



Deep Learning

Lecture 4: Non-convex Optimization for Learning (Part 1)

University of Agder,
Kristiansand, Norway

Prepared by: Hadi Ghauch^{*}, Hossein S. Ghadikolaei[†]

^{*} Telecom ParisTech

[†] Royal Institute of Technology, KTH

<https://sites.google.com/view/fundl/home>

Outline

1. Motivation
2. Optimization of problems without structure
 - First-order methods
 - Successive Approximation
3. Optimization of problems with structure

Outline

1. Motivation
2. Optimization of problems without structure
 - First-order methods
 - Successive Approximation
3. Optimization of problems with structure

Recap of Convex Optimization Algorithm

Strongly convex optimization problem: $\min_{\mathbf{w}} \frac{1}{N} \sum_{i \in [N]} f_i(\mathbf{w}) + r(\mathbf{w})$

Existence of global optimal and efficient solution methods

Deterministic algorithms (Lec 2)

Gradient descent: move along steepest descent (complexity **linear** in N)

Newton's methods: include Hessian (curvature) information in the update.

Not used in learning (complexity **cubic** in d)

Nesterov Acceleration: mix between interpolation and gradient descent.

Need to know, L (complexity **linear** in N)

Momentum Methods (Heavy ball method): add momentum term to reduce the oscillations and improve convergence. Need to know, L, μ (complexity **linear** in N)

Stochastic Gradient Descent

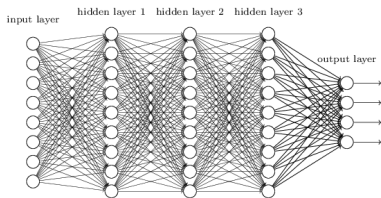
- computational complexity **independent of** N

- SGD family for smooth problems

- noise reduction techniques and adaptive mini-batches, SVRG, and SAGA

Motivation

Resurgence of AI is due to **Deep Neural Networks (DNNs)**
and countless variants (covered later)



DNNs is a composition of **non-linear layers**, $\mathbf{W}_1, \dots, \mathbf{W}_J$:

$$\mathbf{y}_i = \sigma_J(\mathbf{W}_J \cdots \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x}_i))), \quad i \in [N]$$

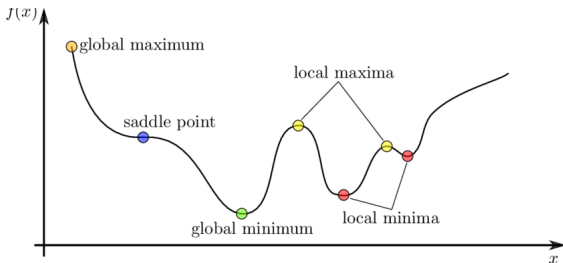
Let $\mathbf{w} \in \mathbb{R}^d$ be the total number of weights in DNN (from all layers)
 $d \geq 10^6$ in current deep learning application

The resulting optimization DNN training

$$\min_{\mathbf{w} \in \mathcal{D}} \frac{1}{N} \sum_{i=1}^N f_i((\mathbf{x}_i, \mathbf{y}_i); \mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 := f(\mathbf{w})$$

non-convex optimization problem

Definitions



w^* **global minimum** iff $\nabla f(w^*) = 0$, $\nabla^2 f(w^*) \succeq 0$, $f(w^*) \leq f(w)$, $\forall w \in \mathcal{D}$

w^* **local minimum** iff $\nabla f(w^*) = 0$, $\nabla^2 f(w^*) \succeq 0$, $f(w^*) \leq f(w)$, $\forall w \in \mathcal{B}$

w^* **non-degenerate saddle points** iff $\nabla f(w^*) = 0$, $\nabla^2 f(w^*)$ has strictly positive and negative eigenvalues

w^* **saddle point** iff $\nabla f(w^*) = 0$, $\nabla^2 f(w^*)$ has (positive and negative eigenvalues)

w^* **stationary point** iff $\nabla f(w^*) = 0$, $w^* \in \mathcal{D}$

decrease in level of restrictiveness

Definitions

Problem: $\arg \min_{\mathbf{w} \in \mathcal{D}} f(\mathbf{w})$

Convex world: any point $\nabla f(\mathbf{w}_k) = 0, \mathbf{w}_k \in \mathcal{D}$ globally optimal point

