

Numerical Methods for Convex Optimization and Their Applications

Jiyue Zeng, supervised by Martin Licht and Mareike Dressler

April 2021

Preface

This thesis mainly studies some numerical methods for constrained and unconstrained minimization problem for convex functions. Topics cover gradient descent method, conjugate gradient method, Newton's method, barrier method and finite element method. This thesis explains the intuition of the design of these methods, provides detailed algorithms, and includes some examples and graphs to illustrate the effect of each method. This thesis provides convergence analysis for some methods, which explains why the algorithm works.

Here is a guide for readers: the first chapter lists all notations that are used frequently in this thesis. The second chapter summarizes background knowledge about convex functions and some good properties they have. This section introduces the object that we will be working on for the rest of the paper. The third chapter briefly introduces the statement of a convex optimization problem and illustrates gradient descent method, conjugate gradient method and their applications. The fourth chapter introduce Newton's descent direction, Newton decrement, two phases convergence analysis and some examples, all of which are key ingredients for Newton's method. In addition, it talks about the application of Newton's method to a specific kind of function, the self-concordant function. The fifth chapter shows the effect of the barrier method(or Interior-point method) for inequality and equality constrained convex optimization problem. Newton's method plays an important role in the development of barrier method. The last chapter briefly talks about the finite element method. This method cuts the domain of the function into segments and approximates the function by linear functions.

Contents

0.1	Introduction	3
0.2	Preliminaries	4
0.2.1	Notation	4
0.2.2	Affine Sets	4
0.2.3	Convex Set	5
0.2.4	Cones	5
0.2.5	The Gradient and Hessian Matrix	5
0.2.6	Convex function	6
0.2.7	Affine function	9
0.3	Convex optimization	10
0.3.1	Introduction	11
0.3.2	Gradient Descent Method	12
0.3.3	Convergence Analysis of Gradient Descent Method	14
0.3.4	Conjugate Gradient Method	17
0.4	Newton's Method	26
0.4.1	The descent direction	27
0.4.2	The Newton Decrement	28
0.4.3	Convergence Analysis of Newton's Method	29
0.4.4	Algorithm	33
0.4.5	Examples	34
0.4.6	Newton's Method for Self-concordant Functions	36
0.5	Interior-point Method	44
0.5.1	Primal and Dual Problem	44
0.5.2	Newton's Method with Equality constraints	46
0.5.3	Barrier Method and Logarithmic Barrier Function	49
0.6	Finite Element Method	53

0.1 Introduction

0.2 Preliminaries

0.2.1 Notation

0.2.2 Affine Sets

Definition 0.2.2.1. We say that a set $C \subseteq \mathbb{R}^n$ is affine if for every $x, y \in C$ and $\lambda \in \mathbb{R}$, $\lambda x + (1 - \lambda)y \in C$, i.e., for some $x, y \in C, \alpha, \beta \in \mathbb{R}, \alpha x + \beta y \in C$ if $\alpha + \beta = 1$.

Lemma 0.2.2.2. More elements can be added as long as the sum of all coefficients is one. Then we can generalize this definition as follows: if C is an affine set, then, for all $x_1, x_2, \dots, x_n \in C$, $\lambda_1 x_1 + \lambda_2 x_2 + \dots + \lambda_n x_n \in C$ if $\lambda_1 + \lambda_2 + \dots + \lambda_n = 1$.

Definition 0.2.2.3. The linear combination $\lambda_1 x_1 + \lambda_2 x_2 + \dots + \lambda_n x_n$ where $\lambda_1 + \lambda_2 + \dots + \lambda_n = 1$ is called as an **affine combination**.

With this definition, we know that an affine set contains all affine combinations of its elements.

Definition 0.2.2.4. A set $S \subseteq \mathbb{R}^n$ is defined to be an **affine subspace** of \mathbb{R}^n if there exist a point $p \in \mathbb{R}^n$ and a subspace $U \subseteq \mathbb{R}^n$ such that

$$S = p + U = \{p + u \mid u \in U\}.$$

Note that since a vector space is closed under addition and scalar multiplication, it is also an affine set. Intuitively speaking, an affine set is developed by shifting the vector space along the direction of a vector away from the origin.

Example 0.2.2.5. If $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$, then $C = \{x \in \mathbb{R}^n \mid Ax = b\}$ is an affine subset of \mathbb{R}^n . To see this, choose $x_1, x_2 \in C$. Then we know that $Ax_1 = b$ and $Ax_2 = b$. For an arbitrary $\alpha \in \mathbb{R}$,

$$A(\alpha x_1 + (1 - \alpha)x_2) = \alpha Ax_1 + (1 - \alpha)Ax_2 = b.$$

Hence, for all $\alpha \in \mathbb{R}$, $\alpha x_1 + (1 - \alpha)x_2 \in C$, which means that C is an affine set.

Theorem 0.2.2.6. Every proper affine subspace V of \mathbb{R}^n is in the form of $\{x \in \mathbb{R}^n \mid Ax = b\}$ for some $A \in \mathbb{R}^{m \times n}$ with linearly independent rows and $b \in \mathbb{R}^m$.

Proof. Since V is an affine subspace, we know that there exists $p \in \mathbb{R}^n$ and a subspace U of \mathbb{R}^n such that $V = p + U$. Suppose U has a basis $\{u_1, u_2, \dots, u_k\}$, $k < n$. We want to find an A such that $Ap = b$ and $AU = 0$. Now, we define a matrix $W = [u_1 \ u_2 \ \dots \ u_k]$. Then we find a basis $\{a_1, a_2, \dots, a_m\}$ with $m = n - k$ of the set $\{a_i \in \mathbb{R}^n \mid W^T a_i = 0\}$. We define $A = [a_1 \ a_2 \ \dots \ a_m]^T$ and $b = Ap$. Finally, we get $AV = A(p + U) = b$, which gives the result. \square

Definition 0.2.2.7. For an arbitrary set $S \subseteq \mathbb{R}^n$, the **affine hull**, denoted $\text{aff}(S)$, is a set containing all affine combinations of elements in S :

$$\text{aff}(S) = \left\{ \sum_{i=1}^k \lambda_i x_i \mid x_1, x_2, \dots, x_k \in S, \lambda_1, \lambda_2, \dots, \lambda_k \in \mathbb{R}, \sum_{i=1}^k \lambda_i = 1 \right\}.$$

The affine hull $\text{aff}(S)$ is the smallest affine set that contains S .

0.2.3 Convex Set

Definition 0.2.3.1. A set S is a **convex set** if for every two points in S the line connecting these two points is contained in S . Mathematically, S is a convex set if for every $x_1, x_2 \in S$, $\lambda x_1 + (1 - \lambda)x_2 \in S$ for every $0 \leq \lambda \leq 1$. To generalize this definition, if S is convex, then for every $x_1, x_2, \dots, x_n \in S$, $\lambda_1 x_1 + \lambda_2 x_2 + \dots + \lambda_n x_n \in S$ with $\lambda_1 + \lambda_2 + \dots + \lambda_n = 1$ and $\lambda_1, \lambda_2, \dots, \lambda_n \geq 0$.

From the definition, an affine set is automatically a convex set.

Definition 0.2.3.2. The linear combination $\lambda_1 x_1 + \lambda_2 x_2 + \dots + \lambda_n x_n$ with $\lambda_1 + \lambda_2 + \dots + \lambda_n = 1$ and $\lambda_i \geq 0$ for $1 \leq i \leq n$ is called a **convex combination**.

Definition 0.2.3.3. The **convex hull** of a set S , denoted $\text{conv}(S)$, is the set of all convex combinations of elements in S . In other words,

$$\text{conv}(S) = \left\{ \sum_{i=1}^n \lambda_i x_i \mid \sum_{i=1}^n \lambda_i = 1, \lambda_i \geq 0, x_i \in S \right\}.$$

The convex hull $\text{conv}(S)$ is the smallest convex set that contains S , i.e., if $S \subseteq U$ and U is a convex set, then we must have $\text{conv}(S) \subseteq U$.

0.2.4 Cones

Definition 0.2.4.1. A set S is a **cone** if for every $x \in S$ and every $\lambda \geq 0$, $\lambda x \in S$. We say a set S is a **convex cone** if it is convex and is a cone, i.e., for every $x_1, x_2 \in S$, $\lambda_1, \lambda_2 \geq 0$, we have $\lambda_1 x_1 + \lambda_2 x_2 \in S$.

Proposition 0.2.4.2. The intersection of two convex cones in the same vector space is a convex cone, but the union may not be.

Proof. Suppose we have two convex cones, S and U . Assume $x, y \in S \cap U$ and $\alpha, \beta \geq 0$. Since S and U are two convex cones, $\alpha x + \beta y \in S$ as well as $\in U$, which means that $\alpha x + \beta y \in S \cap U$. Hence, $S \cap U$ is a convex cone. However, the union may not be. For example, the union of two different lines passing through the origin in \mathbb{R}^2 . \square

0.2.5 The Gradient and Hessian Matrix

Definition 0.2.5.1. The **gradient** of a scalar-valued differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is defined as $\nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$,

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(\mathbf{x}) \\ \frac{\partial f}{\partial x_2}(\mathbf{x}) \\ \vdots \\ \frac{\partial f}{\partial x_n}(\mathbf{x}) \end{bmatrix}.$$

To write the gradient as a linear combination of standard basis \mathbf{e}_i , we get

$$\nabla f(\mathbf{x}) = \frac{\partial f}{\partial x_1}(\mathbf{x})\mathbf{e}_1 + \frac{\partial f}{\partial x_2}(\mathbf{x})\mathbf{e}_2 + \dots + \frac{\partial f}{\partial x_n}(\mathbf{x})\mathbf{e}_n.$$

It measures how fast the function changes with respect to each standard basis vector \mathbf{e}_i .

Theorem 0.2.5.2. *A differentiable function f increases the fastest along the direction of the gradient.*

Proof. Let \mathbf{u} be a unit vector. The dot product $\nabla f(x) \cdot \mathbf{u}$ is the directional derivative of f at \mathbf{x} along \mathbf{u} , which measures the rate of change of f along \mathbf{u} . Since $\nabla f(x) \cdot \mathbf{u} = |\nabla f(x)| |\mathbf{u}| \cos(\theta)$, we choose \mathbf{u} to be the unit vector along the direction of the gradient, that is $\mathbf{u} = \frac{\nabla f(x)}{|\nabla f(x)|}$. Then $\nabla f(x) \cdot \frac{\nabla f(x)}{|\nabla f(x)|}$ reaches its maximum $|\nabla f(x)|$ since $\cos(\theta) = 1$. Hence, $\nabla f(x)$ is the direction that makes the function increases the most. \square

Definition 0.2.5.3. *Suppose we have a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ whose second partial derivatives are defined and continuous all over the domain. The **Hessian matrix** \mathbf{H}_f of f is defined to be*

$$\mathbf{H}_f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{bmatrix}.$$

The Hessian matrix is symmetric because all second partial derivatives are continuous, and \mathbf{H}_f can also be denoted as $\nabla^2 f$.

Definition 0.2.5.4. *Suppose we have a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ whose first derivatives all exist. The **Jacobian matrix** of f is a $m \times n$ matrix defined as*

$$\mathbf{J} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}.$$

The gradient is the Jacobian matrix of a scalar-valued function. The Hessian matrix is the Jacobian matrix of the gradient of f , i.e., $\mathbf{H}_f = \mathbf{J}(\nabla f)$.

0.2.6 Convex function

Definition 0.2.6.1. *Suppose we have a convex set $S \subseteq \mathbb{R}^n$ and a function $f : S \rightarrow \mathbb{R}$. The function f is convex if for all $\lambda \in [0, 1]$ and for all $\mathbf{x}, \mathbf{y} \in S$, we have $f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y})$.*

Note that a function is concave if it satisfies the opposite, i.e., if $f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \geq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y})$, then f is concave.

Theorem 0.2.6.2. *A local minimum of a convex function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is also a global minimum.*

Proof. Suppose we have a convex function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and its local minimum \mathbf{x}^* . Then there exists an open ball $B_r(\mathbf{x}^*)$ centered at \mathbf{x}^* with radius r and for all $\mathbf{x} \in B_r(\mathbf{x}^*)$, $f(\mathbf{x}^*) \leq f(\mathbf{x})$. We choose an arbitrary point $\mathbf{y} \neq \mathbf{x}^*$ from the domain of f . Then we choose a constant $\alpha \in (0, 1)$ such that $\alpha \mathbf{x}^* + (1 - \alpha)\mathbf{y} \in B_r(\mathbf{x}^*)$, which means the following:

$$\begin{aligned} \|\alpha \mathbf{x}^* + (1 - \alpha)\mathbf{y} - \mathbf{x}^*\| &\leq r \\ \implies (1 - \alpha) \|\mathbf{y} - \mathbf{x}^*\| &\leq r \\ \implies \alpha &\geq 1 - \frac{r}{\|\mathbf{y} - \mathbf{x}^*\|}. \end{aligned}$$

Hence, $\alpha \in (0, 1) \cap \left[1 - \frac{r}{\|\mathbf{y} - \mathbf{x}^*\|}, +\infty\right]$. It is easy to see that there exists a possible α . Since f is convex, we know that

$$\begin{aligned} f(\mathbf{x}^*) &\leq f(\alpha\mathbf{x}^* + (1 - \alpha)\mathbf{y}) \\ &\implies f(\mathbf{x}^*) \leq \alpha f(\mathbf{x}^*) + (1 - \alpha)f(\mathbf{y}) \\ &\implies (1 - \alpha)f(\mathbf{x}^*) \leq (1 - \alpha)f(\mathbf{y}) \\ &\implies f(\mathbf{x}^*) \leq f(\mathbf{y}). \end{aligned}$$

Hence, for any $\mathbf{y} \in \text{dom}(f)$, $f(\mathbf{x}^*) \leq f(\mathbf{y})$, which tells us that $f(\mathbf{x}^*)$ is a global minimum. \square

If the function is strictly convex, then it has a unique global minimum. Now we give another way to determine the convexity of a function.

Theorem 0.2.6.3. *A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if and only if f is convex along every line, i.e., $g : \mathbb{R} \rightarrow \mathbb{R}$, defined by $g(t) = f(\mathbf{x} + t\mathbf{v})$, is convex for all $\mathbf{x} \in \text{dom}(f)$, $\mathbf{v} \in \mathbb{R}^n$.*

Proof. (\Rightarrow) Since f is convex over \mathbb{R}^n , f is also convex on a line $\mathbf{x} + t\mathbf{v}$. Hence, g is a convex function.

