

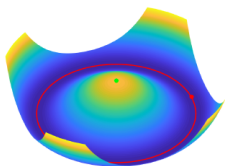
Computational Principles for High-dim Data Analysis

(Lecture Seventeen)

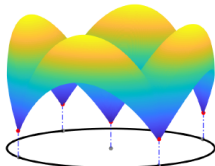
Yi Ma

EECS Department, UC Berkeley

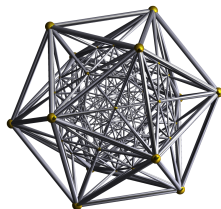
October 28, 2021



Rotational symmetry



Discrete symmetry



Nonconvex Optimization for High-Dim Problems

First Order Methods

- 1 Objectives of Nonconvex Optimization
- 2 Gradient Descent and Newton's Method
- 3 First Order Methods for Nonconvex Problems
 - Gradient and Negative Curvature Descent (Inexact)
 - Negative Curvature and Newton Descent (Inexact)
 - Gradient Descent with Small Random Noise
 - Hybrid Noisy Gradient Descent

"Premature optimization is the root of all evil."
– Donald Knuth, *The Art of Computer Programming*

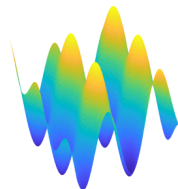
Nonconvex Optimization

Consider the problem of minimizing a general nonlinear function:

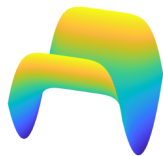
$$\min_{\mathbf{x}} f(\mathbf{x}), \quad \mathbf{x} \in \mathcal{C}. \quad (1)$$

In **the worst case**, even finding a *local* minimizer can be NP-hard¹.

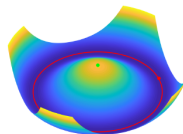
Nonconvex problems that arise from natural physical, geometrical, or statistical origins typically have **nice** structures, in terms of **symmetries!**



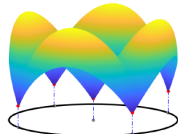
Spurious local minimizers



Flat saddle points



Rotational symmetry



Discrete symmetry

¹Some NP-complete problems in quadratic and nonlinear programming, K.G Murty and S. N. Kabadi, 1987

Objectives

Hence typically people seek to work with relatively benign (gradient Lipschitz continuous) functions:

$$\forall \mathbf{x}, \mathbf{y} \quad \|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2 \leq L_1 \|\mathbf{y} - \mathbf{x}\|_2 \quad (2)$$

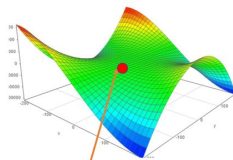
with benign objectives:

- ① convergence to some critical point \mathbf{x}_* such that: $\nabla f(\mathbf{x}_*) = \mathbf{0}$;
- ② the critical point \mathbf{x}_* is second-order stationary: $\nabla^2 f(\mathbf{x}_*) \succeq \mathbf{0}$.

Example: in general f could have irregular second-order stationary points:

Second Order Stationary Points

- $f(w) = \frac{1}{3}(w_1^3 - 3w_1w_2^2)$
- $\nabla f(w) = \begin{bmatrix} w_1^2 - w_2^2 \\ -2w_1w_2 \end{bmatrix}$
- $\nabla^2 f(w) = \begin{bmatrix} 2w_1 & -2w_2 \\ -2w_2 & -2w_1 \end{bmatrix}$
- $\nabla f(0) = 0, \nabla^2 f(0) = 0 \Rightarrow 0$ is SOSP
- $f([\epsilon, \epsilon]) = -\frac{2}{3}\epsilon^3 < f(0)$



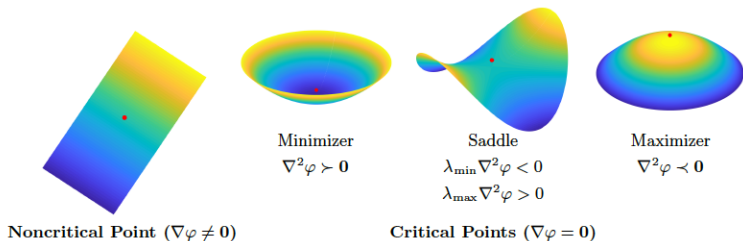
Second Order Stationary Point (SOSP)

Objectives

Hence typically people seek to work with relatively benign (gradient Lipschitz continuous) functions with benign objectives:

- ① convergence to some critical point \mathbf{x}_\star such that: $\nabla f(\mathbf{x}_\star) = \mathbf{0}$;
- ② the critical point \mathbf{x}_\star is second-order stationary: $\nabla^2 f(\mathbf{x}_\star) \succeq \mathbf{0}$.