Nonconvex world: any point $\nabla f(\mathbf{w}_k) = 0, \mathbf{w}_k \in \mathcal{D}$ is stationary point:
global min ? local min ? saddle point ?

global solution hard for general non-convex prob

Goal: bypass stationary points (usually bad) and find local minimum

Second-order necessary (2oN) point: $\|\nabla f(\mathbf{w}_k)\|_2^2 = 0$ & $\nabla^2 f(\mathbf{w}_k) \succeq \mathbf{0}$
- point with zero gradient and positive curvature: local minimum

Approximate 2oN point: $\|\nabla f(\mathbf{w}_k)\|_2^2 \leq \epsilon_g$, $\nabla^2 f(\mathbf{w}_k) \succeq -\epsilon_H \mathbf{I}$ for small positive ϵ_g, ϵ_H

- approximation of local minimum (relaxation of conditions of 2oN pt)

Complexity measure:

gradient evaluations: # calls to incremental first-order oracle with input (\mathbf{w}, i) and output $(f_i(\mathbf{w}), \nabla f_i(\mathbf{w}))$

Basic assumptions

f is bounded below: $f(\mathbf{w}) \geq f_{\inf}$ for all $\mathbf{w} \in \mathcal{D}$

Function, its gradient and Hessian are Lipschitz continuous

$$\begin{aligned}\|f(\mathbf{w}_2) - f(\mathbf{w}_1)\|_2 &\leq L\|\mathbf{w}_2 - \mathbf{w}_1\|_2 \\ \|\nabla f(\mathbf{w}_2) - \nabla f(\mathbf{w}_1)\|_2 &\leq L_g\|\mathbf{w}_2 - \mathbf{w}_1\|_2 \\ \|\nabla^2 f(\mathbf{w}_2) - \nabla^2 f(\mathbf{w}_1)\|_F &\leq L_H\|\mathbf{w}_2 - \mathbf{w}_1\|_2\end{aligned}$$

L , L_g and L_H are Lipschitz const of the function, the gradient and Hessian

Quadratic upper-bounds from Taylor's theorem ($\forall \mathbf{v}, \mathbf{w} \in \mathcal{D}$)

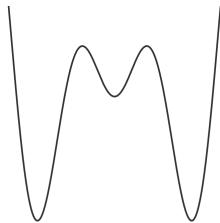
- approximate func at $f(\mathbf{w} + \mathbf{v})$ with Talor, using at grad and Hessian at \mathbf{w} :

$$f(\mathbf{w} + \mathbf{v}) \leq f(\mathbf{w}) + \nabla f(\mathbf{w})^T \mathbf{v} + \frac{1}{2} \mathbf{v}^T \nabla^2 f(\mathbf{w}) \mathbf{v}$$

Non-convex Optimization

$$\min_{\mathbf{w} \in \mathcal{D}} f(\mathbf{w}) \quad (1)$$

$f(\mathbf{w}) : \mathcal{D} \subseteq \mathbb{R}^d \rightarrow \mathbb{R}$ is non-convex
 \mathcal{D} : domain of f (convex or non-convex)



Local optimality **may not** necessarily imply global optimality

Proper initialization is very important in nonconvex optimization

First-order criteria ($\nabla f(\mathbf{w}) = 0$)

necessary and sufficient conditions for convex

only necessary condition for nonconvex

Strategies for Non-convex Optimization

No generic method for all non-convex problems (\neq convex case)
solution approach depends on **structure of the optimization**

1. **No structure on $f(w)$:**
first-order methods may be used. **convergence?**
successive approximation: successively approx (1) with linear/convex bound
2. **$f(w)$ is coordinate separable** into convex subproblems (**k-means, Lyods algo**)
if $f(w) = f(w_1, \dots, w_d)$, where (w_1, \dots, w_d) **coordinates**
solve (1) as series **subproblems**: alternatively min each **subproblems** while fixing all other ones. Strong convexity of subproblems to prove convergence
3. **$f(w)$ is block-coordinate separable** into convex subproblems:
if $f(w) = f(w_1, \dots, w_d)$, where (w_1, \dots, w_d) **block of coordinates**
solve (1) as series **subproblems**: alternatively min each **subproblems** while fixing all other ones. Strong convexity of subproblems to prove convergence
4. **$f(w)$ is block-coordinate separable** into non-convex subproblems:
solve (1) as series **subproblems**: alternatively min each **subproblems** while fixing all other ones. Strong convexity of subproblems not needed to prove convergence

Outline

1. Motivation
2. Optimization of problems without structure
 - First-order methods
 - Successive Approximation
3. Optimization of problems with structure

Can we use a simple GD?