(\Leftarrow) For any two points $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, we can find a vector $\mathbf{v} \in \mathbb{R}^n$ and $t \in \mathbb{R}$ such that $\mathbf{y} = \mathbf{x} + t\mathbf{v}$. Since g is convex, f is also convex. \square

Theorem 0.2.6.4. *Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable. The function f is convex if and only if for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $f(\mathbf{x}) \geq f(\mathbf{y}) + \nabla f(\mathbf{y})(\mathbf{x} - \mathbf{y})$.*

Proof. (\Rightarrow) We choose arbitrary $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. Let $z = \lambda\mathbf{x} + (1 - \lambda)\mathbf{y}$ for some $\lambda \in [0, 1]$. Since f is convex, we know that

$$f(z) = f(\lambda\mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y})$$

Subtracting $f(\mathbf{y})$ on both sides and then simplifying the equation, we get

$$\begin{aligned} f(\mathbf{y} + \lambda(\mathbf{x} - \mathbf{y})) - f(\mathbf{y}) &\leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) - f(\mathbf{y}), \\ f(\mathbf{y} + \lambda(\mathbf{x} - \mathbf{y})) - f(\mathbf{y}) &\leq \lambda(f(\mathbf{x}) - f(\mathbf{y})), \\ \frac{f(\mathbf{y} + \lambda(\mathbf{x} - \mathbf{y})) - f(\mathbf{y})}{\lambda} &\leq f(\mathbf{x}) - f(\mathbf{y}), \text{ for } \lambda \in (0, 1]. \end{aligned}$$

Taking the limit as λ goes to 0, we get

$$\lim_{\lambda \rightarrow 0} \frac{f(\mathbf{y} + \lambda(\mathbf{x} - \mathbf{y})) - f(\mathbf{y})}{\lambda} = \nabla f(\mathbf{y}) \cdot (\mathbf{x} - \mathbf{y}),$$

This computes the directional derivative of f at \mathbf{y} in the direction of vector $\mathbf{x} - \mathbf{y}$. Since the inequality holds for all $\lambda \in (0, 1]$, we get for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $f(\mathbf{x}) \geq f(\mathbf{y}) + \nabla f(\mathbf{y})(\mathbf{x} - \mathbf{y})$.

(\Leftarrow) We choose arbitrary $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and $\lambda \in [0, 1]$. Let $z = \lambda\mathbf{x} + (1 - \lambda)\mathbf{y}$. We know that $f(\mathbf{x}) \geq f(z) + \nabla f(z)(\mathbf{x} - z)$ and $f(\mathbf{y}) \geq f(z) + \nabla f(z)(\mathbf{y} - z)$. We multiply λ to the first inequality and $1 - \lambda$ to the second inequality and then get

$$\begin{aligned} \lambda f(\mathbf{x}) &\geq \lambda f(z) + \lambda(1 - \lambda)(\mathbf{x} - \mathbf{y}) \\ (1 - \lambda)f(\mathbf{y}) &\geq (1 - \lambda)f(z) + \lambda(1 - \lambda)(\mathbf{y} - \mathbf{x}) \end{aligned}$$

Adding up these two inequalities, we get

$$\lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) \geq f(\lambda\mathbf{x} + (1 - \lambda)\mathbf{y}).$$

Hence, f is convex. □

Theorem 0.2.6.5. *Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable. The function f is convex if and only if the Hessian matrix of f is positive semi-definite.*

Proof. (\Rightarrow) Since f is convex, from Theorem 1.5.4 we know that

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, f(\mathbf{x}) \geq f(\mathbf{y}) + \nabla f(\mathbf{y})(\mathbf{x} - \mathbf{y}).$$

Therefore,

$$\begin{aligned} f(\mathbf{y}) + \nabla f(\mathbf{y})(\mathbf{x} - \mathbf{y}) &\leq f(\mathbf{x}) \leq f(\mathbf{y}) + \nabla f(\mathbf{x})(\mathbf{x} - \mathbf{y}) \\ \nabla f(\mathbf{y})(\mathbf{x} - \mathbf{y}) &\leq f(\mathbf{x}) - f(\mathbf{y}) \leq \nabla f(\mathbf{x})(\mathbf{x} - \mathbf{y}) \\ (\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}))(\mathbf{x} - \mathbf{y}) &\geq 0. \end{aligned}$$

Since this inequality holds for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, we choose $\mathbf{x}, \mathbf{x} + \mathbf{h} \in \mathbb{R}^n$ for an arbitrary $\mathbf{h} \in \mathbb{R}^n$ and we have,

$$(\nabla f(\mathbf{x} + \mathbf{h}) - \nabla f(\mathbf{x}))\mathbf{h} \geq 0.$$

Since f is twice continuously differentiable, we have for all $t > 0$,

$$\nabla f(\mathbf{x} + t\mathbf{h}) - \nabla f(\mathbf{x}) = \nabla^2 f(\mathbf{x})t\mathbf{h} + \mathbf{r}(t\mathbf{h}),$$

and

$$\lim_{t \rightarrow 0} \frac{|\mathbf{r}(t\mathbf{h})|}{|t\mathbf{h}|} = 0.$$

Then,

$$\mathbf{h}^T(\nabla f(\mathbf{x} + t\mathbf{h}) - \nabla f(\mathbf{x})) = \mathbf{h}^T \nabla^2 f(\mathbf{x})t\mathbf{h} + \mathbf{h}^T \mathbf{r}(t\mathbf{h}) \geq 0,$$

for all $\mathbf{h} \in \mathbb{R}^n$ and $t > 0$. Taking the limit as t goes to 0, we have

$$\begin{aligned} \lim_{t \rightarrow 0} \frac{\mathbf{h}^T(\nabla f(\mathbf{x} + t\mathbf{h}) - \nabla f(\mathbf{x}))}{t} &= \lim_{t \rightarrow 0} \mathbf{h}^T \nabla^2 f(\mathbf{x})\mathbf{h} + \frac{\mathbf{h}^T \mathbf{r}(t\mathbf{h})}{t} \\ &= \mathbf{h}^T \nabla^2 f(\mathbf{x})\mathbf{h} + \lim_{t \rightarrow 0} \frac{\mathbf{h}^T \mathbf{r}(t\mathbf{h})}{t} \\ &= \mathbf{h}^T \nabla^2 f(\mathbf{x})\mathbf{h} \geq 0. \end{aligned}$$

This tells us that $\nabla^2 f(\mathbf{x}) \succeq 0$, which is the notation for a matrix to be positive semi-definite.

(\Leftarrow) According to the second-order Taylor polynomial for f , we have that $\forall \mathbf{y} \in \mathbb{R}^n$,

$$f(\mathbf{x}) = f(\mathbf{y}) + \nabla f(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^T \nabla^2 f(\mathbf{z})(\mathbf{y} - \mathbf{x}),$$

for some \mathbf{z} between \mathbf{x} and \mathbf{y} . Since $\nabla^2 f$ is positive semi-definite, we know that

$$\frac{1}{2}(\mathbf{y} - \mathbf{x})^T \nabla^2 f(\mathbf{z})(\mathbf{y} - \mathbf{x}) \geq 0,$$

which implies that

$$f(\mathbf{x}) \geq f(\mathbf{y}) + \nabla f(\mathbf{x})(\mathbf{y} - \mathbf{x}).$$

Theorem 1.5.4 tells us that f is a convex function. □

0.2.7 Affine function

Definition 0.2.7.1. A function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is **affine** if there exist a linear function $h : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and a constant vector $b \in \mathbb{R}^m$ such that $f(\mathbf{x}) = h(\mathbf{x}) + b$.

From the definition, an affine function is a linear function plus a constant vector. An affine function is both concave and convex:

$$\begin{aligned} f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) &= h(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) + b \\ &= \lambda h(\mathbf{x}) + (1 - \lambda) h(\mathbf{y}) + \lambda b + (1 - \lambda) b \\ &= \lambda f(x) + (1 - \lambda) f(y). \end{aligned}$$

0.3 Convex optimization

In this section, we introduce what a convex optimization problem is and provides two numerical methods, Gradient descent method and Conjugate gradient method. Gradient descent method is generated to iteratively find the minimum point of a convex and twice continuously differentiable function by using the negative gradient as the descent direction. Conjugate gradient method is an efficient way in solving large systems of linear equations.

0.3.1 Introduction

Definition 0.3.1.1. Suppose we have a convex function $f_0 : \mathbb{R}^n \rightarrow \mathbb{R}$, several other convex functions $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ for $1 \leq i \leq m$, and an affine function $h : \mathbb{R}^n \rightarrow \mathbb{R}^m$.

A **convex optimization problem** has the following form,

$$\begin{aligned} \min \quad & f_0(\mathbf{x}) \\ \text{subject to} \quad & f_i(\mathbf{x}) \leq 0, \text{ for } 1 \leq i \leq m \\ & h(\mathbf{x}) = 0. \end{aligned}$$

We define the **feasible set** of this convex optimization problem, denoted by X , to be

$$X = \left\{ \mathbf{x} \in \left(\bigcap_{i=0}^p \text{dom}(f_i) \right) \cap \text{dom}(h) \mid f_i(\mathbf{x}) \leq 0, \text{ for } 1 \leq i \leq p, \text{ and } h(\mathbf{x}) = 0 \right\},$$

where $\text{dom}(f)$ is the notation for the domain of the function f . The **optimal point** of this standard convex optimization problem is

$$\mathbf{x}^* = \inf \{f_0(\mathbf{x}) \mid \mathbf{x} \in X\}.$$

The problem becomes an unconstrained convex optimization problem if there is no $f_i \leq 0$ and $h = 0$ constraints.

Note that $h(\mathbf{x}) = 0$ can usually be written in the form of $A\mathbf{x} = \mathbf{b}$, or $a_i \cdot \mathbf{x} = b_i$, where a_i is the i -th row of A and b_i denotes the i -th entry of \mathbf{b} .

Since h is affine and f_0, f_i are convex, we know that the whole system is convex, which immediately tells that the a local minimum is the optimal solution to this convex problem.

On top of that, we know that the feasible set X must also be convex since it is an intersection of convex domains.

Theorem 0.3.1.2. Suppose f_0 is differentiable. The point \mathbf{x} is optimal if and only if for all \mathbf{y} in the feasible set X , $\nabla f_0(\mathbf{x})(\mathbf{y} - \mathbf{x}) \geq 0$.

Proof. (\Leftarrow) Since $\nabla f_0(\mathbf{x})(\mathbf{y} - \mathbf{x}) \geq 0$ and f_0 is convex, we know that for all $\mathbf{y} \in X$, $f_0(\mathbf{y}) \geq f_0(\mathbf{x}) + \nabla f_0(\mathbf{x})(\mathbf{y} - \mathbf{x}) \geq f_0(\mathbf{x})$. Hence, \mathbf{x} is optimal.

(\Rightarrow) We prove by contradiction. Suppose \mathbf{x} is optimal and there exists $\mathbf{y} \in X$ such that $\nabla f_0(\mathbf{x})(\mathbf{y} - \mathbf{x}) < 0$. Let $\mathbf{z} = (1 - t)\mathbf{x} + t\mathbf{y}$. Then

$$\lim_{t \rightarrow 0} \frac{f_0(\mathbf{z}) - f_0(\mathbf{x})}{t} = \lim_{t \rightarrow 0} \frac{f_0(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - f_0(\mathbf{x})}{t} = \nabla f_0(\mathbf{x})(\mathbf{y} - \mathbf{x}) < 0.$$

Hence, for a very small t , we can find \mathbf{z} very close to \mathbf{x} such that $f_0(\mathbf{z}) < f_0(\mathbf{x})$, which contradicts with the fact that \mathbf{x} is optimal. \square

Corollary 0.3.1.3. If there is no inequality or equality constraint, i.e., the problem is an unconstrained convex optimization problem, then \mathbf{x} is optimal if $\nabla f_0(\mathbf{x}) = 0$.

Proof. If \mathbf{x} is optimal, from Theorem 2.1.2, we know that for all $\mathbf{y} \in X$, which denotes the feasible set, $\nabla f_0(\mathbf{x})(\mathbf{y} - \mathbf{x}) \geq 0$. Choose $\lambda > 0$ small enough such that $\mathbf{y} = \mathbf{x} - \lambda \nabla f_0(\mathbf{x}) \in X$. Then we know that

$$\nabla f_0(\mathbf{x})(\mathbf{x} - \lambda \nabla f_0(\mathbf{x}) - \mathbf{x}) = -\lambda \|\nabla f_0(\mathbf{x})\|^2 \geq 0.$$

Hence, $\nabla f_0(\mathbf{x}) = 0$. \square

This corollary gives a stopping criterion, that is $\|\nabla f(x)\| \leq \epsilon$. Now we first consider the unconstrained convex optimization problem.

0.3.2 Gradient Descent Method

Suppose we have a convex function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ that is twice continuously differentiable and has a compact domain. We know that f must attain its minimum in the domain. Then how do we find the x that obtains the minimum?

We start from an arbitrary point on the graph and want to find a path that leads us to $\min f$. We can denote this searching path by a sequence of x_i in the domain. We define the path by the following equation,

$$x_{i+1} = x_i + t_i \Delta x_i \text{ for } i \geq 1,$$

where x_i is the current step, x_{i+1} is the next step, and $t_i \Delta x_i$ is the step size. The step size consists of a vector Δx_i that indicates the direction and a constant t_i that decides how much it moves along the direction. The following graph is an example of such an x_i path.

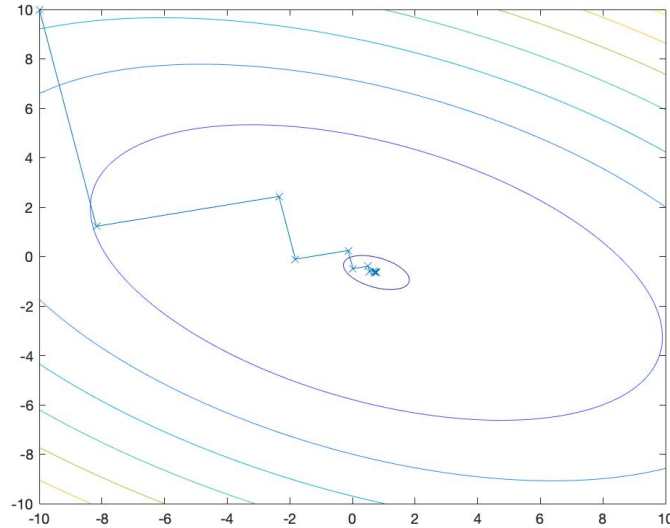


Figure 1: gradient descent with exact line search

The gradient descent method suggests that we should use the negative gradient, $-\nabla f(x_i)$, as the direction, i.e., $\Delta x_i = -\nabla f(x_i)$. This choice makes sense because Theorem 1.4.2 says that f increases the fastest along ∇f , which implies that f decreases the fastest along $-\nabla f$. The fact that f is convex tells that

$$\nabla f(x_i)(x_{i+1} - x_i) \geq 0 \Rightarrow f(x_{i+1}) \geq f(x_i).$$

Therefore, the precondition for obtaining a decreasing sequence is that

$$\nabla f(x_i)(x_{i+1} - x_i) \leq 0. \tag{1}$$

After replacing $(x_{i+1} - x_i)$ by the step size, we have

$$\nabla f(x_i)(-t \nabla f(x_i)) \leq 0$$

and then simplify it,

$$-t \|\nabla f(x_i)\|^2 \leq 0,$$

which is always true if $t \geq 0$. Hence, the precondition for obtaining a decreasing sequence is guaranteed when we use negative gradient. Notice that $-\nabla f(x_i)$ changes at each step depending on which x_i we plug in.

With the direction already defined, what is left is to find t_i for the step size. There are two ways to find the constant t_i .

1st Way 0.3.2.1. Exact Line Search

Define $t = \{t \geq 0 \mid t \text{ minimizes } f(x_i - t\nabla f(x_i))\}$. This guarantees that $f(x_i - t\nabla f(x_i)) \leq f(x_i)$ because t minimizes the function value. To find t , we take the derivative of $f(x_i - t\nabla f(x_i))$ with respect to t ,

$$\begin{aligned} \frac{d}{dt} f(x_i - t\nabla f(x_i)) &= \nabla f(x_i - t\nabla f(x_i)) \frac{d}{dt} (x_i - t\nabla f(x_i)) \\ &= -\nabla f(x_{i+1})^T \nabla f(x_i) \\ &= 0. \end{aligned} \tag{2}$$

Since $\nabla f(x_i)$ is the i -th direction and $\nabla f(x_{i+1})$ is the $i+1$ -th direction, the fact that their dot product is 0 implies that these two directions are orthogonal. Therefore, by using exact line search, we get a zigzag searching path.