Example: a function φ with symmetry only has **regular** critical points:



Gradient Descent (GD)

Function class: ∇f Lipschitz continuous with constant L_1 .

First-order oracle: the gradient $\nabla f(\mathbf{x})$ of the function $f(\mathbf{x})$.

The gradient descent iteration:

$$\mathbf{x}_k = \mathbf{x}_{k-1} - \gamma_k \nabla f(\mathbf{x}_{k-1}). \quad (3)$$

Proposition (Convergence Rate of GD for Nonconvex Functions)

Suppose that $f(\mathbf{x})$ is a (possibly nonconvex) differentiable function with ∇f Lipschitz continuous with constant L_1 . The gradient descent scheme with the step size $\gamma_k = 1/L_1$ converges to a critical point \mathbf{x}_ .*

Furthermore, for the gradient norm at the best iterate $\min_{0 \leq i \leq k-1} \|\nabla f(\mathbf{x}_i)\|_2 \leq \epsilon_g$, the number of iterations $k = O(\epsilon_g^{-2})$.

Newton's Method (strong convex)

Function class: f strongly convex and $\nabla^2 f$ Lipschitz continuous with L_2 .

The second-order oracle: the gradient $\nabla f(\mathbf{x})$ and the Hessian $\nabla^2 f(\mathbf{x})$.

The Newton iteration:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k). \quad (4)$$

Proposition (Convergence Rate of Newton's Method)

Let $f(\mathbf{x})$ be strongly convex, with $\lambda_{\min}(\nabla^2 f(\mathbf{x})) \geq \lambda > 0$ for all \mathbf{x} , and assume that $\nabla^2 f$ is Lipschitz continuous with constant L_2 , and let \mathbf{x}_\star be the (unique) minimizer of f over \mathbb{R}^n . Assuming $\|\mathbf{x}_0 - \mathbf{x}_\star\|_2 < \frac{2\lambda}{L_2}$, the iterates \mathbf{x}_k converge to \mathbf{x}_\star , with quadratic rate.

Unfortunately, for high-dim problems, impossible to compute $\nabla^2 f$.

Cubic Regularized Newton's Method (nonconvex)

Function class: f nonconvex and $\nabla^2 f$ Lipschitz continuous with L_2 .

The second-order oracle: the gradient $\nabla f(\mathbf{x})$ and the Hessian $\nabla^2 f(\mathbf{x})$.
Consider the local cubic surrogate:

$$\hat{f}(\mathbf{y}, \mathbf{x}) \doteq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2} (\mathbf{y} - \mathbf{x})^* \nabla^2 f(\mathbf{x}) (\mathbf{y} - \mathbf{x}) + \frac{L_2}{6} \|\mathbf{y} - \mathbf{x}\|_2^3. \quad (5)$$

The cubic Newton iteration:

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{y}} \hat{f}(\mathbf{y}, \mathbf{x}_k). \quad (6)$$

Theorem (Convergence Rate of Cubic Newton's Method)

Suppose $f(\mathbf{x})$ is bounded from below. Then the sequence $\{\mathbf{x}_k\}$ generated by the cubic regularized Newton step (6) converges to a non-empty set of limit points \mathbf{X}_ of SOS points. For $\|\nabla f(\mathbf{x}_k)\|_2 \leq \epsilon_g$, the number of iterations $k = O(\epsilon_g^{-3/2})$.*

Unfortunately, for high-dim problems, impossible to compute $\nabla_{\equiv}^2 f$.

Gradient and Negative Curvature Descent

Function class: f nonconvex and $\nabla^2 f$ Lipschitz continuous with L_2 .

The oracle: gradient $\nabla f(\mathbf{x})$ and a negative eigenvector \mathbf{e} of $\nabla^2 f(\mathbf{x})$.

Hybrid gradient and negative curvature descent:

- **if** $\|\nabla f(\mathbf{x}_k)\|_2 \geq \epsilon_g = (2L_1\epsilon)^{1/2}$, **then** $\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{1}{L_1}\nabla f(\mathbf{x}_k)$;
- **else if** $-\lambda_k(\nabla^2 f(\mathbf{x})) \geq \epsilon_H = (1.5L_2^2\epsilon)^{1/3}$, **then** $\mathbf{x}_{k+1} = \mathbf{x}_k + \frac{2\lambda_k}{L_2}\mathbf{e}_k$.

