Unraveling the mysteries of stochastic gradient descent on deep neural networks

Pratik Chaudhari



UCLA VISION LAB

The question



Stochastic gradient descent

$$x_{k+1} = x_k - \eta \,\nabla f_{\mathcal{C}}(x_k)$$

Many, many variants:

AdaGrad, rmsprop, Adam, SAG, SVRG, Catalyst, APPA, Natasha, Katyusha...





Empirical evidence: wide "minima"



A bit of statistical physics

Energy landscape of a binary perceptron



Wide minima are a large deviations phenomenon

Tilting the Gibbs measure

Local Entropy [Chaudhari et al., ICLR '17]

$$x^* = \underset{x}{\operatorname{argmin}} f(x)$$

$$= \underset{x}{\operatorname{argmin}} e^{-f(x)}$$

$$\approx \underset{x}{\operatorname{argmin}} - \log \left(G_{\gamma} * e^{-f(x)} \right)$$
Gaussian kernel
of variance γ

Parle: parallelization of SGD

State-of-the-art performance [Chaudhari et al., SysML '18]



Wide-ResNet: CIFAR-10

All-CNN: CIFAR-10 (25% data)

The question

Why is SGD so special?

A continuous-time view of SGD

Diffusion matrix: variance of mini-batch gradients

$$\operatorname{var}(\nabla f_{\mathscr{C}}(x)) = \frac{D(x)}{\mathscr{C}}$$
$$= \frac{1}{\mathscr{C}} \left(\frac{1}{N} \sum_{k=1}^{N} \nabla f_{k}(x) \nabla f_{k}(x)^{\top} - \nabla f(x) \nabla f(x)^{\top} \right)$$

Temperature: ratio of learning rate and step-size

$$\beta^{-1} = \frac{\eta}{2\theta}$$

A continuous-time view of SGD

Continuous-time limit of discrete-time updates

$$dx = -\nabla f(x) \underbrace{dt}_{\triangleq \eta} + \sqrt{2\beta^{-1}D(x)} dW(t)$$

will assume $x \in \Omega \subset \mathbb{R}^{a}$

Fokker-Planck (FP) equation gives the distribution on the weight space induced by SGD

$$\rho_{t} = \operatorname{div}\left(\underbrace{\nabla f \rho}_{\operatorname{drift}} + \underbrace{\beta^{-1}\operatorname{div}(D \rho)}_{\operatorname{diffusion}}\right) \quad \text{where } x(t) \sim \rho(t)$$

Wasserstein gradient flow

► Heat equation $\rho_t = \operatorname{div}(\mathbf{I} \nabla \rho)$ performs steepest descent on the Dirichlet energy

$$\frac{1}{2}\int_{\Omega}\left|\nabla\rho(x)\right|^2 dx$$

It is also the steepest descent in the Wasserstein metric for

$$-H(\rho) = \int_{\Omega} \log \rho \ d\rho$$

$$\rho_{k+1}^{\tau} \in \underset{\rho}{\operatorname{argmin}} \left\{ -H(\rho) + \frac{\mathbb{W}_{2}^{2}(\rho, \rho_{k}^{\tau})}{2\tau} \right\}$$

converges to trajectories
of the heat equation

Negative entropy is a Lyapunov functional for Brownian motion

$$ho_{ ext{heat}}^{ ext{ss}} = rgmin_{
ho} - H(
ho)$$

Wasserstein gradient flow: with drift

• If D = I, the Fokker-Planck equation

$$\rho_t = \operatorname{div} \left(\nabla f \rho + \beta^{-1} I \nabla \rho \right)$$

has the Jordan-Kinderleher-Otto (JKO) functional [Jordan et al., '97]

$$\rho^{\rm ss}(x) = \underset{\rho}{\operatorname{argmin}} \underbrace{\mathbb{E}_{x \sim \rho}[f(x)]}_{\text{energetic term}} - \underbrace{\beta^{-1} H(\rho)}_{\text{entropic term}}$$

as the Lyapunov functional.

