|_2^2 \leq \sigma^2$$

•

$$\int \exp(\|\omega\|^2) p_0(w) dw \leq C.$$

• Spectral gap

$$\lambda := \inf \left\{ \frac{\mathbb{E}_\pi \|\nabla g\|^2}{\mathbb{E}_\pi g^2} : g \in C^1(\mathbb{R}) \cap L_2(\pi), g \neq 0, \int g d\pi = 0 \right\}.$$

Why use this mixing parameter? We have a bound $\frac{1}{\lambda} = \exp(O(\beta + d))$. Previous analyses were asymptotic.

More remarks:

- There was a paper by Yann LeCun on entropy-SGD. They're essentially doing Langevin on a smoothed function

$$\tilde{f}(w) := \frac{1}{2} \|w\|^2 - \ln \int_{\|v\| \leq R} e^{\langle v, w \rangle - \|w\|^2 - \beta f(v)} dv$$

with $\frac{1}{\lambda} \leq O(\exp(\beta R^2))$.

- λ becomes exponential in pretty mild situations, ex. disconnected local optima, $\frac{1}{\lambda} = \exp(\Omega(\beta))$. Google “Eyring-Kramois formula” and “Bovier-Guyard-Klein”. Langevin is “metastable”; it is attracted to local optima.
- You can get around this using homotopy (change function), annealing method (change temperature schedule, iterate at previous time is a warm start), or warm start. Our analysis doesn't include this.

All analyses for Langevin now are fixed-temperature.

- Maybe global optimization is the wrong story, a red herring. (Ex. hitting time analysis.) Wide basins of attraction generalize better?

2.2 Analysis overview

The way this has been analyzed since forever is to look at the continuous time SDE. See Gelfand-Mitter 1991. The step size disappears here.

$$dw(t) := -\nabla f(w(t))dt + \sqrt{\frac{2}{\beta}} d\beta(t) \tag{4}$$

$$W_t := \text{Law}(w_t). \tag{5}$$

The strategy is

$$W_2(\mu_k, \pi) \leq W_2(\mu_k, \nu_{k\eta}) + W_2(\nu_{k\eta}, \pi). \tag{6}$$

For the first term, the discretization error,

$$W_2(\mu_k, \nu_{k\eta}) \leq \widetilde{C}_{\nu_{k\eta}} (\sqrt{D(\mu_k || \nu_{k\eta})} + D(\mu_k || \nu_{k\eta})^{\frac{1}{4}}) \quad (7)$$

This is an entire paper by Bolley and Villani (2005). The fact you can do this at all for Wasserstein is surprising. Here \widetilde{C} is something like $\ln \int \exp(\|w\|^2) dw$.

Dissipativity says that the chain quickly stays within a small ball. Projection would screw up this analysis.

Like the Dalalyan paper (but with some more finickyness) we get $D(\mu_k || \nu_{k\eta})$ small (Use Girsanov's theorem, useful for KL divergence).

For the discretization, this is what you have to do, in addition to algebra.

There are 2 version of the Bolley-Villani inequality, one without $\sqrt{\cdot}$ term; we can't bound the \widetilde{C} in that inequality.

Here

$$\widetilde{C}_\nu := 2 \inf_{\lambda > 0} \left(\frac{1}{\lambda} \left(\frac{3}{2} + \ln \int \exp(\lambda \|w\|^2) d\nu \right) \right).$$

This is hard to understand; we bound with $\lambda = 1$.

The second term we have to control is the diffusion error $W_2(\nu_{k\eta}, \pi)$, which should go to 0 as $k \rightarrow \infty$. This part differs from the Dalalyan paper. We have 3 steps. We have to use a log-Sobolev constant.

1. $\lambda_{LS} \leq O\left(\frac{1}{\beta} + \frac{d+\beta}{\lambda}\right)$. The reverse direction is natural (Taylor expansion). Going this direction is difficult; bounding LS by spectral gap is hard and obscure.
2. (Similar to Dalalyan paper if you are in TV case.)

$$W_2(\nu_{k\eta}, \pi) \leq \sqrt{2\lambda_{LS} D(\nu_{k\eta} || \pi)}.$$

This is the Otto-Villani theorem. See Bakry-Gentil-Ledoux book, Theorem 9.6.1. We need the log-Sobolev constant another time to control mixing time.

3. $D(\nu_{k\eta} || \pi) \leq D(\nu_0 || \pi) \exp\left(-\frac{2k\eta}{\beta\lambda_{LS}}\right)$, BGL Theorem 5.2.1. Entropy decays exponentially fast if you have LS constant.

