

Lecture #23: Nonconvex Optimization

Suppose we want to minimize a nonconvex function

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

but we only want to reach a local minimum

def. x is a local minimum, if for some $\epsilon > 0$,

$$\|x - y\| \leq \epsilon \Rightarrow f(y) \geq f(x)$$

Why would we be content with just a local minimum?

① Minimizing a nonconvex function is NP-hard, so what else can we do?

② the points gradient descent reaches in practice seem to be good enough

In fact the global minimum might not generalize well [Zhang et al]

Theorem [Murty, Trabadi] Deciding if a given point x is a local minimum of a nonconvex function is NP-hard

Why can't you just check ^{the gradient and} if the Hessian has negative eigenvalues?

The connection only goes in one direction:

claim: If x is a local minimum and f is twice differentiable, then $\nabla f(x) = 0$ and $\lambda_{\min}(\nabla^2 f(x)) \geq 0$

Proof: If either condition were violated, there would be a direction where the function decreases. \square

The difficulty in checking if x is a local minimum is:

If $\nabla f(x) = 0$, $\nabla^2 f(x) = 0$ then checking next terms in Taylor expansion \Leftrightarrow tensor optimization

But there is some hope:

Claim: If $\nabla f(x) = 0$ and $\lambda_{\min}(\nabla^2 f(x)) > 0$ then x is a local minimum

Q: Can we impose conditions on f to make finding a local minimum easy?

Aside from local minima/maxima we can also have saddle points

def: x is a saddle point if $\nabla f(x) = 0$ s.t. for any $\varepsilon > 0$ there are y and z with

$$\|x - y\|, \|x - z\| \leq \varepsilon \quad \text{and} \quad f(y) < f(x) < f(z)$$

E.g. $f(x_1, x_2) = x_1^2 - x_2^2$, then $x_1 = x_2 = 0$ is a saddle point

Claim: If $\nabla f(x) = 0$ and $\nabla^2 f(x)$ has positive and negative eigenvalues, x is a saddle point

In some cases, saddle points are inevitable

Imagine we have a function f where

- ① all local minima are global minima
- ② the set of all global minima are discrete and isolated (e.g. due to symmetries in parameterization)

Then we must have saddle points (if not, the function would be convex, which would contradict ②)

Q: Empirically, when gradient descent on a nonconvex function does not work well, is it because we got stuck in a bad local minimum or near a saddle point?

To tame these saddle points, let's assume

def. [Ge et al] A function f has the (θ, γ, δ) strict saddle property if for any x at least one of the following holds

- ① $\|\nabla f(x)\| \geq \theta$
- ② $\lambda_{\min}(\nabla^2 f(x)) \leq -\gamma$
- ③ x is δ -close to a local minimum

Under this condition, there are efficient algorithms for finding a point δ -close to a local minimum

Second order methods: (Trust Region) (Cubic Regularization)

Given a point x either

- ① $\nabla f(x)$ is large, can take a step $-\eta \nabla f(x)$
- ② $\lambda_{\min}(\nabla^2 f(x))$ is negative, can take a step $+\eta y$ where y is the corresponding eigenvector

However computing and storing the Hessian might be too expensive

First Order Methods: Noisy gradient descent

$$x_{t+1} = x_t - \eta \nabla f(x_t) + z_t$$

\uparrow
Gaussian noise

Theorem [Ge et al.] Noisy gradient descent can find a point δ -close to a local minimum in polynomial time

The number of iterations depends on $\theta, r, \delta, \beta, \rho, \frac{\epsilon}{\ln n}$

smoothness of Hessian
↓
smoothness target accuracy

[Lee et al.] use the stable manifold theorem from dynamical systems, but no bound on convergence

Intuition: Low probability over random starting point that trajectory passes too close to a saddle point

there are some cases we can prove no bad local minima:

[Auffinger, Ben-Arous, Cerny]: spherical p-spin glass

e.g. random third order tensor T yields distribution on S^{n-1} where T \uparrow Gaussian entries

$$\Pr[X] = \frac{e^{\sum_{ijk} x_i x_j x_k T_{ijk}}}{Z}$$

Z \uparrow partition function

Through Kac-Rice theorem can calculate distribution of critical points and show

"all local minima are close in objective value to the global minimum"

[Choromanska et al]: deep network \rightsquigarrow spherical p-spin glass

This "reduction" involves expanding each path into a separate random variable and treating them as independent

Note: In a spherical p-spin glass model the objective value of the global minimum is close to the objective value of any other point

If you spike it:

$$T + \lambda u \otimes u \otimes u$$

\uparrow random, Gaussian entries

\uparrow unit vector

then the optimization landscape is completely different

It now resembles [Frieze, Kannan]'s tensors or planted clique, which are believed to be hard to optimize

In another direction, for many matrix models we know characterizations like:

Matrix completion: with high probability the objective function

$$f(u, v) = \sum_{(i,j) \in \Omega} \|M_{ij} - (UV^T)_{ij}\|^2$$

has all local minima satisfy $UV^T = M$ and satisfies the strict saddle property

Also matrix sensing, robust PCA