One-point convexity:

Convexity implies $f(w_k) - f(v) \leq \nabla f(w_k)^T (w_k - v)$ for all v

Then it must hold for a local min: $v = w^*$?

What if w^* is global min? GD converges to global min

Run a GD step: $w_{k+1} = w_k - \frac{\nabla f(w_k)}{L_g}$.

Use def of convexity show $\Rightarrow \Delta_k := f(w_k) - f(w^*) \leq \nabla f(w_k)^T (w_k - w^*)$
 f convex only in neighborhood of w^* : $\Delta_k \leq \mathcal{O}(1/k)$

Linear convergence to local min.

computational complexity: $\mathcal{O}(N)$ gradient evaluations per iteration

assuming one point convexity: func is locally convex around some local min

This condition holds for many functions

Challenge for GD in non-convex setting?

poor scalability with N

convergence to a stationary point, not necessarily a local minimum

results conditioned on showing f satisfies one point convexity condition

Can we use a simple GD?

One-point convexity:

Convexity implies $f(\mathbf{w}_k) - f(\mathbf{v}) \leq \nabla f(\mathbf{w}_k)^T (\mathbf{w}_k - \mathbf{v})$ for all \mathbf{v}

Then it must hold for a local min: $\mathbf{v} = \mathbf{w}^*$?

What if \mathbf{w}^* is global min? GD converges to global min

Run a GD step: $\mathbf{w}_{k+1} = \mathbf{w}_k - \frac{\nabla f(\mathbf{w}_k)}{L_g}$.

Use def of convexity show $\Rightarrow \Delta_k := f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \nabla f(\mathbf{w}_k)^T (\mathbf{w}_k - \mathbf{w}^*)$
 f convex only in neighborhood of \mathbf{w}^* : $\Delta_k \leq \mathcal{O}(1/k)$

Linear convergence to local min.

computational complexity: $\mathcal{O}(N)$ gradient evaluations per iteration

assuming one point convexity: func is locally convex around some local min

This condition holds for many functions

Challenge for GD in non-convex setting?

- poor scalability with N

- convergence to a stationary point, not necessarily a local minimum

results conditioned on showing f satisfies one point convexity condition

Can we use a simple GD?

One-point convexity:

Convexity implies $f(\mathbf{w}_k) - f(\mathbf{v}) \leq \nabla f(\mathbf{w}_k)^T (\mathbf{w}_k - \mathbf{v})$ for all \mathbf{v}

Then it must hold for a local min: $\mathbf{v} = \mathbf{w}^*$?

What if \mathbf{w}^* is global min? GD converges to global min

Run a GD step: $\mathbf{w}_{k+1} = \mathbf{w}_k - \frac{\nabla f(\mathbf{w}_k)}{L_g}$.

Use def of convexity show $\Rightarrow \Delta_k := f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \nabla f(\mathbf{w}_k)^T (\mathbf{w}_k - \mathbf{w}^*)$
 f convex only in neighborhood of \mathbf{w}^* : $\Delta_k \leq \mathcal{O}(1/k)$

Linear convergence to local min.

computational complexity: $\mathcal{O}(N)$ gradient evaluations per iteration

assuming one point convexity: func is locally convex around some local min

This condition holds for many functions

Challenge for GD in non-convex setting?

- poor scalability with N

- convergence to a stationary point, not necessarily a local minimum

results conditioned on showing f satisfies one point convexity condition

Can we use a simple GD?