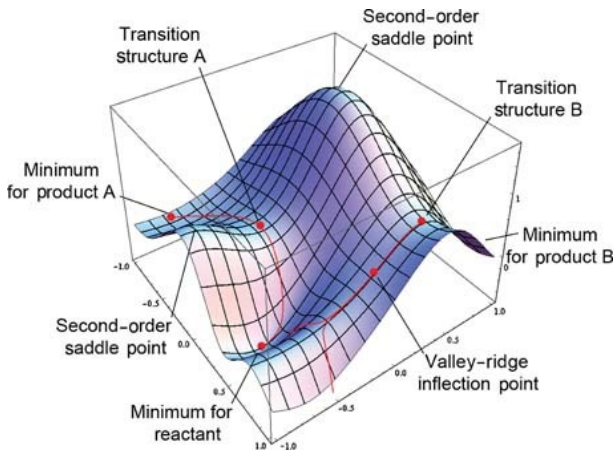
Theorem (Convergence of Gradient and Negative Curvature Descent)

The above hybrid gradient and negative curvature descent scheme converges to a second-order stationary point \mathbf{x}_ with the desired precision in function value ϵ in no more than $k = (f(\mathbf{x}_0) - f(\mathbf{x}_*))/\epsilon$ iterations. Or in terms of $\|\nabla f(\mathbf{x}_k)\|_2 \leq \epsilon_g$, $k = O(\epsilon_g^{-2})$.*

The same convergence rate as GD, but converges to an SOS point!

Gradient and Negative Curvature Descent

An Example: potential energy surface in Chemistry.



Compute Negative Curvature: the Power Iteration

Hessian: $\mathbf{H} \doteq \nabla^2 f(\mathbf{x})$. Want to compute negative curvature direction \mathbf{e} :

$$\mathbf{H}\mathbf{e} = \lambda_{\min}(\mathbf{H})\mathbf{e} \quad \text{or} \quad \mathbf{A}\mathbf{e} = \lambda_{\max}(\mathbf{A})\mathbf{e}, \quad \text{with } \mathbf{A} \doteq \mathbf{I} - L_1^{-1}\mathbf{H} \succ \mathbf{0}.$$

Power iteration:

$$\hat{\lambda}_{k+1} = \frac{\langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle}, \quad \mathbf{x} = \mathbf{A}^k \mathbf{b},$$

where $\mathbf{A}^i \mathbf{b}$ can be approximated for a small $t > 0$ with:

$$\mathbf{A}\mathbf{b} = [\mathbf{I} - L_1^{-1}\mathbf{H}] \mathbf{b} \approx \mathbf{b} - (tL_1)^{-1}(\nabla f(\mathbf{x} + t\mathbf{b}) - \nabla f(\mathbf{x})).$$

Two gradient evaluations per iteration.

Compute Negative Curvature: the Lanczos Method

The Krylov information: $\mathbf{K} \doteq [\mathbf{b}, \mathbf{A}\mathbf{b}, \mathbf{A}^2\mathbf{b}, \dots, \mathbf{A}^k\mathbf{b}]$.

The Lanczos method:

$$\hat{\lambda}_{k+1} = \max_{\mathbf{x}} \frac{\langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle}, \quad \mathbf{x} \in \text{span}(\mathbf{K}).$$

Proposition (Convergence Rate of Lanczos)

Use the Lanczos procedure to find the largest eigenvalue of $\mathbf{I} - L_1^{-1}\mathbf{H}$ starting from a random unit vector. Then, for any $\epsilon_\lambda > 0$ and $\delta \in (0, 1)$, with a probability at least $1 - \delta$ the procedure outputs a unit vector \mathbf{e}' such that $(\mathbf{e}')^ \mathbf{H} \mathbf{e}' \leq \lambda_{\min}(\mathbf{H}) + \epsilon_\lambda$ in at most number of iterations:*

$$\min \left\{ n, \frac{\log(n/\delta^2)}{2\sqrt{2}} \sqrt{\frac{L_1}{\epsilon_\lambda}} \right\}.$$

In terms of the first-order oracle, complexity of the inexact gradient and negative curvature descent is $k \leq O(\epsilon_g^{-2})$.