The exact line search works when the function has such a nice formula for its gradient that we can solve for t explicitly. When the cost of computation is very expensive, we should switch to the backtracking line search.

Algorithm 1: Gradient Descent Method With Exact Line Search

```

graDesExact ( $f, x_0, \epsilon, \text{max\_iter}$ );
Input :  $x_0, \epsilon, \text{max\_iter}$ 
Output:  $x^*$ 
 $x_i = x_0$ ;
 $t = 0$ ;
 $\text{count} = 0$ ;
while  $\|\nabla f(x_i)\| > \epsilon$  &  $\text{count} \leq \text{max\_iter}$  do
     $\text{count} = \text{count} + 1$ ;
    update  $t$  by solving  $\frac{d}{dt} f(x_i - t\nabla f(x_i)) = 0$ ;
     $x_i = x_i - t\nabla f(x_i)$ ;
return  $x_i$ ;

```

2nd Way 0.3.2.2. Backtracking Line Search

Algorithm 2: Backtracking Line Search

```

initialization:  $x_i, x_{i+1} = x_i - t\nabla f(x_i), t = t_0 > 0, \gamma \in (0, 1), \alpha \in (0, \frac{1}{2})$ ;
while  $f(x_i) - f(x_{i+1}) < \alpha t \|\nabla f(x_i)\|^2$ , do
     $t = \gamma t, x_{i+1} = x_i - t\nabla f(x_i)$ ;
end

```

The while loop will quit after finite iterations because of the following.
By the second order Taylor's expansion of the function,

$$f(x_{i+1}) = f(x_i) - t \|\nabla f(x_i)\|^2 + \frac{1}{2}t^2 \nabla f(x_i)^T \nabla^2 f(y) \nabla f(x_i),$$

for some y between x_i and x_{i+1} . When t is very small, we know that the last term is of $\mathcal{O}(t^2)$ and can be ignored. Then, we know for $\alpha < \frac{1}{2}$,

$$f(x_i) - f(x_{i+1}) = t \|\nabla f(x_i)\|^2 > \alpha t \|\nabla f(x_i)\|^2.$$

We know that t can be arbitrarily small within finite iterations because $\gamma < 0$.

To be more precise about the possible range of t , we assume one more condition. Suppose the largest eigenvalue of $\nabla^2 f$ is bounded by M , i.e., $\nabla^2 f \preceq M\mathbb{I}$. The stopping criterion of the while loop can be rewritten as

$$t \|\nabla f(x_i)\|^2 - \frac{1}{2}t^2 \nabla f(x_i)^T \nabla^2 f(y) \nabla f(x_i) \geq \alpha t \|\nabla f(x_i)\|^2.$$

After simplifying the inequality, we get

$$(1 - \alpha)t \|\nabla f(x_i)\|^2 \geq \frac{1}{2}t^2 \nabla f(x_i)^T \nabla^2 f(y) \nabla f(x_i).$$

With the upper bound on $\nabla^2 f$, we know that

$$\frac{1}{2}t^2 \nabla f(x_i)^T \nabla^2 f(y) \nabla f(x_i) \leq \frac{1}{2}t^2 \|\nabla f(x_i)\|^2 M.$$

$$(1 - \alpha)t \|\nabla f(x_i)\|^2 \geq \frac{1}{2}t^2 \|\nabla f(x_i)\|^2 M,$$

$$t \leq \frac{2(1 - \alpha)}{M}. \quad (3)$$

Therefore, the while loop must end when $t \in (0, \min\{\frac{2(1-\alpha)}{M}, t_0\})$, which implies that $t \geq \frac{2(1-\alpha)\gamma}{M}$. The algorithm guarantees that there is a decent amount of decrease in f and the step size is not too big.

0.3.3 Convergence Analysis of Gradient Descent Method

Assume the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is strongly convex, which means that we can bound its Hessian matrix, $m\mathbb{I} \preceq \nabla^2 f \preceq M\mathbb{I}$. As before, we let $f(x^*)$ denote the minimum of f . By the second order Taylor's expansion, we have for all $x, y \in \text{dom}(f)$,

$$\begin{aligned} f(y) &= f(x) + \nabla f(x)(y - x) + \frac{1}{2}(y - x)^T \nabla^2 f(x)(y - x) \\ &\geq f(x) + \nabla f(x)(y - x) + \frac{1}{2}m \|y - x\|^2. \end{aligned}$$

Algorithm 3: Gradient Descent Method With Backtracking Line Search

graDesBack ($f, x_0, \epsilon, \text{max_iter}$);

Input : $x_0, \epsilon, \text{max_iter}$

Output: x^*

$x_i = x_0$;

$x_{i+1} = x_i - t \nabla f(x_i)$;

$t = 1$;

$\text{count} = 0$;

$\alpha = 0.25$;

$\gamma = 0.5$;

while $\|\nabla f(x_i)\| > \epsilon$ & $\text{count} \leq \text{max_iter}$ **do**

$\text{count} = \text{count} + 1$;

while $f(x_i) - f(x_{i+1}) < \alpha t \|\nabla f(x_i)\|^2$, **do**

$t = \gamma t$;

$x_{i+1} = x_i - t \nabla f(x_i)$;

$x_i = x_{i+1}$;

return x_i ;

Consider the righthand side as a quadratic function with variable $y - x$. It reaches the minimum when $y - x = \frac{1}{m} \nabla f(x)$. After taking the minimum on both sides, we get

$$\begin{aligned} f(x^*) &\geq f(x) + \nabla f(x) \left(\frac{1}{m} \nabla f(x) \right) + \frac{m}{2} \left\| \frac{1}{m} \nabla f(x) \right\|^2 \\ &= f(x) - \frac{1}{2m} \|\nabla f(x)\|^2. \end{aligned}$$

Therefore,

$$\|\nabla f(x)\|^2 \geq 2m(f(x) - f(x^*)) \quad (4)$$

From the second order Taylor's expansion, we get

$$f(x_i - t \nabla f(x_i)) = f(x_i) - t \|\nabla f(x_i)\|^2 + \frac{1}{2} t^2 \nabla f(x_i)^T \nabla^2 f(x_i) \nabla f(x_i). \quad (5)$$

Applying the upper bound M on (5),

$$f(x_i - t \nabla f(x_i)) \leq f(x_i) - t \|\nabla f(x_i)\|^2 + \frac{1}{2} t^2 M \|\nabla f(x_i)\|^2. \quad (6)$$

Consider both sides as functions of t . Then the righthand side is simply a real-valued single variable quadratic function whose minimum is obtained at $t = \frac{1}{M}$. The minimum of the lefthand side via exact line search is just $f(x_{i+1})$. Hence, after minimizing both sides over the variable t and applying (3), we have

$$\begin{aligned} f(x_{i+1}) &\leq f(x_i) - \frac{1}{2M} \|\nabla f(x_i)\|^2 \\ &\leq f(x_i) - \frac{m}{M} (f(x_i) - f(x^*)). \end{aligned}$$

Then subtract $f(x^*)$ on both sides,

$$f(x_i + 1) - f(x^*) \leq \left(1 - \frac{m}{M}\right) (f(x_i) - f(x^*)). \quad (7)$$

Since $m/M < 1$, we know that the sequence $\{f(x_i) - f(x^*)\}_{i=1}^{\infty}$ converges to 0. By recursively applying (7), we get

$$f(x_i) - f(x^*) \leq \left(1 - \frac{m}{M}\right)^i (f(x_0) - f(x^*)).$$

Things are slightly different when the backtracking line search is used. Based on the while condition from the algorithm of backtracking line search, we know that

$$f(x_i) - f(x_{i+1}) \geq \alpha t \|\nabla f(x_i)\|^2.$$

We can replace t by $\frac{2(1-\alpha)\gamma}{M}$ from (3),

$$\begin{aligned} f(x_i) - f(x_{i+1}) &\geq \frac{2\alpha(1-\alpha)\gamma}{M} \|\nabla f(x_i)\|^2 \\ &\geq \frac{4m\alpha(1-\alpha)\gamma}{M} (f(x_i) - f(x^*)). \end{aligned}$$

Subtracting $f(x^*)$ on both sides,

$$\begin{aligned} f(x_i) - f(x^*) &\geq f(x_{i+1}) - f(x^*) + \frac{4m\alpha(1-\alpha)\gamma}{M} (f(x_i) - f(x^*)), \\ \left(1 - \frac{4m\alpha(1-\alpha)\gamma}{M}\right) (f(x_i) - f(x^*)) &\geq f(x_{i+1}) - f(x^*). \end{aligned} \quad (8)$$

We know that $\frac{4m\alpha(1-\alpha)\gamma}{M} < 1$ because $\alpha(1-\alpha) \leq \frac{1}{4}$ and $\frac{m}{M}, \gamma < 1$. Hence, after applying (8) recursively, we get

$$f(x_i) - f(x^*) \leq \left(1 - \frac{4m\alpha(1-\alpha)\gamma}{M}\right)^i (f(x_0) - f(x^*)).$$

Now, we apply the gradient descent method to a quadratic function with exact line search and backtracking line search respectively.

Example 0.3.3.1. Suppose we have a convex quadratic function $f(\mathbf{x}) = \mathbf{x}^T A \mathbf{x} + b\mathbf{x} + c$ where $A = \begin{bmatrix} 2 & 1 \\ 5 & 7 \end{bmatrix}$, $b = \begin{bmatrix} 4 \\ 6 \end{bmatrix}$, $c = 5$. The optimal point x^* that gives the minimum value of f is $x^* = \begin{bmatrix} -1 \\ 0 \end{bmatrix}$.

We fix a starting point $x_0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$. Figure 1, 2 are graphs of gradient descent method with exact line search and backtracking line search respectively. The blue lines on the graph represent the path of x_i starting from $[1, 1]$, ending at $[0, -1]$. Figure 3, 4 are the projections of Figure 1, 2 onto the x, y plane.

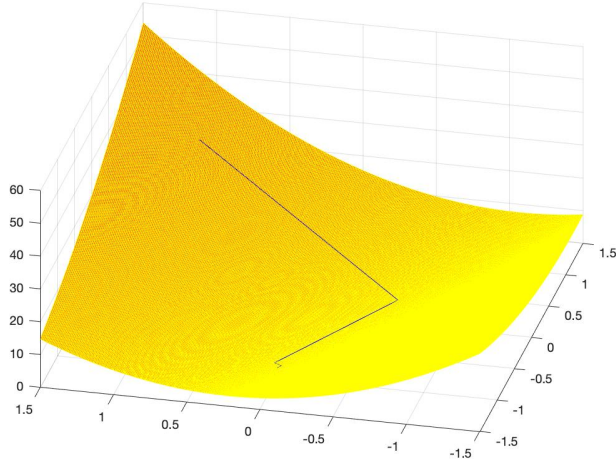


Figure 2: gradient descent with exact line search

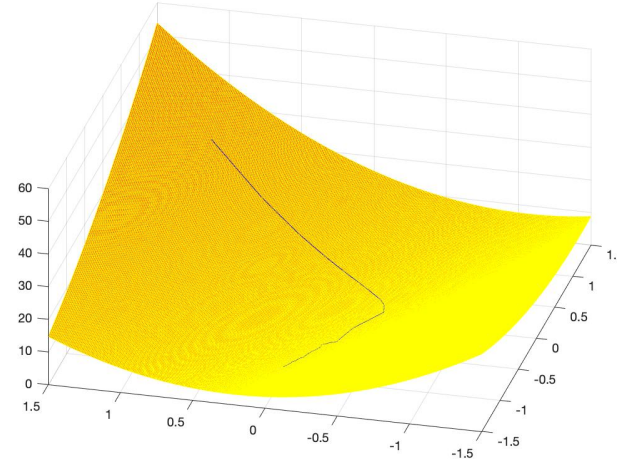


Figure 3: gradient descent with back-tracking line search

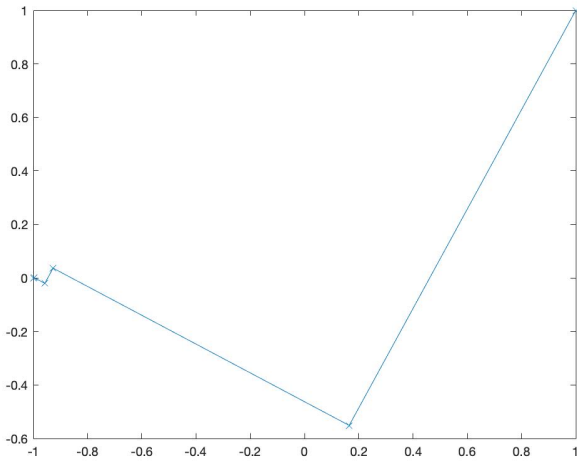


Figure 4: The searching path with exact line search

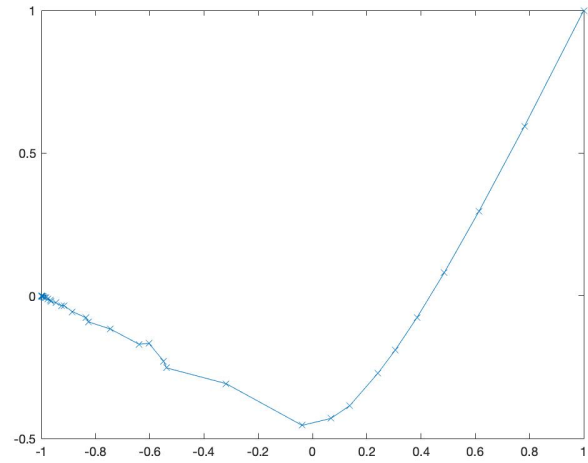


Figure 5: The searching path with back-tracking line search

From Figure 1,3, we see that the searching path of the exact line search has a zigzag shape. Comparing Figure 3,4, we also see that exact line search takes much fewer steps than backtracking line search. This is because on each step, exact line search finds the best x_{i+1} that minimize the function value while backtracking line search takes a small range of x_{i+1} that provides some decrease.

0.3.4 Conjugate Gradient Method

We narrow down the problem from minimizing a twice continuously differentiable convex function to minimizing a quadratic function, $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T A \mathbf{x} + b\mathbf{x} + c$. Before presenting the conjugate gradient method, we introduce the concept of A-conjugacy.

Definition 0.3.4.1. Suppose we have a positive definite square matrix $A \in \mathbb{R}^{n \times n}$. Two nonzero vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ are **A-conjugate** if $\mathbf{u}^T A \mathbf{v} = 0$.

See Figure 6 for geometric explanations.

Theorem 0.3.4.2. Suppose A is a positive definite $n \times n$ square matrix. If \mathbf{u}, \mathbf{v} are nonzero vectors that are mutually A-conjugate, then \mathbf{u}, \mathbf{v} are linearly independent.