One-point convexity:

Convexity implies $f(\mathbf{w}_k) - f(\mathbf{v}) \leq \nabla f(\mathbf{w}_k)^T (\mathbf{w}_k - \mathbf{v})$ for all \mathbf{v}

Then it must hold for a local min: $\mathbf{v} = \mathbf{w}^*$?

What if \mathbf{w}^* is global min? GD converges to global min

Run a GD step: $\mathbf{w}_{k+1} = \mathbf{w}_k - \frac{\nabla f(\mathbf{w}_k)}{L_g}$.

Use def of convexity show $\Rightarrow \Delta_k := f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \nabla f(\mathbf{w}_k)^T (\mathbf{w}_k - \mathbf{w}^*)$
 f convex only in neighborhood of \mathbf{w}^* : $\Delta_k \leq \mathcal{O}(1/k)$

Linear convergence to local min.

computational complexity: $\mathcal{O}(N)$ gradient evaluations per iteration

assuming one point convexity: func is locally convex around some local min

This condition holds for many functions

Challenge for GD in non-convex setting?

- poor scalability with N

- convergence to a stationary point, not necessarily a local minimum

- results conditioned on showing f satisfies one point convexity condition

One-point convexity [Allen-Zhu, ICML, 2017]

Convex f

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{v} \rangle \geq f(\mathbf{w}) - f(\mathbf{v})$$

μ -strongly convex

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \mu \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq 2\mu (f(\mathbf{w}) - f(\mathbf{w}^*))$$

μ -strongly convex and L_g -smooth

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \frac{1}{2L_g} \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

If f **nonconvex** but satisfies

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

If f **nonconvex** but satisfies

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

two-layer NN [Li-Yuan, 2017]

$$\|\nabla f(\mathbf{w})\|_2^2 \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

(known as Polyak-Lojasiewicz condition)

finite sum minimization [Reddi-Sra-Poczos-Smola, 2016]

If f **nonconvex** but satisfies

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\beta}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \gamma \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

dictionary learning [Arora-Ge-Ma-Moitra, 2015]

One-point convexity [Allen-Zhu, ICML, 2017]

Convex f

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{v} \rangle \geq f(\mathbf{w}) - f(\mathbf{v})$$

μ -strongly convex

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \mu \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq 2\mu (f(\mathbf{w}) - f(\mathbf{w}^*))$$

μ -strongly convex and L_g -smooth

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \frac{1}{2L_g} \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

If f **nonconvex** but satisfies

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

If f **nonconvex** but satisfies

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

two-layer NN [Li-Yuan, 2017]

$$\|\nabla f(\mathbf{w})\|_2^2 \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

(known as Polyak-Lojasiewicz condition)

finite sum minimization [Reddi-Sra-Poczos-Smola, 2016]

If f **nonconvex** but satisfies

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\beta}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \gamma \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

dictionary learning [Arora-Ge-Ma-Moitra, 2015]

One-point convexity, how to solve?

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{v} \rangle \geq f(\mathbf{w}) - f(\mathbf{v})$$

GD/SGD converges in $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

GD converges in $\mathcal{O}\left(\frac{1}{\epsilon}\right)$ for smooth f

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \mu \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

GD/SGD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq 2\mu (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$ for smooth f

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \frac{1}{2L_g} \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

GD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD/SGD converges in $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

GD converges in $\mathcal{O}\left(\frac{1}{\epsilon}\right)$ for smooth f

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

GD/SGD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$ for smooth f

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\beta}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \gamma \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

GD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

One-point convexity, how to solve?

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{v} \rangle \geq f(\mathbf{w}) - f(\mathbf{v})$$

GD/SGD converges in $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

GD converges in $\mathcal{O}\left(\frac{1}{\epsilon}\right)$ for smooth f

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \mu \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

GD/SGD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq 2\mu (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$ for smooth f

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \frac{1}{2L_g} \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

GD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD/SGD converges in $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

GD converges in $\mathcal{O}\left(\frac{1}{\epsilon}\right)$ for smooth f

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

GD/SGD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$ for smooth f

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\beta}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \gamma \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

GD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

One-point convexity, how to solve?