Sketch: Let n be strong convexity.

$$TV(\nu_{k\eta}, \pi) \leq \frac{1}{2} \chi^2(\nu_{k\eta} || \pi) \exp(-k\eta \cdot m/2), \quad (8)$$

The RHS is BGL Corollary 4.8.2. You get Poincaré constant from strong convexity.

We want to control

$$\mathbb{E}f(w_k) - \inf_w f(w) \leq (\mathbb{E}f(w_k) - \mathbb{E}_\pi f(w)) + (\mathbb{E}_\pi f(w) - \inf_w f(w)) \quad (9)$$

$$\leq O(W_2(\mu_*, \pi)) \quad (10)$$

Here is why we use W_2 , to control functions which are smooth. There's an interesting trick here. To analyze this term, expand into entropy and log-partition function

$$\int f(w) d\pi = \frac{1}{\beta} (\text{Ent}(\pi) - \ln \int \exp(-\beta f(w)) dw)$$

W_2 implies weak convergence and 2nd moment convergence. (For W_p , p th moment.) For the first term: Upper-bound variance of $\text{Ent}(\nu_k)$ to upper bound $\text{Ent}(\pi)$, then show it's maximized for gaussians. The second term is bounded by algebra.

In summary we compare to Dalalyan:

1. Use W_2 to control variation of smooth functions.
2. Rest looks similar, convert W_2 to KL by heavy-lifting Villani papers.
3. What's scary is controlling log-Sobolev constant and using the high-power theorems involving it. Our assumption was on $\frac{1}{\lambda} = \exp(O(\beta+d))$ to give a completely nonasymptotic analysis.

2.3 Tricks from the analysis

How to do the discretization? Something I tried which was a bad idea: analyze W_2 without any technology.

$$\mu_{k+1} := \mu_k - \eta g_{k+1} + \sqrt{\frac{2\eta}{\beta}} \xi_{k+1} \quad (11)$$

$$\nu_{k+1} := \nu_k - \nabla f(w_k) + \sqrt{\frac{2\eta}{\beta}} \xi_{k+1} \quad (12)$$

Use explicit coupling that makes things cancel: the same Gaussian. This doesn't work—these things really start to drift.

Instead, interpolate between w_k and $w(k\eta)$. Fix stochastic gradients for chunks of time. Define a fixed stochastic gradient where I made randomness concrete: $\hat{g}(w, l)$ a deterministic function where l is the source of randomness. Then

$$\mathbb{E}_l(\hat{g}(w, l)) = \nabla f(w).$$

Here l_{k+1} is randomness at time k , $l(t) = l_k$ for $t \in [k\eta, (k+1)\eta)$.

$$w_{k+1} := w_k - \eta \hat{g}(w_k, l_{k+1}) + \sqrt{\frac{2\eta}{\beta}} \xi_{k+1} \quad (13)$$

$$w(t) = w_0 - \eta \int_0^t \hat{g}(u(\lfloor \frac{s}{\eta} \rfloor \eta), l(s)) ds + \sqrt{\frac{2}{\beta}} \int_0^t d\beta(s) \quad (14)$$

$$v(t) = w_0 - \eta \mathbb{E}[\hat{g}(\lfloor \frac{s}{\eta} \rfloor \eta, l(s)) | u(s) = v(s)] + \sqrt{\frac{2}{\beta}} \int_0^t d\beta(s). \quad (15)$$

Problem: $w(t)$ is not a Markov chain. Find a 2nd chain where we get rid of the stochastic gradient.

Let

$$p_v^* := \text{Law}(v(s) : 0 \leq s \leq t) \quad (16)$$

$$p_w^* := \text{Law}(w(s) : 0 \leq s \leq t) \quad (17)$$

This looks gross, but you can algebra it away.

Cf. conditional independence given iterate, $\langle \nabla f - g, w_k - \bar{w} \rangle$. Use tower property. We use an “infinitesimal” tower property.

Get

$$D(p_v^* || p_w^*) = \frac{\beta}{4} \int_0^t \mathbb{E} \left\| \nabla f(u(s)) - \mathbb{E}(\hat{g}(u(\left\lfloor \frac{s}{\eta} \right\rfloor \eta), l(s)) | U(s) = v(s)) \right\|^2 ds \quad (18)$$

$$D(\mu_k || \nu_{k\eta})^{\frac{1}{4}} \leq D(p_v^t || p_w^t) \quad (19)$$

This is a superpowered data-processing inequality.

Cf. SGD noise: to kill it need decreasing step size. In Dalalyan, the assumption on initial distribution is only used for the mixing time.