Proof. In order to show linear independence, it suffice to show that $c_1 \mathbf{u} + c_2 \mathbf{v} = 0$ implies $c_1 = 0$, $c_2 = 0$. Multiplying $A \mathbf{v}$ on both sides, we get $c_1 \mathbf{u}^T A \mathbf{v} + c_2 \mathbf{v}^T A \mathbf{v} = 0$. Since A is positive definite, \mathbf{v} is not a zero vector, we get $\mathbf{v}^T A \mathbf{v} > 0$. Since \mathbf{u}, \mathbf{v} are A-conjugate, $\mathbf{u}^T A \mathbf{v} = 0$. Hence, $c_2 = 0$ and then $c_1 = 0$. \square

If A is the identity matrix, then A-conjugacy is the same as orthogonality. If the matrix A is instead a symmetric positive definite $n \times n$ square matrix and if we have a set of n mutually A-conjugate vectors $\{\mathbf{u}_i\}_{i=1}^n \in \mathbb{R}^n$ with respect to A , then $\{\mathbf{u}_i\}_{i=1}^n$ form a basis of \mathbb{R}^n . Any vector $\mathbf{x} \in \mathbb{R}^n$ can be written as $\mathbf{x} = \sum_{i=1}^n \alpha_i \mathbf{u}_i$, $\alpha_i \in \mathbb{R}$.

Theorem 0.3.4.3. Gram-Schmidt Algorithm

Suppose A is a symmetric positive definite $n \times n$ square matrix. If we start with a sequence of linearly independent vectors $\{\mathbf{v}_i\}_{i=1}^n \in \mathbb{R}^n$, we can generate a sequence of mutually A-conjugate vectors $\{\mathbf{u}_k\}_{k=1}^n \in \mathbb{R}^n$ from it.

Proof. Set $\mathbf{u}_1 = \mathbf{v}_1$. Let $\mathbf{u}_2 = \mathbf{v}_2 + \gamma_2 \mathbf{u}_1$, for some $\gamma \in \mathbb{R}$. Since we want $\mathbf{u}_2, \mathbf{u}_1$ to be mutually A-conjugate, we multiply $A \mathbf{u}_1$ on both sides:

$$\mathbf{u}_2^T A \mathbf{u}_1 = \mathbf{v}_2^T A \mathbf{u}_1 + \gamma_2 \mathbf{u}_1^T A \mathbf{u}_1$$

$$0 = \mathbf{v}_2^T A \mathbf{u}_1 + \gamma_2 \mathbf{u}_1^T A \mathbf{u}_1$$

$$\gamma = -\frac{\mathbf{v}_2^T A \mathbf{u}_1}{\mathbf{u}_1^T A \mathbf{u}_1}.$$

To finish the induction process, let $\mathbf{u}_k = \mathbf{v}_k + \sum_{j=1}^{k-1} \gamma_j \mathbf{u}_j$. We multiply $A \mathbf{u}_m$, for some $m \in \{1, 2, \dots, k-1\}$, on both sides of the equation:

$$\mathbf{u}_k^T A \mathbf{u}_m = \mathbf{v}_k^T A \mathbf{u}_m + \sum_{j=1}^{k-1} \gamma_j \mathbf{u}_j^T A \mathbf{u}_m.$$

We want \mathbf{u}_k to be A-conjugate to all previous $\mathbf{u}_i, i \in \{1, 2, \dots, k-1\}$ and $\mathbf{u}_i, \mathbf{u}_j$ are already mutually A-conjugate. Therefore, we obtain

$$0 = \mathbf{v}_k^T A \mathbf{u}_m + \gamma_m \mathbf{u}_m^T A \mathbf{u}_m.$$

$$\gamma_m = -\frac{\mathbf{v}_k^T A \mathbf{u}_m}{\mathbf{u}_m^T A \mathbf{u}_m}. \quad (9)$$

We can get all coefficients $\gamma_m, m \in \{1, 2, \dots, k-1\}$ by doing the same algebra. Hence, we have created a sequence of mutually A-conjugate vectors $\{\mathbf{u}_k\}_{k=1}^n$. \square

Now, we explain the main idea of the conjugate gradient descent method. Suppose we want to minimize a quadratic function $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T A \mathbf{x} + \mathbf{b}\mathbf{x} + c$, where A is a symmetric positive definite $n \times n$ square matrix. Here we use the same notations x^*, x_i, x_0 as in the last section. We generate a set of n mutually A-conjugate vectors $\{\mathbf{u}_k\}_{k=1}^n \in \mathbb{R}^n$ via Gram-Schmidt Algorithm. We define the **error** $e_i = x_i - x_0$. We can write e_i as a linear combination of $\{\mathbf{u}_k\}_{k=1}^n$, i.e., $e_i = \sum_{k=1}^n \alpha_k \mathbf{u}_k$.

The conjugate gradient method says that we just need to remove the error in one direction \mathbf{u}_k at each step, that is

$$\begin{aligned} x_0 &= x^* + e_0 = x^* + \sum_{k=1}^n \alpha_k \mathbf{u}_k, \\ x_1 &= x^* + e_1 = x^* + \sum_{k=2}^n \alpha_k \mathbf{u}_k \\ &\vdots \\ x_{n-1} &= x^* + e_{n-1} = x^* + \alpha_n \mathbf{u}_n \\ x_n &= x^*. \end{aligned}$$

After n steps, we get all errors removed since we traversed through all directions \mathbf{u}_k . One question left to be discussed is: Why is orthogonal basis not a good choice? The reason is that in order to compute coefficients α_k , we need A-conjugacy.

Suppose instead we have a set of orthogonal basis $\{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n\}$, $\mathbf{b}_i \cdot \mathbf{b}_j = 0$, for any $i \neq j$. Now, e_0 can be written as a linear combination of this set, i.e., $e_0 = \sum_{i=1}^n \beta_i \mathbf{b}_i$. To compute one coefficient β_i , we multiply \mathbf{b}_i^T on both sides:

$$\begin{aligned} \mathbf{b}_i^T e_0 &= \sum_{i=1}^n \beta_i \mathbf{b}_i^T \mathbf{b}_i, \\ \beta_i &= -\frac{\mathbf{b}_i^T e_0}{\mathbf{b}_i^T \mathbf{b}_i}. \end{aligned}$$

Here, the problem is that we have no idea what e_0 is. If we know it, $x_0 - e_0$ already gives the optimal solution x^* . Using A-conjugacy solves this problem.

Now, we can ask a new question. From which set of linear independent vectors should we generate $\{\mathbf{u}_k\}_{k=1}^n$? The conjugate gradient method suggests using gradients:

$$B_{n-1} = \{\nabla f(x_0), \nabla f(x_1), \dots, \nabla f(x_{n-1})\}.$$

The gradient can be written in several forms: for some $1 \leq i \leq n-1$,

$$\nabla f(x_i) = A x_i + \mathbf{b} = A x_i - A x^* = A(x_i - x^*) = A e_i. \quad (10)$$

We will explain how to generate the A-conjugate basis from B_{n-1} via (15). We first explain the mechanism of this method.

At each step, we choose one direction and minimize f along that direction. The iteration formula for x_i is $x_{i+1} = x_i + \beta_{i+1}\mathbf{u}_{i+1}$. To minimize, we set

$$\begin{aligned} \frac{d}{d\beta_{i+1}} f(x_{i+1}) &= \nabla f(x_{i+1})^T \mathbf{u}_{i+1} \\ &= \mathbf{u}_{i+1}^T \nabla f(x_i + \beta_{i+1}\mathbf{u}_{i+1}) \\ &= \mathbf{u}_{i+1}^T (A(x_i + \beta_{i+1}\mathbf{u}_{i+1}) - x^*) \\ &= \mathbf{u}_{i+1}^T \nabla f(x_i) + \beta_{i+1} \mathbf{u}_{i+1}^T A \mathbf{u}_{i+1} \\ &= 0. \end{aligned} \tag{11}$$

$$\beta_{i+1} = -\frac{\mathbf{u}_{i+1}^T \nabla f(x_i)}{\mathbf{u}_{i+1}^T A \mathbf{u}_{i+1}}. \tag{12}$$

Since $e_i = \sum_{k=1}^n \alpha_k \mathbf{u}_k$, we multiply $\mathbf{u}_{i+1}^T A$ on both side,

$$\begin{aligned} \mathbf{u}_{i+1}^T A e_i &= \sum_{k=1}^n \alpha_k \mathbf{u}_{i+1}^T A \mathbf{u}_k \\ \implies \mathbf{u}_{i+1}^T A e_i &= \alpha_{i+1} \mathbf{u}_{i+1}^T A \mathbf{u}_{i+1} \\ \implies \mathbf{u}_{i+1}^T A e_i &= \alpha_{i+1} \mathbf{u}_{i+1}^T A \mathbf{u}_{i+1} \\ \implies \alpha_{i+1} &= \frac{\mathbf{u}_{i+1}^T A e_i}{\mathbf{u}_{i+1}^T A \mathbf{u}_{i+1}} = \frac{\mathbf{u}_{i+1}^T \nabla f(x_i)}{\mathbf{u}_{i+1}^T A \mathbf{u}_{i+1}}. \end{aligned} \tag{13}$$

From (12) and (13), we see that $\alpha_{i+1} = -\beta_{i+1}$. This guarantees that once we minimize the function along one direction, we remove the error in that direction completely. Hence, if we set $\alpha_i = \frac{\mathbf{u}_i^T \nabla f(x_{i-1})}{\mathbf{u}_i^T A \mathbf{u}_i}$, we can finish the minimization process in at most n steps.

Some Properties

Claim 0.3.4.4. $\nabla f(x_i)$ is orthogonal to \mathbf{u}_j for all $j \leq i$, $i \in \{0, 1, \dots, n-1\}$, $j \in \{1, \dots, i\}$.

Proof. We know that $e_i = \sum_{k=i+1}^n \alpha_k \mathbf{u}_k$, $i \in \{0, 1, \dots, n-1\}$. We multiply $\mathbf{u}_j^T A$, for some $j \leq i$, on both sides,

$$\mathbf{u}_j^T \nabla f(x_i) = \mathbf{u}_j^T A e_i = \sum_{k=i+1}^n \alpha_k \mathbf{u}_j^T A \mathbf{u}_k = 0.$$

□

Claim 0.3.4.5. Let $B_{i-1} = \{\nabla f(x_0), \nabla f(x_1), \dots, \nabla f(x_{i-1})\}$, $D_i = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_i\}$, for any $i \in \{1, 2, \dots, n\}$. Then, $\text{span}\{B_{i-1}\} = \text{span}\{D_i\}$

Proof. First, we prove that $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_i\} \subseteq \text{span}\{\nabla f(x_0), \nabla f(x_1), \dots, \nabla f(x_{i-1})\}$.

$$\mathbf{u}_1 = \nabla f(x_0) \in \text{span}\{B_0\},$$

$$\mathbf{u}_2 = \nabla f(x_1) + \gamma_1 \mathbf{u}_1 = \nabla f(x_1) + \gamma_1 \nabla f(x_0) \in \text{span}\{B_1\},$$

$$\begin{aligned} & \vdots \\ \mathbf{u}_i &= \nabla f(x_{i-1}) + \sum_{k=1}^{i-1} \gamma_k \mathbf{u}_k \in \text{span}\{B_{i-1}\}. \end{aligned}$$

By definition, \mathbf{u}_i is a linear combination of $\nabla f(x_i)$ and $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{i-1}$. By recursion each $\mathbf{u}_i \in \text{span}\{B_{i-1}\}$. Hence, forward direction is proved.

Next, we prove that $\{\nabla f(x_0), \nabla f(x_1), \dots, \nabla f(x_{i-1})\} \subseteq \text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_i\}$. This direction is obvious if we move some terms and express $\nabla f(x_i)$ as a linear combination of \mathbf{u}_i :

$$\begin{aligned} \nabla f(x_0) &= \mathbf{u}_1 \in \text{span}\{\mathbf{u}_1\}, \\ \nabla f(x_1) &= \mathbf{u}_2 - \gamma_1 \mathbf{u}_1 \in \text{span}\{\mathbf{u}_1, \mathbf{u}_2\}, \\ & \vdots \\ \nabla f(x_{i-1}) &= \mathbf{u}_i - \sum_{k=1}^{i-1} \gamma_k \mathbf{u}_k \in \text{span}\{\mathbf{u}_1, \dots, \mathbf{u}_i\}. \end{aligned}$$

□

From claim 0.3.4.4 and claim 0.3.4.5, we know that $\nabla f(x_i)$ is orthogonal to B_{i-1} .

Claim 0.3.4.6. *Let $E_{i-1} = \{\nabla f(x_0), A\nabla f(x_0), \dots, A^{i-1}\nabla f(x_0)\}$. Then $\text{span}\{E_{i-1}\} = \text{span}\{B_{i-1}\}$, for any $i \in \{1, 2, \dots, n\}$.*

Proof. From (15), we get that

$$\begin{aligned} \nabla f(x_i) &= A(x_i - x^*) \\ &= A(x_{i-1} + \beta_i \mathbf{u}_i - x^*) \\ &= \nabla f(x_{i-1}) + \beta_i A \mathbf{u}_i \end{aligned} \tag{14}$$

From (0.3.4.5), we know that

$$\nabla f(x_i) = \nabla f(x_{i-1}) + \beta_i A \sum_{k=0}^{i-1} \nabla f(x_k)$$

We just need to prove the base cases:

$$\nabla f(x_0) = \mathbf{u}_1 \in \text{span}\{E_0\},$$

$$\nabla f(x_1) = \nabla f(x_0) + \beta_1 A \mathbf{u}_1 = \nabla f(x_0) + \beta_1 A \nabla f(x_0) \in \text{span}\{E_1\}.$$

Using recursion, finishes the proof. □

From claim 0.3.4.5 and claim 0.3.4.6, we get the result that $\text{span}\{E_{i-1}\} = \text{span}\{D_i\}$ and $\text{span}\{A * D_i\} \subseteq \text{span}\{E_i\} = \text{span}\{B_i\}$. Therefore, $\nabla f(x_{i+1})$ is orthogonal to $\text{span}\{A * D_i\}$ for $i \in \{1, 2, \dots, n-2\}$.

Next, we have proved a key fact that will be very helpful when we generate $\{\mathbf{u}_k\}_{k=1}^n$. By using claim 0.3.4.3 and claim 0.3.4.6, we know that $\nabla f(x_k)^T A \mathbf{u}_m = 0$ for all $m \in \{1, 2, \dots, k-1\}$. Therefore, we can determine all coefficients in Gram-Schmidt Algorithm:

$$\begin{aligned}\gamma_1 &= \gamma_2 = \dots = \gamma_{k-1} = 0, \\ \gamma_k &= -\frac{\nabla f(x_k)^T A \mathbf{u}_k}{\mathbf{u}_k^T A \mathbf{u}_k}, \\ \mathbf{u}_{k+1} &= \nabla f(x_k) + \gamma_k \mathbf{u}_k.\end{aligned}\tag{15}$$

Now we present the algorithm of conjugate gradient method. The following algorithm is written with respect to the function $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A \mathbf{x} + b \mathbf{x} + c$.