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{v} \rangle \geq f(\mathbf{w}) - f(\mathbf{v})$$

GD/SGD converges in $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

GD converges in $\mathcal{O}\left(\frac{1}{\epsilon}\right)$ for smooth f

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \mu \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

GD/SGD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq 2\mu (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$ for smooth f

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \frac{1}{2L_g} \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

GD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD/SGD converges in $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

GD converges in $\mathcal{O}\left(\frac{1}{\epsilon}\right)$ for smooth f

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

GD/SGD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$ for smooth f

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\beta}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \gamma \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

GD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

One-point convexity, how to solve?

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{v} \rangle \geq f(\mathbf{w}) - f(\mathbf{v})$$

GD/SGD converges in $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

GD converges in $\mathcal{O}\left(\frac{1}{\epsilon}\right)$ for smooth f

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \mu \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

GD/SGD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq 2\mu (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$ for smooth f

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \frac{1}{2L_g} \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

GD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD/SGD converges in $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

GD converges in $\mathcal{O}\left(\frac{1}{\epsilon}\right)$ for smooth f

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

GD/SGD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

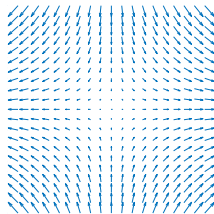
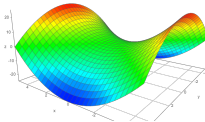
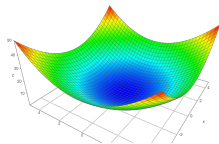
$$\|\nabla f(\mathbf{w})\|_2^2 \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$ for smooth f

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\beta}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \gamma \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

GD converges in $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

How to escape non-degenerate saddle points



Find an approximate 2oN point: $\|f(\mathbf{w}_k)\|_2^2 \leq \epsilon_g, \nabla^2 f(\mathbf{w}_k) \succeq -\epsilon_H \mathbf{I}$
- approximate local min (zero gradient with positive curvature)

Gradient-based

- 1) Perturbed GD: $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k) + \text{noise}$
 $\log(d)$ complexity on the parameter size [Jin-Ge-Netrapalli-Kakade-Jordan, 2017]
- 2) SGD

Newton's method: too expensive

Hessian-based

Faster reaction to saddle points using curvature information, expensive iterations could be very efficient [Allen-Zhu, 2018]

A generic Hessian-based algorithm

1. If $\|\nabla f(\mathbf{w}_k)\|_2 > \epsilon_g$, run your favorite algorithm (say GD with step-size $1/L_g$ to get closer to a stationary point
2. else, if $\nabla^2 f(\mathbf{w}_k) \not\preceq \epsilon_h \mathbf{I}$ there is at least one negative eigenval
find the eigenvector, \mathbf{v} , of the smallest eigenvalue (λ_{\min}^k) of $\nabla^2 f(\mathbf{w}_k)$, namely

$$\|\mathbf{v}_k\|_2^2 = 1, \quad \mathbf{v}^T \nabla^2 f(\mathbf{w}_k) \mathbf{v} = \lambda_{\min}^k, \quad \nabla f(\mathbf{w}_k)^T \mathbf{v} \leq 0$$

3. Move toward that direction. **why? \mathbf{v} is a descent direction, since $\nabla f(\mathbf{w}_k)^T \mathbf{v} \leq 0$**

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \frac{2|\lambda_{\min}^k|}{L_H} \mathbf{v}$$

4. Otherwise, terminate.

\Rightarrow The number of iterations is at most $\max\left(\frac{2L_g}{\epsilon_g^2}, \frac{3L_H^2}{2\epsilon_h^3}\right) (f(\mathbf{w}_0) - f_{\inf})$

Notice from Taylor expansion that $f(\mathbf{w}_k + \mathbf{v}) \approx f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T \mathbf{v} + \frac{1}{2} \mathbf{v}^T \nabla^2 f(\mathbf{w}_k) \mathbf{v}$.
Step 2 finds a direction (\mathbf{v}) that gives the highest reduction to f . Do we need to really find the minimum eigenvalue? or a strong negative one is enough?

