

**A GRADUATE COURSE
IN
OPTIMIZATION**

A GRADUATE COURSE
IN
OPTIMIZATION
CIE6010 Notebook

Prof. Yin Zhang

The Chinese University of Hong Kong, Shenzhen



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

Contents

Acknowledgments	vii
Notations	ix
1 Week1	1
1.1 Monday	1
1.1.1 Introduction to Optimizaiton	1
1.2 Wednesday	2
1.2.1 Reviewing for Linear Algebra	2
1.2.2 Reviewing for Calculus	2
1.2.3 Introduction to Optimization	3
2 Week2	7
2.1 Monday	7
2.1.1 Reviewing and Announments	7
2.1.2 Quadratic Function Case Study	8
2.2 Wednesday	11
2.2.1 Convex Analysis	11
3 Week3	17
3.1 Wednesday	17
3.1.1 Convex Analysis	17
3.1.2 Iterative Method	18
3.2 Thursday	22
3.2.1 Announcement	22
3.2.2 Sparse Large Scale Optimization	22

4	Week4	27
4.1	Wednesday	27
4.1.1	Local Convergence Rate	27
4.1.2	Newton's Method	29

Acknowledgments

This book is from the CIE6010 in fall semester, 2018.

CUHK(SZ)

Notations and Conventions

X	Set
$\inf X \subseteq \mathbb{R}$	Infimum over the set X
$\mathbb{R}^{m \times n}$	set of all $m \times n$ real-valued matrices
$\mathbb{C}^{m \times n}$	set of all $m \times n$ complex-valued matrices
x_i	i th entry of column vector \mathbf{x}
a_{ij}	(i, j) th entry of matrix \mathbf{A}
\mathbf{a}_i	i th column of matrix \mathbf{A}
\mathbf{a}_i^T	i th row of matrix \mathbf{A}
\mathbb{S}^n	set of all $n \times n$ real symmetric matrices, i.e., $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $a_{ij} = a_{ji}$ for all i, j
\mathbb{H}^n	set of all $n \times n$ complex Hermitian matrices, i.e., $\mathbf{A} \in \mathbb{C}^{n \times n}$ and $\bar{a}_{ij} = a_{ji}$ for all i, j
\mathbf{A}^T	transpose of \mathbf{A} , i.e, $\mathbf{B} = \mathbf{A}^T$ means $b_{ji} = a_{ij}$ for all i, j
\mathbf{A}^H	Hermitian transpose of \mathbf{A} , i.e, $\mathbf{B} = \mathbf{A}^H$ means $b_{ji} = \bar{a}_{ij}$ for all i, j
$\text{trace}(\mathbf{A})$	sum of diagonal entries of square matrix \mathbf{A}
$\mathbf{1}$	A vector with all 1 entries
$\mathbf{0}$	either a vector of all zeros, or a matrix of all zeros
\mathbf{e}_i	a unit vector with the nonzero element at the i th entry
$\mathcal{C}(\mathbf{A})$	the column space of \mathbf{A}
$\mathcal{R}(\mathbf{A})$	the row space of \mathbf{A}
$\mathcal{N}(\mathbf{A})$	the null space of \mathbf{A}
$\text{Proj}_{\mathcal{M}}(\mathbf{A})$	the projection of \mathbf{A} onto the set \mathcal{M}

Chapter 4

Week4

4.1. Wednesday

You need to take care of several things during your assignment:

1. Do Not Repeat Computation! e.g., if you enter

$$|f(x^{k+1}) - f(x^k)| \leq 10^{-\varepsilon} |f(x^k)|,$$

this is very bad, since you evaluated this function three times.

2. Arrange Computation Properly. e.g., compute

$$\mathbf{x}\mathbf{x}^T\mathbf{A}$$

is very expensive, but $\mathbf{x}(\mathbf{x}^T\mathbf{A})$ is not. Be aware of the size of matrices / vectors.

3. Appreciate sparsity, e.g., to compute $\mathbf{A}\mathbf{D}\mathbf{A}^T$ is bad if using $\mathbf{D} = \text{diag}(\mathbf{d})$, while using $\mathbf{D} = \text{sparse}(1:n, 1:n, \mathbf{d})$ is better.

Your code should be at least faster than the testing script.

4.1.1. Local Convergence Rate

Definition 4.1 [Q_1 Factor] Let $x^k \rightarrow x^*$ and $e_k := \|x^k - x^*\| \rightarrow 0$. The Q_1 factor of $\{x^k\}$ is given as:

$$Q_1 = \limsup_{k \rightarrow \infty} \frac{\|e^{k+1}\|}{\|e^k\|}$$

- The sequence $\{x^k\}$ is convergent if $Q_1 \in [0,1]$.

•

$$Q_1 = \begin{cases} 0, & \text{Q-super linear convergence} \\ \rho \in (0,1), & \text{Q-linear convergence} \\ 1, & \text{Q-sublinear} \end{cases}$$

R Q-linear convergence is not always so good, e.g., $\rho = .999$ may require 20000 iterations, while $\rho = .1$ may only require 20.