Algorithm 4: Conjugate Gradient Method

```
conjugateGradient (A, b, x0, max_iter);
Input : A, b, x0, max_iter
Output: x*
β = 0;
γ = 0;
count = 0;
n = size(A);
x = x0;
∇f(x) = Ax + b;
u = Ax + b;
while count ≤ n - 1 and count ≤ max_iter do
    count = count + 1;
    β = - $\frac{\mathbf{u}^T \nabla f(x)}{\mathbf{u}^T A \mathbf{u}}$ ;
    x = x + βu;
    ∇f(x) = Ax + b;
    γ = - $\frac{\nabla f(x)^T A \mathbf{u}}{\mathbf{u}^T A \mathbf{u}}$ ;
    u = ∇f(x) + γu;
return x;
```

Simplification of the Algorithm

We can simplify the coefficients β and γ by modifying the following equations:

$$\begin{aligned}x_{i+1} &= x_i + \beta_{i+1} \mathbf{u}_{i+1} \\ Ax_{i+1} + b &= Ax_i + b + \beta_{i+1} A \mathbf{u}_{i+1} \\ \nabla f(x_{i+1}) &= \nabla f(x_i) + \beta_{i+1} A \mathbf{u}_{i+1}\end{aligned}$$

Multiplying $\nabla f(x_{i+1})$ and $\nabla f(x_i)$ to the equation yields,

$$\begin{aligned}\|\nabla f(x_{i+1})\|^2 &= \beta_{i+1} \nabla f(x_{i+1})^T A \mathbf{u}_{i+1} \\ \|\nabla f(x_i)\|^2 &= -\beta_{i+1} \nabla f(x_i)^T A \mathbf{u}_{i+1}\end{aligned}\tag{16}$$

From (15), we know that $\nabla f(x_i) = \mathbf{u}_{i+1} + \gamma_i \mathbf{u}_i$. Since \mathbf{u}_{i+1} is A-conjugate to \mathbf{u}_i , we combine (15) and (16):

$$\|\nabla f(x_i)\|^2 = -\beta_{i+1} \mathbf{u}_{i+1}^T A \mathbf{u}_{i+1}.$$

The variable γ used in the algorithm can be computed by:

$$\gamma = \frac{\|\nabla f(x_i)\|^2}{\|\nabla f(x_{i-1})\|^2}.$$

From claim 0.3.4.4 we know that $\nabla f(x_{i-1})$ is orthogonal to \mathbf{u}_{i-1} and from (15), we get $\mathbf{u}_i = \nabla f(x_{i-1}) - \gamma_i \mathbf{u}_{i-1}$. Then the variable β used in the algorithm can be computed in the following way:

$$\beta = -\frac{\nabla f(x_{i-1})^T (\nabla f(x_{i-1}) - \gamma_i \mathbf{u}_{i-1})}{\mathbf{u}_i^T A \mathbf{u}_i} = -\frac{\|\nabla f(x_{i-1})\|^2}{\mathbf{u}_i^T A \mathbf{u}_i}.$$

Notice that we can not change the denominator into $\|\nabla f(x_i)\|^2$ because by the time we compute β , \mathbf{x}_i has not been updated yet. i.e., we still do not know the value of $\|\nabla f(x_i)\|^2$.

Hence, we can simplify the above algorithm by replacing β, γ by what we just computed. Since all the setup of function and variables remain the same, we only present the new version of the while loop:

Algorithm 5: Conjugate Gradient Method's While Loop

while $count \leq n - 1$ and $count \leq max_iter$ **do**

$count = count + 1;$
 $\beta = -\frac{\|\nabla f(x)\|^2}{\mathbf{u}^T A \mathbf{u}};$
 $x = x + \beta \mathbf{u};$
 $\nabla f(x)_{pre} = \nabla f(x);$
 $\nabla f(x) = Ax + b;$
 $\gamma = \frac{\|\nabla f(x)\|^2}{\|\nabla f(x)_{pre}\|^2};$
 $\mathbf{u} = \nabla f(x) + \gamma \mathbf{u};$

return $x;$

Next, we give two example of applying conjugate gradient method.

Example 0.3.4.7. In the first example, we use conjugate gradient method to find the minimum point of a quadratic function with a symmetric positive definite matrix. Suppose $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A \mathbf{x} + b \mathbf{x} + c$, where $A = \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix}$, $b = \begin{bmatrix} -2 \\ 5 \end{bmatrix}$, $c = 3$. $\nabla f(x) = A\mathbf{x} + b$. The optimal point x^* that gives the minimum value of f is $x^* = \begin{bmatrix} -11 \\ 7 \\ 19 \\ 14 \end{bmatrix}$. We choose a initial starting point $x = \begin{bmatrix} -3 \\ 5 \end{bmatrix}$.

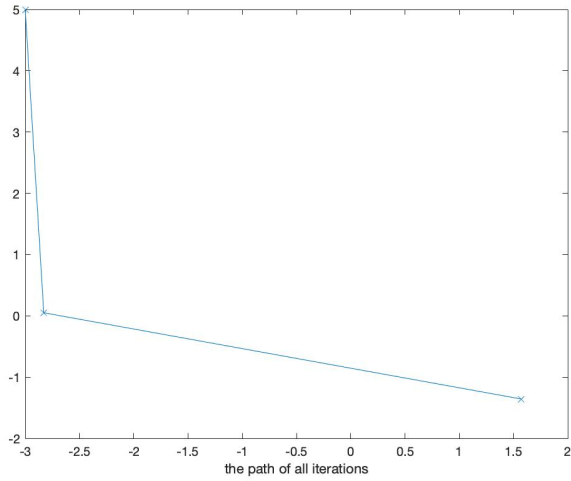


Figure 6: The conjugate gradient descent method

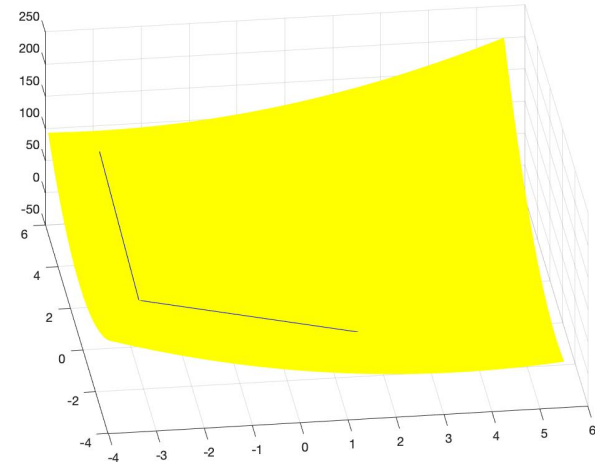


Figure 7: The conjugate gradient descent method search path

The graph shows that conjugate gradient method only takes 2 steps because the matrix A is a 2 matrix. Figure 6 shows a search path with 2 directions. If we imagine two unit vectors in these two directions, the Figure 6 shows the geometric explanation of A -conjugacy. Notice these two directions would be orthogonal to each other if $A = \mathbb{I}$.

Example 0.3.4.8. In the second example, we use a higher dimensional matrix A . The following graph shows an example of using conjugate gradient method to solve the minimum value of the function $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T A \mathbf{x} + b\mathbf{x} + c$.

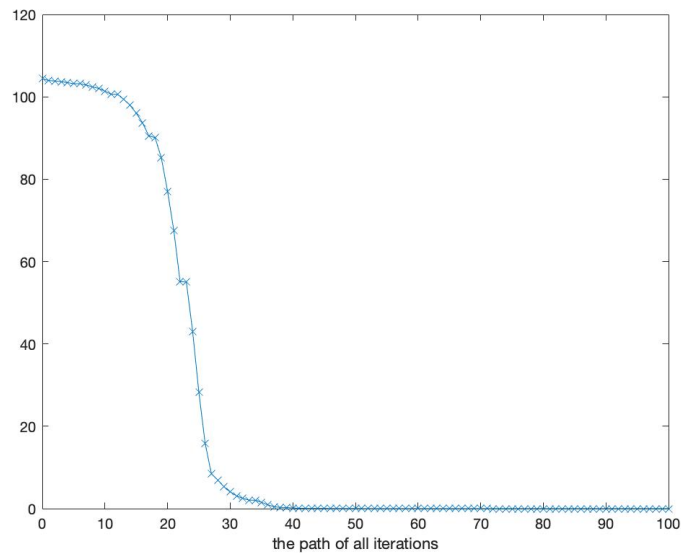


Figure 8: The conjugate gradient descent method in a higher dimension

Suppose we have an arbitrary 100×100 symmetric positive definite matrix A . We can construct

a random vector b with 100 entries and an arbitrary starting point \mathbf{x}_0 as a vector of 100 entries. Figure 9 is the graph of the change of e_i with respect to each iteration step \mathbf{x}_i . It shows that if we have a matrix with good eigenvalues, then conjugate gradient method takes much fewer steps than the maximum number of iterations, i.e, the number of dimension of the matrix.

Here is a brief explanation of how the above matrix is designed. We first create a random matrix A with size 100×100 . Create a new symmetric matrix B by adding A with A^T . Do eigendecomposition on B to get its eigenvalue matrix D and eigenvector matrix V . Apply absolute value to the eigenvalue matrix and add one, which guarantees that all eigenvalues are positive. We denote the new eigenvalue matrix as D' . Finally, we can get a random symmetric positive definite matrix by multiplying $A^T D' A$.

Algorithm 6: Create a random symmetric positive definite matrix

```

while  $count \leq n - 1$  and  $count \leq max\_iter$  do
     $A = rand(100, 100);$ 
     $B = A + A^T;$ 
     $[D, V] = eig(B);$ 
     $D' = abs(D) + 1;$ 
     $B' = A^T D' A;$ 

```

In summary, we make a conclusion on several advantages of conjugate gradient method:

- The update formula of the conjugate gradient method is simple.
- For quadratic functions, the conjugate gradient method always converges in a finite number of iterations.
- When the matrix has a very high dimension n but its eigenvalues are clustered, then the method converges much faster than n steps.
- The algorithm need no storage of the matrix A , hence, memory efficient.

0.4 Newton's Method

From the previous two sections, we learned the Gradient Descent Method, which uses linear approximation or the gradient as the descent direction to minimize the function, and the Conjugate Gradient Method, which solves large systems of linear equations. In this section, we are going to present Newton's Method which minimizes the function by minimizing a quadratic approximation with the help of both the gradient and the Hessian matrix. We are going to introduce Newton's descent direction, Newton decrement, two phases convergence analysis and some examples. Then we are going to take a look at the application of Newton's method to a specific kind of function, the self-concordant function.

0.4.1 The descent direction

Suppose the objective function $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable and its Hessian matrix is positive definite. Notice that second order differentiability guarantees the symmetry of its Hessian matrix and positive definiteness guarantees the invertibility of the matrix. The second-order Taylor approximation of this function at x is

$$f(x + \epsilon) = f(x) + \nabla f(x)^T \epsilon + \frac{1}{2} \epsilon^T \nabla^2 f(x) \epsilon + \mathcal{O}(\|\epsilon\|^3), \epsilon > 0.$$

We assume the approximation is local and ϵ is not too big. Therefore, we can ignore the last term that contains a large power of ϵ . Since $\nabla^2 f(x)$ is positive definite, we can view the approximate function as a convex and quadratic function of the variable ϵ . To achieve its minimum, we choose $\epsilon = -\nabla^2 f(x)^{-1} \nabla f(x)$. The intuitive idea is that we try to minimize the approximate function that satisfies the first and second order derivative of the objective function and the minimizer is good enough to serve as the descent direction for the objective function. For future notation, we denote this quadratic approximation function as $\hat{f}(a+x) := f(a) + \nabla f(a)^T x + \frac{1}{2} x^T \nabla^2 f(a) x$, for some fixed point a .

Here is a picture illustrating the second order Taylor approximation of the objective function $f(x, y) = x^2 + y^2 + e^{x^2+y^2}$ at the origin. In this case, the minimum of the quadratic approximation function is the same as the minimum of the objective function.

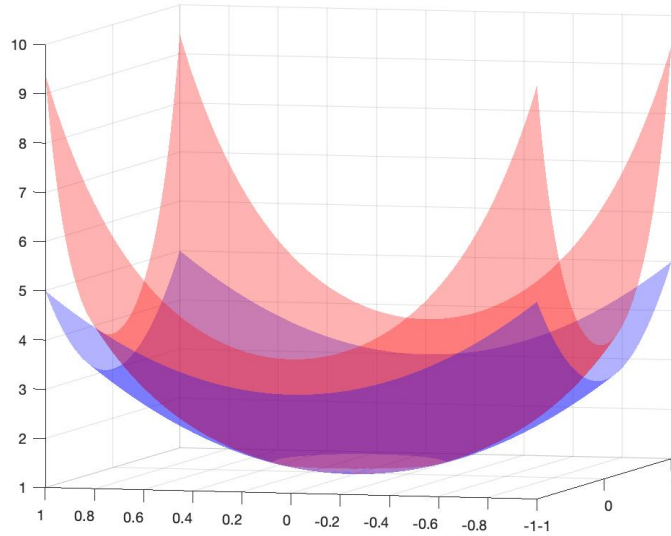


Figure 9: 2nd order Taylor Approximation

Another intuition for Newton's Method is the linear approximation of the gradient. The function reaches its minimum specifically when $\nabla f(x) = 0$. By linear approximation, $\nabla f(x + \epsilon) \approx \nabla f(x) + \nabla^2 f(x) \epsilon = 0$. This suggests the same idea: $\epsilon = -\nabla^2 f(x)^{-1} \nabla f(x)$.

The above explanation implies the idea of Newton's Method. The descent direction for Newton's Method is $\Delta x = -\nabla^2 f(x)^{-1} \nabla f(x)$. The iteration step is $x_{i+1} = x_i + t \Delta x_i$, where $\Delta x_i = -\nabla^2 f(x)^{-1} \nabla f(x)$. As before, t is a variable that determines the step size.

Theorem 0.4.1.1. Suppose we have a twice continuously differentiable function f whose Hessian matrix is positive definite. Assume we have x_i . By using Newton's Method Newton's descent direction, we can find $x_{i+1} = x_i - t\nabla^2 f(x)^{-1}\nabla f(x)$ with appropriate t such that $f(x_{i+1}) < f(x_i)$.

Proof. Denote $\mathbf{y} = -\nabla^2 f(x)^{-1}\nabla f(x)$. Define a function $g(t) = f(x_i + t\mathbf{y})$. The derivative of this function is $g'(t) = \nabla f(x_i + t\mathbf{y}) \cdot \mathbf{y}$. Plug in $t = 0$:

$$g'(0) = \nabla f(x_i) \cdot \mathbf{y} = -\nabla f(x_i)^T \nabla^2 f(x)^{-1} \nabla f(x).$$

Since $\nabla^2 f(x)$ is positive definite, we know that its inverse is also positive definite, which means that $g'(0) < 0$. Then in a small neighborhood $(-\epsilon, \epsilon)$, we know the function g is decreasing. There exists $t \in (-\epsilon, \epsilon)$ such that $g(t) < g(0)$, which is the same as saying $f(x_i) > f(x_i - t\nabla^2 f(x)^{-1}\nabla f(x))$. This guarantees that with a proper value of t , the Newton's iteration step gives a sequence of points with decreasing function value. \square

0.4.2 The Newton Decrement

The Newton decrement of f at x is defined to be

$$\lambda(x) = (\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x))^{\frac{1}{2}} = (-\nabla f(x)^T \Delta x)^{\frac{1}{2}}.$$

The $\lambda(x_i)$ is obtained by computing the difference between $f(x_i)$ and its quadratic approximation:

$$\begin{aligned} f(x_i) - \inf_{\epsilon} \hat{f}(x + \epsilon) &= f(x_i) - \hat{f}(x_i + \Delta x_i) \\ &= f(x_i) - (f(x_i) + \nabla f(x_i)^T \Delta x_i + \frac{1}{2} \Delta x_i^T \nabla^2 f(x_i) \Delta x_i) \\ &= -\nabla f(x_i)^T \Delta x_i - \frac{1}{2} \Delta x_i^T \nabla^2 f(x_i) \Delta x_i \\ &= \lambda(x_i)^2 - \frac{1}{2} \lambda(x_i)^2 \\ &= \frac{1}{2} \lambda(x_i)^2 \end{aligned} \tag{17}$$

Therefore, $\lambda(x)$ measures the difference between the objective function and the minimum of the quadratic approximation at each x . Even though the minimum of the quadratic approximation is not exactly the minimum of f , it is a good approximation. When x is close enough to the optimal point x^* , the difference between $\inf f$ and $\inf \hat{f}$ should be fairly small, which is just saying that $\lambda(x^*)$ is very small. Hence, $\lambda(x)$ is a good stopping criterion for Newton's Method. As $\lambda(x)$ gets small enough, we can say that x is very close to x^* .