A generic Hessian-based algorithm

1. If $\|\nabla f(\mathbf{w}_k)\|_2 > \epsilon_g$, run your favorite algorithm (say GD with step-size $1/L_g$ to get closer to a stationary point
2. else, if $\nabla^2 f(\mathbf{w}_k) \not\preceq \epsilon_h \mathbf{I}$ there is at least one negative eigenval
find the eigenvector, \mathbf{v} , of the smallest eigenvalue (λ_{\min}^k) of $\nabla^2 f(\mathbf{w}_k)$, namely

$$\|\mathbf{v}_k\|_2^2 = 1, \quad \mathbf{v}^T \nabla^2 f(\mathbf{w}_k) \mathbf{v} = \lambda_{\min}^k, \quad \nabla f(\mathbf{w}_k)^T \mathbf{v} \leq 0$$

3. Move toward that direction. **why? \mathbf{v} is a descent direction, since $\nabla f(\mathbf{w}_k)^T \mathbf{v} \leq 0$**

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \frac{2|\lambda_{\min}^k|}{L_H} \mathbf{v}$$

4. Otherwise, terminate.

\Rightarrow The number of iterations is at most $\max\left(\frac{2L_g}{\epsilon_g^2}, \frac{3L_H^2}{2\epsilon_h^3}\right) (f(\mathbf{w}_0) - f_{\inf})$

Notice from Taylor expansion that $f(\mathbf{w}_k + \mathbf{v}) \approx f(\mathbf{w}_k) + \nabla f(\mathbf{w})^T \mathbf{v} + \frac{1}{2} \mathbf{v}^T \nabla^2 f(\mathbf{x}) \mathbf{v}$.
Step 2 finds a direction (\mathbf{v}) that gives the highest reduction to f . Do we need to really find the minimum eigenvalue? or a strong negative one is enough?

Outline

1. Motivation
2. Optimization of problems without structure
 - First-order methods
 - Successive Approximation
3. Optimization of problems with structure

Successive Approximation Methods

$$\arg \min_{\boldsymbol{w} \in \mathcal{D}} f(\boldsymbol{w})$$

Idea: minimize an **approximate surrogate** fnc.

Update approximation iteratively to make it locally tight:

$$\boldsymbol{w}_{k+1} = \arg \min_{\boldsymbol{w} \in \mathcal{D}} \tilde{f}(\boldsymbol{w}; \boldsymbol{w}_k)$$

- $\tilde{f}(\boldsymbol{w}; \boldsymbol{w}_k)$ is the approximation of f at \boldsymbol{w}_k

Two intuitive choices for $\tilde{f}(\boldsymbol{w}; \boldsymbol{w}_k)$:

- **Gradient:** Successive Linear Approximation
- **Hessian:** Successive Convex Approximation

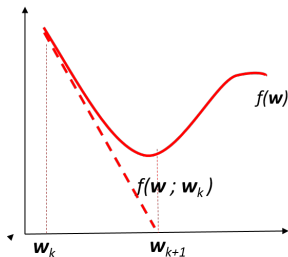
Successive Linear Approximation

Successive Linear Approximation (SLA)

$$\mathbf{w}_{k+1} = \arg \min_{\mathbf{w} \in \mathcal{D}} \tilde{f}(\mathbf{w}; \mathbf{w}_k) \quad (2)$$

$\tilde{f}(\mathbf{w}; \mathbf{w}_k)$ is first order approx of f at \mathbf{w}_k :

$$\begin{aligned} \tilde{f}(\mathbf{w}; \mathbf{w}_k) := & f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T (\mathbf{w} - \mathbf{w}_k) \\ & + (\gamma/2) \|\mathbf{w} - \mathbf{w}_k\|_2^2 \end{aligned}$$



Theorem 1: suppose that f is differentiable and Lipschitz continuous. Thus, the updates in (2) converge to a **stationary point** of f , as $k \rightarrow \infty$

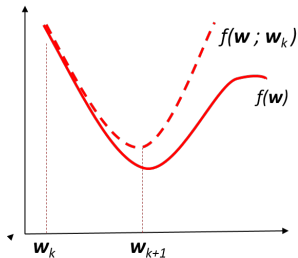
Successive Convex Approximation

Successive Convex Approximation (SCA)

$$\mathbf{w}_{k+1} = \arg \min_{\mathbf{w} \in \mathcal{D}} \tilde{f}(\mathbf{w}; \mathbf{w}_k) \quad (3)$$