Negative Curvature and Newton Descent

Quadratic regularized Newton:

$$\mathbf{s}_k = \arg \min_{\mathbf{s}} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{s} \rangle + \frac{1}{2} \mathbf{s}^* \nabla^2 f(\mathbf{x}_k) \mathbf{s} + \frac{\lambda}{2} \|\mathbf{s}\|_2^2 \quad (7)$$

$$= -[\nabla^2 f(\mathbf{x}_k) + \lambda \mathbf{I}]^{-1} \nabla f(\mathbf{x}_k). \quad (8)$$

The Levenberg-Marquardt iteration:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - [\nabla^2 f(\mathbf{x}_k) + \lambda \mathbf{I}]^{-1} \nabla f(\mathbf{x}_k). \quad (9)$$

LM is popular for solving nonlinear least squares problems.

Negative Curvature and Newton Descent

Function class: f nonconvex and $\nabla^2 f$ Lipschitz continuous with L_2 .

The oracle: gradient $\nabla f(\mathbf{x})$ and $\nabla^2 f(\mathbf{x})$ (to be approximated).

Hybrid gradient and negative curvature descent:

- **if** $-\lambda_k(\nabla^2 f(\mathbf{x})) \geq \epsilon_H = (3L_2^2\epsilon)^{1/3}$, **then** $\mathbf{x}_{k+1} = \mathbf{x}_k + \frac{2\lambda_k}{L_2} \mathbf{e}_k$;
- **else if** $\|\nabla f(\mathbf{x}_k)\|_2 \geq \epsilon_g = 3^{8/3}L_2^{1/3}\epsilon^{2/3}/2$, **then** $\mathbf{x}_{k+1} = \mathbf{x}_k + \gamma_k \mathbf{s}_k$.

Theorem

Assume $\{\mathbf{x}_k\}$ are generated by the hybrid negative curvature and Newton descent. Then in at most

$$k \leq \frac{f(\mathbf{x}_0) - f(\mathbf{x}_*)}{\epsilon} \quad (10)$$

iterations, \mathbf{x}_k will be an approximate second-order stationary point such that $\|\nabla f(\mathbf{x}_k)\|_2 \leq \epsilon_g$, $\lambda_{\min}(\nabla^2 f(\mathbf{x}_k)) \geq -\epsilon_H$.

Conjugate Gradient Descent

Need to compute e_k and s_k without knowing $\nabla^2 f(\mathbf{x})$. Notice that, similar to e_k , to find s_k we need solve: $\underbrace{[\nabla^2 f(\mathbf{x}_k) + \lambda \mathbf{I}]}_{\mathbf{A}} \mathbf{s}_k = \underbrace{-\nabla f(\mathbf{x}_k)}_{\mathbf{y}}$.

A special case of the quadratic minimization problem: $\min_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$.

Conjugate gradient descent:² Initialize the residual \mathbf{r}_i and descent direction \mathbf{d}_i as: $\mathbf{d}_0 = \mathbf{r}_0 = \mathbf{y} - \mathbf{A}\mathbf{x}_0$. Then:

$$\text{Conjugate Gradient: } \begin{cases} \alpha_i &= \frac{\mathbf{r}_i^* \mathbf{r}_i}{\mathbf{d}_i^* \mathbf{A} \mathbf{d}_i}, \\ \mathbf{x}_{i+1} &= \mathbf{x}_i + \alpha_i \mathbf{d}_i, \\ \mathbf{r}_{i+1} &= \mathbf{r}_i - \alpha_i \mathbf{A} \mathbf{d}_i, \\ \beta_{i+1} &= \frac{\mathbf{r}_{i+1}^* \mathbf{r}_{i+1}}{\mathbf{r}_i^* \mathbf{r}_i}, \\ \mathbf{d}_{i+1} &= \mathbf{r}_{i+1} + \beta_{i+1} \mathbf{d}_i. \end{cases} \quad i = 0, 1, 2, \dots \quad (11)$$

²An introduction to the conjugate gradient method without the agonizing pain, Jonathan Shewchuk, Technical report, Carnegie Mellon University, 1994.

Negative Curvature and Newton Descent: Complexity

Theorem (Complexity of Approximate Conjugate Gradient)

To solve $\mathbf{A}\mathbf{s} = \mathbf{y}$ with $\epsilon_H \mathbf{I} \preceq \mathbf{A} \preceq (L_1 + 2\epsilon_H)\mathbf{I}$, the conjugate gradient method computes an \mathbf{s}' that satisfies

$$\|(\nabla^2 f(\mathbf{x}_k) + 2\epsilon_H)\mathbf{s}_k + \nabla f(\mathbf{x}_k)\|_2 \leq \frac{1}{2}\epsilon_H \|\mathbf{s}_k\|_2$$

in at most $O(\epsilon_H^{-1/2} \log(\frac{1}{\epsilon_H}))$ iterations.