On top of that, the Newton decrement will be used in the Backtracking Line Search. Recall that the while condition of the Backtracking Line Search is

$$f(x + t\Delta x) \leq f(x) + \alpha t \nabla f(x)^T \Delta x, \tag{18}$$

where $\nabla f(x)^T \Delta x = -\lambda(x)^2$.

However, why is the Exact Line Search not a good approach? Remember that for the Exact Line Search we have to compute t such that $\frac{d}{dt}f(x_i + t\Delta x_i) = 0$. Therefore, we need

$$\begin{aligned}\frac{d}{dt}f(x_i + t\Delta x_i) &= \nabla f(x_i + t\Delta x_i)^T \Delta x_i \\ &= -\nabla f(x_i + t\Delta x_i)^T \nabla^2 f(x)^{-1} \nabla f(x) = 0\end{aligned}$$

Since both $\nabla^2 f(x)^{-1}$ and $\nabla f(x)$ cannot be zero. Therefore, we have $\nabla f(x_i + t\Delta x_i) = 0$. Solving this equation exactly is difficult unless you have a nice formula for ∇f .

Next, we talk a little bit about the affine invariance of Newton descent direction.

Suppose $A \in \mathbb{R}^{n \times n}$ is a nonsingular matrix. Suppose we have a function $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$. Define a new function $f_1(y) = f(Ay)$, such that $x = Ay$. Then

$$\nabla f_1(y) = A^T \nabla f(x), \nabla^2 f_1(y) = A^T \nabla^2 f(x) A. \quad (19)$$

Then Newton's step for this new function f_1 is

$$\begin{aligned}\Delta y &= -\nabla^2 f_1(y)^{-1} \nabla f_1(y) \\ &= -(A^T \nabla^2 f(x) A)^{-1} A^T \nabla f(x) \\ &= -A^{-1} \nabla^2 f(x)^{-1} \nabla f(x) \\ &= A^{-1} \Delta x.\end{aligned} \quad (20)$$

Hence, $A\Delta y = \Delta x$. The descent directions of f_1 and f follow the same affine transformation. For Newton steps,

$$x + t\Delta x = A(y + t'\Delta y).$$

0.4.3 Convergence Analysis of Newton's Method

We assume a few conditions for the objective function f :

- f is twice continuously differentiable.
- $\nabla^2 f(x) \succeq m\mathbb{I}, m > 0$, the smallest eigenvalue of $\nabla^2 f(x)$ is at least m .
- $\nabla^2 f(x) \preceq M\mathbb{I}, M > 0$, the largest eigenvalue of $\nabla^2 f(x)$ is at most M .
- $\nabla^2 f(x)$ is a Lipschitz function with constant L , i.e., $\|\nabla^2 f(x) - \nabla^2 f(y)\|_2 \leq L\|x - y\|_2$ for any $x, y \in \text{dom}(f)$.

The Lipschitz condition provides a bound on the third derivative of f . When the function is quadratic, we have $\|\nabla^2 f(x) - \nabla^2 f(y)\|_2 = 0$ and L can be chosen to be 0. So when L is small or close to zero, the function can be well approximated by a quadratic function. As we will show in the next section, if f is a quadratic function, Newton's Method takes only one step to reach the optimal point x^* . Therefore, Newton's Method works well for functions that have small Lipschitz constant L .

Choose a number $\eta \in (0, 3(1 - 2\alpha)\frac{m^2}{L})$ where α is the coefficient used in the Backtracking Line Search as written in (18), m and L are defined above. Since α is an arbitrary number less than $\frac{1}{2}$,

we assume α is not too small, $\alpha \in (\frac{1}{3}, \frac{1}{2})$. We know that $3(1 - 2\alpha)\frac{m^2}{L}$ is greater than 0 because α is chosen to be less than $\frac{1}{2}$ and m, L are both positive.

There are two phases for the convergence of Newton's Method. The first phase is called the **Damped Newton Phase** where $\|\nabla f(x)\|_2 \geq \eta$. The second phase is called the **Quadratically Convergent Phase** where $\|\nabla f(x)\|_2 < \eta$.

Damped Newton Phase

In this phase, Newton's Method uses the Backtracking Line Search to determine the step size t . We are going to show that there exists a number γ such that each iteration step results in a decrease of at least γ in the objective function. We assume that $\|\nabla^2 f(x)\| \geq \eta$. By the Taylor's expansion theorem,

$$f(x + t\Delta x) = f(x) + t\nabla f(x)^T \Delta x + \frac{1}{2}t^2 \Delta x^T \nabla^2 f(y) \Delta x, \quad (21)$$

for some y between x and $x + t\Delta x$. Since $\nabla^2 f(x) \preceq M\mathbb{I}$, we get

$$f(x + t\Delta x) \leq f(x) + t\nabla f(x)^T \Delta x + \frac{1}{2}t^2 M \|\Delta x\|_2^2. \quad (22)$$

By definition, $\nabla f(x)^T \Delta x = -\lambda(x)^2$ and

$$\lambda(x)^2 = \Delta x^T \nabla^2 f(x) \Delta x \geq m \|\Delta x\|_2^2. \quad (23)$$

Therefore, (22) can be simplified into:

$$f(x + t\Delta x) \leq f(x) - t\lambda(x)^2 + \frac{t^2 M}{2} \frac{\lambda(x)^2}{m} \quad (24)$$

Since (24) holds for all t , we choose $\tilde{t} = \frac{m}{M}$ which is the minimizer of the right hand side.

$$\begin{aligned} f(x + \frac{m}{M}\Delta x) &\leq f(x) - \frac{m}{M}\lambda(x)^2 + \frac{m}{2M}\lambda(x)^2 \\ &= f(x) - \frac{m}{2M}\lambda(x)^2 \\ &= f(x) - \alpha \frac{m}{M}\lambda(x)^2 \end{aligned} \quad (25)$$

Thus, $\tilde{t} = \frac{m}{M}$ satisfies the exit condition of the while loop in the Backtracking Line Search. We know that the Backtracking Line Search returns a step size $t \geq \gamma \frac{m}{M}$. Recall that γ is a variable used to shrink t at each iteration.

Notice that by definition,

$$\lambda(x)^2 = \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) \geq \frac{1}{M} \|\nabla f(x)\|_2^2$$

Hence, plugging in $t \geq \gamma \frac{m}{M}$ and reorganizing (25) yield:

$$\begin{aligned} f(x_i) - f(x_{i+1}) &\geq \alpha \gamma \frac{m}{M} \lambda(x)^2 \\ &\geq \alpha \gamma \frac{m}{M^2} \|\nabla f(x)\|_2^2 \\ &\geq \alpha \gamma \frac{m}{M^2} \eta^2. \end{aligned} \quad (26)$$

Therefore, at each iteration step, the amount of decrease is at least $\alpha\gamma\frac{m}{M^2}\eta^2$. The total number of iterations will be bounded by

$$\frac{f(x_0) - f(x^*)}{\alpha\gamma\frac{m}{M^2}\eta^2},$$

where x_0 is the starting point and x^* is the optimal point.

Quadratically Convergent Phase

In the Quadratically Convergent Phase, Newton's Method needs no Backtracking Line Search because we are going to prove that $t = 1$ satisfies the quit condition for the while loop. More importantly, in this phase, f decreases quadratically at each step.

Remember we have two assumptions: $\eta \leq 3(1 - 2\alpha)\frac{m^2}{L}$ and $\|\nabla f(x)\|_2 < \eta$. We have not used the Lipschitz condition in the Damped Newton Phase and we will use it here.

By the Lipschitz condition,

$$\|\nabla^2 f(x + t\Delta x) - \nabla^2 f(x)\|_2 \leq L \|\Delta x\|_2. \quad (27)$$

Define a new function $\tilde{f}(t) = f(x + t\Delta x)$. Then $\tilde{f}''(t) = \Delta x^T \nabla^2 f(x + t\Delta x) \Delta x$. This suggests that we can multiply Δx on both sides of (27) and get an inequality involving the second derivative of \tilde{f} :

$$\begin{aligned} |\Delta x^T [\nabla^2 f(x + t\Delta x) - \nabla^2 f(x)] \Delta x| &\leq tL \|\Delta x\|_2^3 \\ |\Delta x^T \nabla^2 f(x + t\Delta x) \Delta x - \Delta x^T \nabla^2 f(x) \Delta x| &\leq tL \|\Delta x\|_2^3 \\ |\tilde{f}''(t) - \tilde{f}''(0)| &\leq tL \|\Delta x\|_2^3 \\ |\tilde{f}''(t) - \tilde{f}''(0)| &\leq tL \frac{\lambda(x)^3}{m^{\frac{3}{2}}}, \end{aligned} \quad (28)$$

where the last inequality follows by (23).

We can get an upper bound on \tilde{f} by computing the integral twice. We know that

$$\tilde{f}''(0) = \Delta x^T \nabla^2 f(x) \Delta x = \lambda(x)^2,$$

and

$$\tilde{f}'(0) = \nabla f(x)^T \Delta x = -\lambda(x)^2, \tilde{f}(0) = f(x),$$

which are conditions for determining constants in the indefinite integral. Then with the help of (23) and above conditions, we integrate (28) one time and get:

$$\begin{aligned} \tilde{f}''(t) &\leq \tilde{f}''(0) + tL \|\Delta x\|_2^3 \\ \tilde{f}'(t) &\leq t\tilde{f}''(0) + \frac{t^2 L}{2m^{\frac{3}{2}}} \lambda(x)^3 - \lambda(x)^2, \end{aligned} \quad (29)$$

where $\tilde{f}'(0) = -\lambda(x)^2$ and $m^{\frac{3}{2}}$ comes from $\lambda(x)^2 \geq m\|\Delta x\|_2^2$.

$$\begin{aligned} \tilde{f}(t) &\leq \frac{t^2}{2} \lambda(x)^2 + \frac{t^3 L}{6m^{\frac{3}{2}}} \lambda(x)^3 - \lambda(x)^2 t + f(x) \\ &= -\frac{t^2}{2} \lambda(x)^2 + \frac{t^3 L}{6m^{\frac{3}{2}}} \lambda(x)^3 + f(x) \\ &= -\left(\frac{t^2}{2} - \frac{t^3 L \lambda(x)}{6m^{\frac{3}{2}}}\right) \lambda(x)^2 + f(x). \end{aligned} \quad (30)$$

Notice the above inequality looks similar to the while condition in the Backtracking Line Search except that the coefficient is more complicated than α . Now, we need to use our two assumptions.

$$\lambda(x) \leq \frac{1}{\sqrt{m}} \|\nabla f(x)\|_2 \leq 3(1 - 2\alpha) \frac{m^{\frac{3}{2}}}{L},$$

and

$$\alpha \leq \frac{1}{2} - \frac{L\lambda(x)}{6m^{\frac{3}{2}}}.$$

Therefore, when $t = 1$,

$$\begin{aligned} f(x + \Delta x) &\leq -\left(\frac{1}{2} - \frac{L\lambda(x)}{6m^{\frac{3}{2}}}\right)\lambda(x)^2 + f(x) \\ &\leq -\alpha\lambda(x)^2 + f(x). \end{aligned}$$

This means that the step size $t = 1$ satisfies the Backtracking Line Search exit condition. Next, we prove that the rate of convergence is quadratic. Assume $t = 1$.

Theorem 0.4.3.1. *If g is differentiable and ∇g is a Lipschitz function with constant L , i.e., $\|\nabla g(x) - \nabla g(y)\|_2 \leq L \|x - y\|_2$. Then*

$$\|g(x + y) - g(x) - \nabla g(x)^T y\|_2 \leq \frac{L}{2} \|y\|_2^2.$$

Proof. Change the left hand side into an integral,

$$\begin{aligned} g(x + y) - g(x) - \nabla g(x)^T y &= \int_0^1 (\nabla g(x + ty) - \nabla g(x)) y \, dt \\ &\leq \int_0^1 L \|y\|_2^2 t \, dt = \frac{L}{2} \|y\|_2^2. \end{aligned}$$

□

By using the above theorem with $g = \nabla f$, $y = \Delta x$, we get

$$\begin{aligned} \|\nabla f(x + \Delta x)\|_2 &= \|\nabla f(x + \Delta x) - \nabla f(x) - \nabla^2 f(x)^T \Delta x\|_2 \\ &\leq \frac{L}{2} \|\Delta x\|_2^2 \\ &\leq \frac{L}{2m^2} \|\nabla f(x)\|_2^2. \end{aligned} \tag{31}$$

Therefore, we get the desired result, $\|\nabla f(x_{i+1})\|_2 \leq \frac{L}{2m^2} \|\nabla f(x_i)\|_2^2$.

When $\|\nabla f(x_i)\|_2 < \eta$,

$$\|\nabla f(x_{i+1})\|_2 \leq \frac{L}{2m^2} \eta^2 \leq 3(1 - 2\alpha) \frac{m^2}{L} \frac{L}{2m^2} \eta \leq \frac{1}{2} \eta,$$

where the last inequality comes from the fact that $\alpha \in (\frac{1}{3}, \frac{1}{2})$ implies $3(1 - 2\alpha) < 1$.

Therefore, once the condition $\|\nabla f(x_i)\|_2 < \eta$ holds for some x_i , it is going to be true for all steps afterwards. This means that after x_i the method moves into Quadratically convergent phase and $t = 1$ is always satisfied.

Recall from the Gradient Descent Method, we get an inequality (4):

$$\|\nabla f(x)\|^2 \geq 2m(f(x) - f(x^*)).$$

Here, we can reuse this inequality because the objective function f satisfies all the conditions to derive this inequality. Readers can refer to the previous section to see the proof.