$\tilde{f}(\mathbf{w}; \mathbf{w}_k)$ is convex approx of f at \mathbf{w}_k :

$$\begin{aligned} \tilde{f}(\mathbf{w}; \mathbf{w}_k) := & f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T (\mathbf{w} - \mathbf{w}_k) \\ & + (\mathbf{w} - \mathbf{w}_k)^T \nabla^2 f(\mathbf{w}_k) (\mathbf{w} - \mathbf{w}_k) \end{aligned}$$



Theorem 2: suppose that f is differentiable and Lipschitz continuous. Thus, the updates in (3) converge to a **stationary point** of f , as $k \rightarrow \infty$

Successive Approximation Methods

Take-home messages:

- In general, **no guarantees** on global convergence only convergence to stationary point
- Quality of solution depends on initial point
- SCA converges quadratically, SLA converges linearly
- Computational complexity too high for SCA: $\mathcal{O}(d^3)$ from inverting Hessian matrix

Outline

1. Motivation
2. Optimization of problems without structure
 - First-order methods
 - Successive Approximation
3. Optimization of problems with structure

Formulation into equivalent problem

Some nonconvex problems are **equivalently reformulated** into convex form.
-see equivalence between optimization problems (Chapter 3, [Boyd, 2004])

Rayleigh-Ritz quotient:

basis of **linear discriminant analysis** in dimensionality reduction
given **two classes of points** generated by two distributions with cov matrices, where $A \succeq 0, B \succ 0$

find the hyperplane, w , maximizing expected dist among two classes after projection on w

$$\max_{\|w\|_2=1} \frac{w^T A w}{w^T B w}$$

the optimization problem is equivalent to the min eigenvector. How ?

Minimum eigenvector of a symmetric matrix A : $\min_{\|w\|_2=1} w^T A w$

Nonconvex in the Euclidean space. why ?

Globally optimal solution by equivalently formulating the problem on the Riemann manifold

In general no systematic approach (rather problem specific)

Conclusions

Solution method depends on structure of problem

Non-convex optimization w/out structure:

- first-order methods: with tweaks to escape saddle pts
- Successive linear/convex approximation

Non-convex optimization with structure: (Next)

Some references

- Stephen J. Wright, "Coordinate descent algorithm" , Math. Program., 2015.
- S. Reddi, S. Sra, B. Póczos, and A. Smola, "Proximal stochastic methods for nonsmooth nonconvex finite-sum optimization," , NIPS, 2016.
- S. Shalev-Shwartz and T. Zhang, "Stochastic dual coordinate ascent methods for regularized loss minimization," JMLR, 2012.
- Y. Li and Y. Yuan, "Convergence analysis of two-layer neural networks with ReLU activation," , NIPS, 2017.
- M. Razaviyayn, M. Hong, and Z.-Q. Luo, "A unified convergence analysis of block successive minimization methods for nonsmooth optimization" , SIAM J. Optim., 2013.
- M. Hong, M. Razaviyayn, Z. Q. Luo and J. S. Pang, "A unified algorithmic framework for block-structured optimization involving big data: With applications in machine learning and signal processing," IEEE SPM, 2016.
- Z. Allen-Zhu, "Recent advances in stochastic convex and non-convex optimization," ICML Tutorial, 2017.
- S. Arora, R. Ge, T. Ma, and A. Moitra, "Simple, efficient, and neural algorithms for sparse coding," JMLR, 2015.
- C. Jin, R. Ge, P. Netrapalli, S.M. Kakade, and M.I. Jordan, "How to escape saddle points efficiently," ICML, 2017.
- Z. Allen-Zhu, "Natasha 2: Faster non-convex optimization than SGD," NIPS, 2018.
- Y. Zhang, P. Liang, and M. Charikar, "A hitting time analysis of stochastic gradient Langevin dynamics," COLT, 217.



Deep Learning

Lecture 4: Non-convex Optimization for Learning (Part 1)

University of Agder,
Kristiansand, Norway

Prepared by: Hadi Ghauch^{*}, Hossein S. Ghadikolaei[†]

^{*} Telecom ParisTech

[†] Royal Institute of Technology, KTH

<https://sites.google.com/view/fundl/home>