FP is the steepest descent on JKO in the Wasserstein metric

What happens for non-isotropic noise?

$$\rho_{t} = \operatorname{div}\left(\underbrace{\nabla f \rho}_{\operatorname{drift}} + \underbrace{\beta^{-1}\operatorname{div}(D \rho)}_{\operatorname{diffusion}}\right)$$

FP monotonically minimizes the free energy

$$\rho^{\rm ss}(x) = \underset{\rho}{\operatorname{argmin}} \mathbb{E}_{x \sim \rho} \left[\Phi(x) \right] - \beta^{-1} H(\rho)$$

Rewrite as

$$F(\rho) = \beta^{-1} \mathsf{KL}(\rho \parallel \rho^{\mathsf{ss}})$$

compare with $|x - x^*|$ for deterministic optimization.

SGD performs variational inference

Theorem [Chaudhari & Soatto, ICLR '18]

The functional

$$F(\rho) = \beta^{-1} \mathsf{KL}(\rho \parallel \rho^{\mathsf{ss}})$$

is minimized monotonically by trajectories of the Fokker-Planck equation

$$\rho_t = \operatorname{div} \left(\nabla f \rho + \beta^{-1} \operatorname{div} \left(D \rho \right) \right)$$

with ρ^{ss} as the steady-state distribution. Moreover,

$$\Phi = -\beta^{-1} \log \rho^{\rm ss}$$

up to a constant.

Some implications

Learning rate should scale linearly with batch-size

$$\beta^{-1} = \frac{\eta}{2\ell}$$
 should not be small

also generalizes better

Sampling with replacement regularizes better than without

$$\beta_{\rm w/o\ replacement}^{-1} = \frac{\eta}{2\ell} \left(1 - \frac{\ell}{N}\right)$$

Minimize mutual information of the representation with the training data [Tishby '99, Achille & Soatto '17]

$$\mathsf{IB}_{\beta}(\theta) = \mathbb{E}_{x \sim \rho_{\theta}} \left[f(x) \right] - \beta^{-1} \mathsf{KL}(\rho_{\theta} \parallel \mathsf{prior})$$

Minimizing these functionals is hard, SGD does it naturally

Potential Phi vs. original loss f

The solution of the variational problem is

$$\rho^{\rm ss}(x) = \frac{1}{Z_{\beta}} e^{-\beta \Phi(x)}$$

Key point

$$\rho^{\rm ss}(x) \neq \frac{1}{Z_{\beta}'} e^{-\beta f(x)}$$

Most likely locations of SGD are not the critical points of the original loss

The two losses are equal if and only if noise is isotropic

$$D(x) = I \quad \Leftrightarrow \quad \Phi(x) = f(x)$$

Deep networks have highly non-isotropic noise



Evaluate neural architectures using the diffusion matrix

How different are cats and dogs, really?





SGD converges to limit cycles

Theorem [Chaudhari & Soatto, ICLR '18]

The most likely trajectories of SGD are

 $\dot{x}=j(x),$

where the "leftover" vector field

$$j(x) = -\nabla f(x) + D(x) \nabla \Phi(x) - \beta^{-1} \operatorname{div} D(x)$$

is such that

$$\operatorname{div} j(x) = 0.$$

Trajectories of SGD

▶ Run SGD for 10⁵ epochs



An example



j(x) = 0

|j(x)| is small

very large |j(x)|

Theorem [Chaudhari & Soatto, ICLR '18]

The Ito SDE

$$dx = -\nabla f \ dt + \sqrt{2\beta^{-1}D} \ dW(t)$$

is equivalent to an A-type SDE

$$dx = -(D+Q) \nabla \Phi dt + \sqrt{2\beta^{-1}D} dW(t)$$

with the same steady-state $ho^{
m ss} \propto e^{-eta \Phi(x)}$ if

$$\nabla f = (D+Q) \nabla \Phi - \beta^{-1} \operatorname{div} (D+Q).$$

Knots in our understanding



Punchline

Is SGD special?

arXiv:1710.11029, ICLR '18

Stochastic gradient descent performs variational inference, converges to limit cycles for deep networks, Pratik Chaudhari and Stefano Soatto.



www.pratikac.info

Thank you, questions?