Suppose that starting from x_i the method gets into the quadratically convergent phase. Then by applying (31) recursively, we get

$$\begin{aligned} f(x_k) - f(x^*) &\leq \frac{1}{2m} \|\nabla f(x_k)\|_2^2 \\ \|\nabla f(x_k)\|_2 &\leq \frac{L}{2m^2} \|\nabla f(x_{k-1})\|_2^2 \leq \left(\frac{L}{2m^2} \|\nabla f(x_i)\|_2\right)^{2^{k-i}} \end{aligned} \quad (32)$$

We can simplify this even further:

$$\begin{aligned} \frac{L}{2m^2} \|\nabla f(x_i)\| &\leq \frac{L}{2m^2} 3(1 - 2\alpha) \frac{m^2}{L} \leq \frac{1}{2}, \\ \|\nabla f(x_k)\|_2^2 &\leq \left(\frac{1}{2}\right)^{2^{k-i}}. \end{aligned}$$

Therefore, we finally get

$$f(x_k) - f(x^*) \leq \frac{1}{2m} \left(\frac{1}{2}\right)^{2^{k-i}}.$$

Suppose we want to obtain x_k such that $f(x_k) - f(x^*) \leq \epsilon$. We just need to make the upper bound less than ϵ :

$$\begin{aligned} \frac{1}{2m} \left(\frac{1}{2}\right)^{2^{k-i}} &\leq \epsilon \\ 2^{k-i} &\leq \log_2 (2m\epsilon)^{-1} \\ k - i &\leq \log_2 \log_2 (2m\epsilon)^{-1}. \end{aligned}$$

In practice, we only need 5 or 6 steps because $\left(\frac{1}{2}\right)^{2^6}$ is already $5.42 * 10^{-20}$. Therefore, combining two phases, the total number of iterations are bounded by

$$\frac{f(x_0) - f(x^*)}{\alpha \gamma \frac{m}{M^2} \eta^2} + 6.$$

0.4.4 Algorithm

The following algorithm uses the Backtracking Line Search to determine each iteration step size.

Algorithm 7: Newton's Method

NewtonMethod ($f, x_0, \epsilon, \alpha, \gamma$);

Compute the initial λ^2 and Δx for x_0 ;

while $\frac{\lambda^2}{2} \geq \epsilon$ **do**

$x_{i+1} = x_i + \Delta x$;

while $f(x_i) - f(x_{i+1}) < -\alpha t \nabla f(x_i)^T \Delta x$, **do**

$t = \gamma t$;

$x_{i+1} = x_i + t \Delta x$;

$x_i = x_{i+1}$;

$\Delta x_i = -\nabla^2 f(x_i)^{-1} \nabla f(x_i)$;

$\lambda^2 = \nabla f(x_i)^T \nabla^2 f(x_i)^{-1} \nabla f(x_i)$;

return x_i ;

0.4.5 Examples

Quadratic function

Consider a quadratic function $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A \mathbf{x} + b \mathbf{x} + c$ where A is symmetric and invertible. We know that $\Delta x = -A^{-1}(Ax + b) = -x - A^{-1}b$ and $\nabla f(x) = Ax + b$. Thus,

$$\begin{aligned} x_1 &= x_0 + \Delta x_0 \\ &= x_0 - x_0 - A^{-1}b \\ &= -A^{-1}b, \end{aligned}$$

which gives the minimum value for f . Hence, for any quadratic function, Newton's Method converges in one step.

Exponential Function in \mathbb{R}^2

Consider a convex function

$$f(x_1, x_2) = e^{x_1+3x_2-0.1} + e^{x_1-3x_2-0.1} + e^{-x_1-0.1}.$$

We can compute its gradient and Hessian matrix directly:

$$\nabla f = e^{x_1+3x_2-0.1} \begin{bmatrix} 1 \\ 3 \end{bmatrix} + e^{x_1-3x_2-0.1} \begin{bmatrix} 1 \\ -3 \end{bmatrix} + e^{-x_1-0.1} \begin{bmatrix} -1 \\ 0 \end{bmatrix}$$

and

$$\nabla^2 f = e^{x_1+3x_2-0.1} \begin{bmatrix} 1 & 3 \\ 3 & 9 \end{bmatrix} + e^{x_1-3x_2-0.1} \begin{bmatrix} 1 & -3 \\ -3 & 9 \end{bmatrix} + e^{-x_1-0.1} \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}.$$

Therefore, we know Δx and λ . By using Newton's Method with a starting point $x_0 = \begin{bmatrix} -5 \\ -5 \end{bmatrix}$, we get the following two graphs:

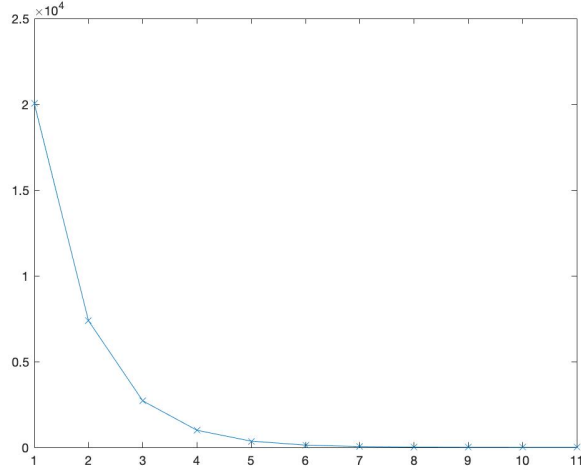


Figure 10: Newton's Method Search Path

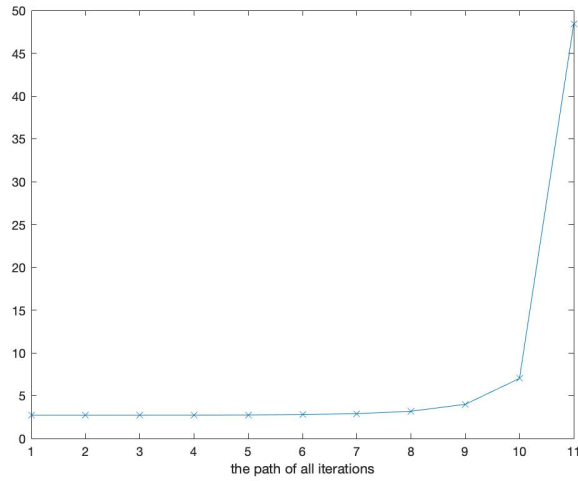


Figure 11: Damped Phase and Quadratically Convergent Phase

Figure 10 shows how $f(x_i) - f(x^*)$ decreases at each step. Here we see that Newton's Method gives a convergent sequence of points which leads to x^* . In Figure 11, the vertical axis represents $g(i) = \frac{f(x_i) - f(x^*)}{f(x_{i+1}) - f(x^*)}$. Once the algorithm gets into the quadratically convergent phase (starting from index 9), its error decreases quadratically, from approximately 3 to 9 to 49. This means that $(f(x_i) - f(x^*))^2 \approx f(x_{i+1}) - f(x^*)$ and $g(i)^2 \approx g(i+1)$. From the graph, we see that at step $i = 10$ and $i = 11$, $g(10)^2 \approx g(11)$. This matches our convergence analysis.

In summary, once the iteration step moves into the Quadratically Convergent Phase, Newton's Method converges rapidly within 5 or 6 steps with high precision. However, storing the Hessian matrix for each iteration is space inefficient and the cost of matrix and vector multiplication is high if the function has large dimensions.

In practice, we usually do not have precise estimates for constants, m, M, L , which are used to

set up a range for η and separate two convergence phases. However, that is a theoretically correct analysis. Next, we are going to introduce a new type of function, called self-concordant function, whose convergence analysis is independent from these constants.

0.4.6 Newton's Method for Self-concordant Functions

Definition 0.4.6.1. A function $f(x) : \mathbb{R} \rightarrow \mathbb{R}$ that has a third derivative and satisfies for all $x \in \text{dom}(f)$,

$$|f'''(x)| \leq 2f''(x)^{\frac{3}{2}}$$

is called a self-concordant function.

From the definition, a self-concordant function f is convex because $f''(x) \geq 0$.

Affine Invariant Property

Theorem 0.4.6.2. Suppose we have a self-concordant function $f(x) : \mathbb{R} \rightarrow \mathbb{R}$. Define a new function $\tilde{f}(x) = f(ax + b)$ for some constant $a \neq 0$ and b . Then $\tilde{f}(x)$ is also self-concordant.

Proof. Since \tilde{f} is an affine transformation of f , \tilde{f} is also convex and has its third derivative. $\tilde{f}''' = a^3 f'''(ax + b)$ and $\tilde{f}'' = a^2 f''(ax + b)$. Since f is self-concordant and $ax + b \in \text{dom}(f)$, we get

$$|a^3 f'''(ax + b)| \leq 2(a^2 f''(x))^{\frac{3}{2}},$$

$$|\tilde{f}'''(x)| \leq 2\tilde{f}''(x)^{\frac{3}{2}}.$$

Therefore, \tilde{f} is also self-concordant. □

Theorem 0.4.6.3. Self-concordance is preserved under scalar multiplication if the constant $c \geq 1$ and addition.

Proof. The first part is straightforward. If $c \geq 1$, then $c^{\frac{3}{2}} \geq c$. So $c|f'''(x)| \leq 2(cf''(x))^{\frac{3}{2}}$. Next, suppose we have two self-concordant functions $f_1, f_2 : \mathbb{R} \rightarrow \mathbb{R}$. By definition and triangle inequality, we have

$$\begin{aligned} |f_1'''(x) + f_2'''(x)| &\leq |f_1'''(x)| + |f_2'''(x)| \\ &\leq 2(f_1''(x))^{\frac{3}{2}} + 2(f_2''(x))^{\frac{3}{2}} \\ &\leq 2(f_1''(x) + f_2''(x))^{\frac{3}{2}}. \end{aligned}$$

□

Definition 0.4.6.4. A function $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ is self-concordant if it is self-concordant along every line in the domain, i.e., $\hat{f}(t) = f(\mathbf{x} + t\mathbf{v})$ is a self-concordant function of t for all directions \mathbf{v} and for all $\mathbf{x} \in \text{dom}(f)$. In other words, $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfies

$$\left| \frac{d}{dt} \nabla^2 f(\mathbf{x} + t\mathbf{v}) \right|_{t=0} \preceq 2\sqrt{\mathbf{v}^T \nabla^2 f(\mathbf{x}) \mathbf{v}} \nabla^2 f(\mathbf{x}).$$

In higher dimension cases, the affine invariant property becomes the following:

Suppose we have a self-concordant function $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$. Let $A \in \mathbb{R}^{n \times m}$, $\mathbf{b} \in \mathbb{R}^n$. Then $f(A\mathbf{x} + \mathbf{b})$ is also self-concordant.

Example 0.4.6.5. Consider the negative logarithm.

Let $f(x) = -\log(x)$, $\text{dom}(f) = \{x > 0\}$. Second derivative $f''(x) = \frac{1}{x^2}$. Third derivative $f'''(x) = -\frac{1}{x^3}$. Obviously, $f'''(x) \leq 2f''(x)$. From here, we know that functions $f(x) = \sum_{k=1}^n -\log(b_k - a_k x)$ are all self-concordant because self-concordance is preserved under affine transformation and addition. We can then generalize to functions of higher dimensions $f(\mathbf{x}) = \sum_{k=1}^n -\log(\mathbf{b}_k - \mathbf{a}_k^T \mathbf{x})$.

Example 0.4.6.6. Consider another function $f(x) = x \log(x) - \log(x)$.

The domain of the function is $\text{dom}(f) = \{x > 0\}$. Second derivative $f''(x) = \frac{x+1}{x^2} > 0$. Third derivative $f'''(x) = -\frac{x+2}{x^3}$. Then

$$\left| \frac{f'''(x)}{2f''(x)^{\frac{3}{2}}} \right| = \frac{x+2}{x^3} * \frac{x^3}{2(x+1)^{\frac{3}{2}}} = \frac{x+2}{2(x+1)^{\frac{3}{2}}} = \frac{1}{2} \left(\frac{1}{(x+1)^{\frac{1}{2}}} + \frac{1}{(x+1)^{\frac{3}{2}}} \right),$$

which reaches its maximum at $x = 0$. Therefore,

$$\left| \frac{f'''(x)}{2f''(x)^{\frac{3}{2}}} \right| = 1$$

and f is self-concordant.

We know that if the Hessian matrix of a function f is positive definite, then the function is strictly convex. However, the converse is not true. For example, $f(x) = x^4$. The second derivative at 0 is $f''(0) = 0$, which is not positive definite. It can be proved that the Hessian matrix of a strictly convex self-concordant function is positive definite everywhere. Readers who are interested can refer to ***.

Recall from the previous section that the Newton decrement is

$$\lambda(x) = (\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x))^{\frac{1}{2}}$$

Here, we are going to prove another way to express the Newton decrement.

Theorem 0.4.6.7. Let v be a descent direction. From (1), we know that v satisfies $v^T \nabla f \leq 0$. The Newton decrement can be written as

$$\lambda(x) = \sup_{v \neq 0} \frac{-v^T \nabla f(x)}{(v^T \nabla^2 f(x) v)^{\frac{1}{2}}}.$$

Proof. We know that $\nabla^2 f(x)$ is positive definite. Define $w = (\nabla^2 f(x))^{\frac{1}{2}} v$, then $v = (\nabla^2 f(x))^{-\frac{1}{2}} w$ and $(\|w\|_2)^2 = v^T \nabla^2 f(x) v$. Then

$$\begin{aligned} \sup_{v^T \nabla^2 f(x) v = 1} -v^T \nabla f(x) &= \sup_{\|w\|_2 = 1} -w^T (\nabla^2 f(x))^{-\frac{1}{2}} \nabla f(x) \\ &= \left\| (\nabla^2 f(x))^{-\frac{1}{2}} \nabla f(x) \right\|_2 \\ &= (\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x))^{\frac{1}{2}} \\ &= \lambda(x). \end{aligned}$$

The second line comes from

$$w = \frac{-(\nabla^2 f(x))^{-\frac{1}{2}} \nabla f(x)}{\left\| (\nabla^2 f(x))^{-\frac{1}{2}} \nabla f(x) \right\|_2}.$$

Since w can be any unit vector, to get the supremum of the right hand side, we let w point in the direction of $\nabla^2 f(x)^{-\frac{1}{2}} \nabla f(x)$. Therefore,

$$\lambda(x) = \sup_{v^T \nabla^2 f(x) v = 1} -v^T \nabla f(x) = \sup_{v \neq 0} \frac{-v^T \nabla f(x)}{(v^T \nabla^2 f(x) v)^{\frac{1}{2}}}.$$

The last equality is obvious because $v^T \nabla^2 f(x) v = 1$. □

From the previous theorem, we immediately get an inequality

$$\lambda(x) \geq \frac{v^T \nabla f(x)}{(v^T \nabla^2 f(x) v)^{\frac{1}{2}}} \quad (33)$$

since $\lambda(x)$ is the supremum. The equality is obtained when $v = -\Delta x = \nabla^2 f(x)^{-1} \nabla f(x)$.

Theorem 0.4.6.8. *Suppose f is a strictly convex self-concordant function. Then the self-concordance inequality can be rewritten as*

$$\left| \frac{d}{dt} (f''(t))^{-\frac{1}{2}} \right| \leq 1,$$

for all $t \in \text{dom}(f)$.