With the first-order oracle, complexity of the inexact negative curvature and newton descent achieves **the best known rate**: $k \leq O(\epsilon_g^{-7/4})$.

Gradient Descent with Small Random Noise

Function class: f nonconvex and $\nabla^2 f$ Lipschitz continuous with L_2 .

The oracle: gradient $\nabla f(\mathbf{x})$ and small random noise.

The **Langevin dynamics** with noisy gradient flow:

$$\dot{\mathbf{x}}(t) = -\frac{1}{2}\nabla f(\mathbf{x}(t)) + \sqrt{\lambda}\mathbf{n}(t), \quad (12)$$

Probability density of \mathbf{x} converges to the **Gibbs measure**:

$$p^\lambda(\mathbf{x}) = C^\lambda \exp\left(-\frac{1}{\lambda}f(\mathbf{x})\right). \quad (13)$$

Lemma (Laplace's Method: Scalar Case)

Suppose $f(x)$ is a twice continuously differentiable function with a unique maximizer x_0 and $f''(x_0) < 0$. Then we have

$$\lim_{\lambda \rightarrow 0} \int e^{\frac{1}{\lambda}f(x)} dx = e^{\frac{1}{\lambda}f(x_0)} \sqrt{\frac{2\pi\lambda}{-f''(x_0)}} \propto \int e^{\frac{1}{\lambda}f(x)} \delta(x - x_0) dx. \quad (14)$$

The Laplace Method

Theorem (Laplace Method: Multivariate and Multiple Global Minimizers)

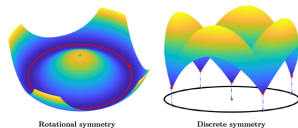
Let $f(\mathbf{x})$ be a function with at least quadratic growth as $\mathbf{x} \rightarrow \infty$. Suppose $f(\mathbf{x})$ has multiple global non-degenerate minimizers at $\mathbf{x}_*^1, \dots, \mathbf{x}_*^N$ and they are all non-degenerate. Then in the limit $\lambda \downarrow 0$, the density $p^\lambda(\mathbf{x})$ converges to

$$p^0(\mathbf{x}) = \frac{\sum_{i=1}^N a_i \delta(\mathbf{x} - \mathbf{x}_*^i)}{\sum_{i=1}^N a_i}, \quad \text{with} \quad a_i = \det[\mathbf{H}(\mathbf{x}_*^i)]^{-1/2}, \quad (15)$$

where $\mathbf{H}(\mathbf{x}) = \nabla^2 f(\mathbf{x})$ is the Hessian of the function $f(\mathbf{x})$.

When all global minimizers make a continuous submanifold \mathcal{M} , $p^\lambda(\mathbf{x})$ converges to a density on \mathcal{M} given by:

$$p^0(\mathbf{x}) = \frac{\det[\mathbf{H}(\mathbf{x})]^{-1/2}}{\int_{\mathcal{M}} \det[\mathbf{H}(\mathbf{y})]^{-1/2} d\mathbf{y}}, \quad \mathbf{x} \in \mathcal{M}.$$



Noisy Gradient with Langevin Monte Carlo

Function class: ∇f Lipschitz continuous with constant L_1 .

First-order oracle: the gradient $\nabla f(\mathbf{x})$ and small noise \mathbf{n} .

Langevin Monte Carlo:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{1}{L_1} \nabla f(\mathbf{x}_k) + \sqrt{2\lambda/L_1} \mathbf{n}_k. \quad (16)$$

Proposition (Noisy Gradient Descent)

Considering the above noisy gradient descent scheme (16), if $\|\nabla f(\mathbf{x}_k)\|_2 \geq (2L_1\epsilon)^{1/2}$, then we have

$$\mathbb{E}[f(\mathbf{x}_{k+1}) \mid \mathbf{x}_k] \leq f(\mathbf{x}_k) - \epsilon + \lambda. \quad (17)$$

Descent when $\|\nabla f(\mathbf{x}_k)\|_2 > (2L_1\lambda)^{1/2}$; explore stability otherwise.

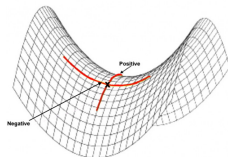
Effect of Noisy Gradient around a Saddle Point

Consider a standard quadratic function:

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^* \mathbf{H} \mathbf{x} \text{ for a constant } \mathbf{H} \in \mathbb{R}^{n \times n},$$

with the smallest eigenvalue $\lambda_{\min} < 0$,

and the Lipschitz constant $L_1 = \max_i |\lambda_i(\mathbf{H})|$.