Linear is always better than sublinear; super linear is always better than linear. e.g.,

$$\{e_1^k\} = \{\beta^k\}, \quad \{e_2^k\} = \left\{\frac{1}{k^p}\right\} \text{ for } p = 1/2, 1, 2$$

Then $Q_1(e_1^k) = \beta \in (0,1)$; and $Q_1(e_2^k) = \lim_{k \rightarrow \infty} \frac{(k+1)^p}{k^p} = 1$. When will $e^k \leq \varepsilon$?

$$\beta^k \leq \varepsilon \implies k \geq \left(\frac{-1}{\ln \beta}\right) \ln \frac{1}{\varepsilon} = O(\ln \frac{1}{\varepsilon})$$

$$\frac{1}{k^p} \leq \varepsilon \implies k \geq \left(\frac{1}{\varepsilon}\right)^{1/p}$$

If $\varepsilon = 10^{-8}$, $O(\ln \frac{1}{\varepsilon})$ is not growing very fast. For $p = 1$ case, $k \geq O(\frac{1}{\varepsilon})$

We have a faster type of convergence:

Definition 4.2 [Q_2 Factor]

$$Q_2 = \limsup_{k \rightarrow \infty} \frac{e^{k+1}}{(e^k)^2}$$

If $Q_2 = M < +\infty$, i.e., $e^{k+1} = O((e^k)^2)$, then $\{x^k\} \rightarrow x^*$ Q-quadratically. ■

Newton's method generally gives us the quadratic convergence.

4.1.2. Newton's Method

The newton's method requires to solve a non-linear system of equations: $(\nabla f(x) = 0)$

$$F : \mathbb{R}^n \rightarrow \mathbb{R}^n, F(x) = 0.$$

We don't know how to solve non-linear system in general. But Newton gives us the remediation: do the linearization.

$$F(x + d) \approx F(x) + \langle F(x), \mathbf{d} \rangle = 0$$

To get the solution, it suffices to solve

$$\mathbf{d} = -(F'(x))^{-1}F(x),$$

and hence $\mathbf{x} \leftarrow \mathbf{x} + \alpha \mathbf{d}$.

Algorithm. Choose \mathbf{x}^0 .

For $k = 0, 1, 2, \dots$, solve the system w.r.t. \mathbf{d} :

$$\nabla^2 f(\mathbf{x}^k) \mathbf{d} = -\nabla f(\mathbf{x}^k)$$

Set $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha^k \mathbf{d}^k$.

End.

Interpretation. In order to minimize a strictly convex function $f(x)$, we find

$$f(\mathbf{x} + \mathbf{d}) \approx f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^T \nabla^2 f(\mathbf{x}) \mathbf{d} := q(\mathbf{d})$$

It suffices to minimize $q(\mathbf{d})$:

$$\nabla q(\mathbf{d}) = 0 \iff \nabla f(\mathbf{x}) + \nabla^2 f(\mathbf{x}) \mathbf{d} = 0$$

How to guarantee \mathbf{d} is the descent direction? It becomes an art.

Convergence of Rate. For $F(x) = 0$, suppose $\{\mathbf{x}^k\}$ is generated by Newton. As $\mathbf{x}^k \rightarrow \mathbf{x}^*$, $F(\mathbf{x}^*) = 0$. From

$$\mathbf{x}^{k+1} = \mathbf{x}^k - [F'(\mathbf{x}^k)]^{-1}F(\mathbf{x}^k),$$

we find

$$\begin{aligned}\mathbf{x}^{k+1} - \mathbf{x}^* &= \mathbf{x}^k - \mathbf{x}^* - [F'(\mathbf{x}^k)]^{-1} \left(F(\mathbf{x}^k) - F(\mathbf{x}^*) \right) \\ &= [F'(\mathbf{x}^k)]^{-1} \left(F(\mathbf{x}^*) - F(\mathbf{x}^k) - F'(\mathbf{x}^k)(\mathbf{x}^* - \mathbf{x}^k) \right)\end{aligned}$$

It follows that

$$\|\mathbf{x}^{k+1} - \mathbf{x}^*\| \leq \|[F'(\mathbf{x}^k)]^{-1}\| O(\|\mathbf{x}^k - \mathbf{x}^*\|^2) = O(\|\mathbf{x}^k - \mathbf{x}^*\|^2)$$

During this deviation we assume two things:

- Step-size is 1!
- The limit exists.

Hence, in practice, to implement Newton's method, 1 is the first choice; but gradient descent method is not.

Newton's method is good for nice problems, i.e., the function is convex, and the inverse of gradient is easy to solve.

In machine learning most time we implement the gradient descent method.

Learn and implement these things by yourself:

- Luo's note: Lecture #4, P7, Nestorov Accelerated Method

$\{a_r\}, r = 0, 1, 2, \dots$ we have

$$a_r = \frac{1}{2}(1 + \sqrt{1 + 4a_{r-1}^2}) \text{ with } a_0 = 0$$

From a_r we generate $t^r = (a_r - 1)/a_{r+1}$

Algorithm. Set $x^0 = x^1 = 0$.

for $r = 1, 2, \dots$

compute a^{r+1}, t^r

$$\mathbf{y}^{r+1} = (1 + t^r)\mathbf{x}^r - t^r\mathbf{x}^{r-1}$$

$$\mathbf{x}^{r+1} = \mathbf{y}^{r+1} - \frac{1}{L}\nabla f(\mathbf{y}^{r+1})$$

end