The Langevin dynamics becomes:

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - \frac{1}{L_1} \nabla f(\mathbf{x}_k) + \sqrt{2\lambda/L_1} \mathbf{n}_k \\ &= \underbrace{(\mathbf{I} - L_1^{-1} \mathbf{H})}_{\mathbf{A}} \mathbf{x}_k + \underbrace{\sqrt{2\lambda/L_1} \mathbf{n}_k}_{\mathbf{b}}. \end{aligned} \quad (18)$$

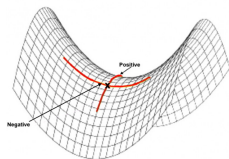
This is **an unstable linear dynamic system** with random noise as the input:

$$\mathbf{x}_{k+1} = \mathbf{A} \mathbf{x}_k + \mathbf{b} \mathbf{n}_k. \quad (19)$$

Escaping Saddle Point

Therefore, the accumulated dynamics:

$$\mathbf{x}_{k+1} = \mathbf{A}^{k+1} \mathbf{x}_0 + b \sum_{i=0}^k \mathbf{A}^{k-i} \mathbf{n}_i. \quad (20)$$



$\mathbf{A}^{k+1} \mathbf{x}_0$ and $\mathbf{A}^{k-i} \mathbf{n}_i$ are **powers** of the matrix \mathbf{A} applied to random vectors (assuming \mathbf{x}_0 random too).

Question: which direction survives in power iteration?

Proposition (Escaping Saddle Point via Noisy Gradient Descent)

Consider the noisy gradient descent via the Langevin dynamics (18) for the function $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^* \mathbf{H} \mathbf{x}$, starting from $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$. Then after $k \geq \frac{\log n - \log(|\lambda_{\min}|/L_1)}{2 \log(1 + |\lambda_{\min}|/L_1)}$ steps, we have

$$\mathbb{E}[f(\mathbf{x}_{k+1}) - f(\mathbf{x}_0)] \leq -\lambda. \quad (21)$$

Hybrid Noisy Gradient Descent

Function class: f nonconvex and $\nabla^2 f$ Lipschitz continuous with L_2 .

The oracle: gradient $\nabla f(\mathbf{x})$ and small noise \mathbf{n} .

Hybrid noisy gradient descent:

- if $\|\nabla f(\mathbf{x}_k)\|_2 \geq \epsilon_g$, then $\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{1}{L_1} \nabla f(\mathbf{x}_k)$;
- else $\mathbf{x}_k^0 = \mathbf{x}_k$, and negative curvature descent with noisy gradients:
for $i = 0, 1, 2, \dots, k_{\max} = O(\log n)$

$$\mathbf{x}_k^{i+1} = \mathbf{x}_k^i - \frac{1}{L_1} \nabla f(\mathbf{x}_k^i) + \sqrt{2\epsilon/L_1} \mathbf{n}^i,$$

where $\mathbf{n}^i \sim \mathcal{N}(0, \mathbf{I})$.

Complexity: To guarantee $\|\nabla f(\mathbf{x})\| \leq \epsilon_g$, the number of total gradient evaluations needed is $O(\epsilon_g^{-2})$, up to a $\log(n)$ factor.³

³Perturbed **accelerated** gradient descent reduces to $O(\epsilon_g^{-7/4})$.

Summary

Table: Oracles and complexities (up to log factors) of different optimization methods. “Stat. Point” stands for the type of stationary point \mathbf{x}_* to which the method guarantees to converge. Complexity is measured in terms of the number of oracles accessed before attaining a prescribed accuracy $\|\nabla f(\mathbf{x}_*)\| \leq \epsilon_g$.

Method	Oracle	Stat. Point	Complexity
Vanilla gradient descent	first-order	first-order	$O(\epsilon_g^{-2})$
Cubic Regularized Newton	second-order	second-order	$O(\epsilon_g^{-1.5})$
Gradient/negative curvature	first-order	second-order	$O(\epsilon_g^{-2})$
Negative curvature/Newton	first-order	second-order	$O(\epsilon_g^{-1.75})$
Hybrid noisy gradient	first-order	second-order	$O(\epsilon_g^{-2})$
Perturbed accelerated gradient	first-order	second-order	$O(\epsilon_g^{-1.75})$

Assignments

- Reading: Section 9.1 - 9.5 of Chapter 9.
- Programming Homework #3.