Proof. The proof is very easy. The main point of this proof is to develop the upper and lower bounds on $f''(t)$.

$$\left| \frac{d}{dt} (f''(t))^{-\frac{1}{2}} \right| = \left| -\frac{1}{2} (f''(t))^{-\frac{3}{2}} f'''(t) \right| \leq 1,$$

which is just another way of saying $|f'''(x)| \leq 2f''(x)^{\frac{3}{2}}$. Assume $t \geq 0$ and the interval $[0, t]$ is contained in $\text{dom}(f)$. Then we can integrate the derivative between 0 and t :

$$\int_0^t \frac{d}{dx} (f''(x))^{-\frac{1}{2}} dx = f''(t)^{-\frac{1}{2}} - f''(0)^{-\frac{1}{2}} \in [-t, t].$$

Hence, we get

$$-t \leq f''(t)^{-\frac{1}{2}} - f''(0)^{-\frac{1}{2}} \leq t$$

We can isolate $f''(t)$ to get:

$$\begin{aligned} f''(0)^{-\frac{1}{2}} - t &\leq f''(t)^{-\frac{1}{2}} \leq t + f''(0)^{-\frac{1}{2}}, \\ \frac{f''(0)}{(1 + t\sqrt{f''(0)})^2} &\leq f''(t) \leq \frac{f''(0)}{(1 - t\sqrt{f''(0)})^2}. \end{aligned} \quad (34)$$

The right hand side inequality is valid when $f''(0)^{-\frac{1}{2}} - t \geq 0$, that is $0 \leq t \leq f''(0)^{-\frac{1}{2}}$. □

Bound on $f(x) - f(x^*)$

Assume the function $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ is a strictly convex self-concordant function and v is a descent direction but does not need to be the Newton direction. Let $\tilde{f}(t) : \mathbb{R} \rightarrow \mathbb{R}$ be $\tilde{f}(t) = f(x + tv)$, which is also strictly convex and self-concordant.

By using the lower bound of (34), we can integrate to get that

$$\begin{aligned}
\tilde{f}'(t) - \tilde{f}'(0) &= \int_0^t \tilde{f}''(x) dx \\
&\geq \int_0^t \frac{\tilde{f}''(0)}{(1 + x\sqrt{\tilde{f}''(0)})^2} dx \\
&= -\frac{\sqrt{\tilde{f}''(0)}}{1 + x\sqrt{\tilde{f}''(0)}} \Big|_0^t \\
&= \sqrt{\tilde{f}''(0)} - \frac{\sqrt{\tilde{f}''(0)}}{1 + t\sqrt{\tilde{f}''(0)}}. \\
\tilde{f}'(t) &\geq \tilde{f}'(0) + \sqrt{\tilde{f}''(0)} - \frac{\sqrt{\tilde{f}''(0)}}{1 + t\sqrt{\tilde{f}''(0)}}.
\end{aligned} \tag{35}$$

We integrate (39) again and get

$$\tilde{f}(t) \geq \tilde{f}(0) + t\tilde{f}'(0) + t\sqrt{\tilde{f}''(0)} - \log\left(1 + t\sqrt{\tilde{f}''(0)}\right). \tag{36}$$

Notice that the right hand side is a convex function of t . Then we can find t^* that reaches its minimum. We need to solve for t when its derivative is 0. We get

$$t^* = \frac{-\tilde{f}'(0)}{\tilde{f}''(0) + \sqrt{\tilde{f}''(0)}\tilde{f}'(0)}.$$

Since (39) is always true for $t \geq 0$, we can plug in t^* and get:

$$\begin{aligned}
\inf_{t \geq 0} \tilde{f}(t) &\geq \tilde{f}(0) + t^* \tilde{f}'(0) + t^* \sqrt{\tilde{f}''(0)} - \log \left(1 + t^* \sqrt{\tilde{f}''(0)} \right) \\
&= \tilde{f}(0) + \left(\tilde{f}'(0) + \sqrt{\tilde{f}''(0)} \right) \frac{-\tilde{f}'(0)}{\tilde{f}''(0) + \sqrt{\tilde{f}''(0)} \tilde{f}'(0)} \\
&\quad - \log \left(1 - \frac{\tilde{f}'(0)}{\tilde{f}''(0) + \sqrt{\tilde{f}''(0)} \tilde{f}'(0)} \sqrt{\tilde{f}''(0)} \right) \\
&= \tilde{f}(0) + \frac{-\tilde{f}'(0)}{\sqrt{\tilde{f}''(0)}} - \log \left(\frac{\tilde{f}''(0)}{\tilde{f}''(0) + \sqrt{\tilde{f}''(0)} \tilde{f}'(0)} \right) \\
&= \tilde{f}(0) - \frac{\tilde{f}'(0)}{\sqrt{\tilde{f}''(0)}} + \log \left(1 + \frac{-\tilde{f}'(0)}{\sqrt{\tilde{f}''(0)}} \right)
\end{aligned} \tag{37}$$

We did all these calculations because now we can treat $\frac{\tilde{f}'(0)}{\sqrt{\tilde{f}''(0)}}$ as a variable. Notice that $\tilde{f}'(0) = v^T \nabla f(x)$ and $\tilde{f}''(0) = v^T \nabla^2 f(x) v$. Recall the inequality we got for the Newton decrement. Here, (33) can be rewritten as

$$\lambda(x) \geq \frac{\tilde{f}'(0)}{\sqrt{\tilde{f}''(0)}}.$$

Consider the function $g(x) = x + \log(1 - x)$, (\log is of base e). We know that $g(0) = 0$ and $g'(x) = \frac{-x}{1-x}$, which is negative on $(0, 1)$. So the function g is decreasing on $(0, 1)$. Then the following inequality is true for any descent direction v provided that $\lambda(x) < 1$,

$$-\frac{\tilde{f}'(0)}{\sqrt{\tilde{f}''(0)}} + \log \left(1 + \frac{-\tilde{f}'(0)}{\sqrt{\tilde{f}''(0)}} \right) \geq \lambda(x) + \log(1 - \lambda(x)), \tag{38}$$

$$f(x^*) = \inf_{t \geq 0} \tilde{f}(t) \geq \tilde{f}(0) + \lambda(x) + \log(1 - \lambda(x)). \tag{39}$$

On the last line, $f(x^*) = \inf_{t \geq 0} \tilde{f}(t)$ is because we can choose v to be any descent direction. Notice that $\lambda(x) + \log(1 - \lambda(x)) \geq -\lambda(x)^2$ on the interval $(0, 0.68)$. Then (39) can be simplified into

$$f(x^*) \geq f(x) - \lambda(x)^2, \lambda(x)^2 \geq f(x) - f(x^*), \tag{40}$$

provided that $\lambda(x) < 0.68$. Here we get the desired result.

Remember the termination condition for the general Newton's method is $\frac{\lambda(x)^2}{2} \leq \epsilon$. If the objective function we are dealing with is self-concordant, then we can double the value still gives us a valid upper bound. In conclusion, the termination condition of Newton's method for self-concordant functions become $\lambda(x)^2 < \epsilon$ where $\epsilon < 0.68^2$.

Convergence Analysis of Newton's Method for Self-concordant Functions

Assume the objective function is strictly convex and self-concordant. Now we need neither upper or lower bound on $\nabla^2 f(x)$ nor the Lipschitz condition. Instead, we only use the assumption: self-concordance and the Newton decrement will replace $\|\nabla f(x)\|_2$. Similar to the classic analysis of the Newton method, there are two phases: **Damped Newton Phase** where $\lambda(x_k) > \eta$ and **Quadratically Convergent Phase** where $\lambda(x_k) \leq \eta$, $\eta \in (0, 1/4)$.

Damped Newton Phase

Define $\tilde{f}(t) = f(x + t\Delta x)$. So far, we have not used the upper bound (34) for the second derivative of self-concordant functions. We are going to use it now. Similar idea as before, we integrating the upper bound of (34):

$$\begin{aligned}
 \tilde{f}'(t) - \tilde{f}'(0) &= \int_0^t \tilde{f}''(x) dx \\
 &\leq \int_0^t \frac{\tilde{f}''(0)}{(1 - x\sqrt{\tilde{f}''(0)})^2} dx \\
 &= \frac{\sqrt{\tilde{f}''(0)}}{1 - x\sqrt{\tilde{f}''(0)}} \Big|_0^t \\
 &= \frac{\sqrt{\tilde{f}''(0)}}{1 - t\sqrt{\tilde{f}''(0)}} - \sqrt{\tilde{f}''(0)}. \\
 \tilde{f}'(t) &\leq \tilde{f}'(0) + \frac{\sqrt{\tilde{f}''(0)}}{1 - t\sqrt{\tilde{f}''(0)}} - \sqrt{\tilde{f}''(0)}.
 \end{aligned} \tag{41}$$

We integrate (41) again and get

$$\tilde{f}(t) \leq \tilde{f}(0) + t\tilde{f}'(0) - t\sqrt{\tilde{f}''(0)} - \log \left(1 - t\sqrt{\tilde{f}''(0)} \right).$$

Plugging in $\tilde{f}'(0) = -\lambda(x)^2$ and $\tilde{f}''(0) = \lambda(x)^2$, we get

$$\tilde{f}(t) \leq \tilde{f}(0) - t\lambda(x)^2 - t\lambda(x) - \log(1 - t\lambda(x)). \tag{42}$$

Remember the above inequality is valid when $0 \leq t \leq f''(0)^{-\frac{1}{2}} = \frac{1}{\lambda(x)}$.

Claim 0.4.6.9. *The Backtracking Line Search always ends up with a step size*

$$t \geq \frac{\gamma}{1 + \lambda(x)}.$$

Proof. First of all, $\frac{\gamma}{1+\lambda(x)} < \frac{1}{\lambda(x)}$ because $\gamma < 1$, which is within the possible range for t . Let $\tilde{t} = \frac{1}{1+\lambda(x)}$. Plugging in \tilde{t} to (42), we get

$$\begin{aligned}\tilde{f}(\tilde{t}) &\leq \tilde{f}(0) - \tilde{t}\lambda(x)^2 - \tilde{t}\lambda(x) - \log(1 - \tilde{t}\lambda(x)) \\ &= \tilde{f}(0) - \lambda(x) + \log(1 + \lambda(x)).\end{aligned}$$

Now consider the function $h(x) = -x + \log(1 + x) + \frac{x^2}{2(x+1)}$. We know that for $x > -1$,

$$h(0) = 0, h'(x) = \frac{-x^2}{2(1+x)^2} \leq 0.$$

Since the derivative is always less than 0, we know h is a monotonically decreasing function. Since $h(0) = 0$, we know that

$$h(x) = -x + \log(1 + x) + \frac{x^2}{2(x+1)} \leq 0.$$

Replace x by $\lambda(x)$, we get that

$$\begin{aligned}h(x) &= -\lambda(x) + \log(1 + \lambda(x)) + \frac{\lambda(x)^2}{2(\lambda(x) + 1)} \leq 0, \\ -\lambda(x) + \log(1 + \lambda(x)) &\leq \frac{-\lambda(x)^2}{2(\lambda(x) + 1)} \\ &\leq \frac{-\alpha\lambda(x)^2}{(\lambda(x) + 1)} \\ &= \frac{-\alpha\lambda(x)^2}{\tilde{t}}.\end{aligned}\tag{43}$$

Therefore, $f(x) - f(x_+) = \tilde{f}(0) - \tilde{f}(\tilde{t}) \geq \frac{-\alpha\lambda(x)^2}{\tilde{t}}$, where $f(x_+)$ denotes the next step. \square

In conclusion, at the end of the Backtracking Line Search, we have $t \geq \frac{\gamma}{1+\lambda(x)}$. At each iteration step of the Damped Newton Phase, the function value decreases at least

$$\alpha\gamma \frac{\lambda(x)^2}{1 + \lambda(x)} \geq \alpha\gamma \frac{\eta^2}{1 + \eta}$$

since $\frac{x^2}{1+x}$ is an increasing function for $x \geq 0$.

Quadratically Convergent Phase

In order to show that the unit step size is always valid for the Backtracking Line Search, we need to restrict η to a smaller range. Take $\eta = \frac{1-2\alpha}{4} \leq \frac{1}{4}$. Then $\lambda(x) \leq \frac{1-2\alpha}{4}$. Plugging in $t = 1$ to (42), we get

$$\begin{aligned}\tilde{f}(1) &\leq \tilde{f}(0) - \lambda(x)^2 - \lambda(x) - \log(1 - \lambda(x)) \\ &\leq \tilde{f}(0) - \frac{1}{2}\lambda(x)^2 + \lambda(x)^3\end{aligned}\tag{44}$$

The above inequality comes from the fact that $-x - \log(1 - x) \leq \frac{1}{2}x^2 + x^3$ for $x \in [0, 0.81]$ and $\lambda(x) \leq \frac{1-2\alpha}{4} < 0.81$. We can simplify the inequality even further:

$$\begin{aligned}
\lambda(x) &\leq \frac{1-2\alpha}{4}, \\
\frac{1}{2} - \alpha - \lambda(x) &\geq 0, \\
\lambda(x)^2 \left(\frac{1}{2} - \alpha - \lambda(x) \right) &\geq 0, \\
\frac{1}{2} \lambda(x)^2 - \lambda(x)^3 &\geq \alpha \lambda(x)^2, \\
-\frac{1}{2} \lambda(x)^2 + \lambda(x)^3 &\leq -\alpha \lambda(x)^2
\end{aligned} \tag{45}$$

Therefore, (44) can be simplified into

$$\tilde{f}(1) \leq \tilde{f}(0) - \alpha \lambda(x)^2,$$

or

$$f(x) - f(x_+) = \tilde{f}(0) - \tilde{f}(1) \geq \alpha \lambda(x)^2.$$

From this, we know that the unit step size $t = 1$ satisfies the exit condition of the Backtracking Line Search. To prove that the convergent rate is quadratic, we need the following inequality:

$$\lambda(x_+) \leq \left(\frac{\lambda(x)}{1 - \lambda(x)} \right)^2,$$

which is true for $\lambda(x) < 1$. The proof will not be presented here but can be found in ***. Since we have an even smaller upper bound $\lambda(x) \leq 1/4$, we get that

$$\begin{aligned}
\frac{1}{(1 - \lambda(x))^2} &\leq 2, \\
\lambda(x_+) &\leq 2\lambda(x)^2.
\end{aligned} \tag{46}$$

Hence, by (40), $f(x_+) - f(x^*) \leq \lambda(x_+)^2 \leq (2\lambda(x))^2$. Applying (46) recursively, we get that

$$\begin{aligned}
f(x_k) - f(x^*) &\leq \lambda(x_k)^2 \leq \left(\frac{1}{2} \right)^{1-2^k} (\lambda(x_0))^{2^k} \\
&\leq \left(\frac{1}{2} \right)^{1-2^k} \left(\frac{1}{4} \right)^{2^k} \\
&= \left(\frac{1}{2} \right)^{2^k+1}
\end{aligned} \tag{47}$$

Lastly, similar to the classic Newton's method, we can find an upper bound for the total number of iterations:

$$(f(x_0) - f(x^*)) \frac{1 + \eta}{\alpha \gamma \eta^2} + 6.$$

0.5 Interior-point Method

In this section, we are going to introduce the Interior-point method which solves inequality and equality constrained convex optimization problems. Here is the statement of the problem:

$$\begin{aligned} & \text{minimize} && f_0(\mathbf{x}) \\ & \text{subject to} && f_i(\mathbf{x}) \leq 0, \text{ for } 1 \leq i \leq m, \\ & && h_j(\mathbf{x}) = 0, \text{ for } 1 \leq j \leq n, \end{aligned} \tag{48}$$

where $f_0, f_1, \dots, f_m : \mathbb{R}^n \rightarrow \mathbb{R}$ are convex and twice continuously differentiable functions and $h_1, \dots, h_n : \mathbb{R}^n \rightarrow \mathbb{R}$ are affine functions. Let \mathcal{D} be the domain of this optimization problem and let $\mathcal{X} \subseteq \mathcal{D}$ is the set of feasible points, that is for all $\mathbf{x} \in \mathcal{X}$, $f_i(\mathbf{x}) \leq 0$ and $h_i(\mathbf{x}) = 0$. Let f^* be the optimal value and x^* be the point that gives $f_0(x^*) = f^*$. Before we go into the details, we first introduce some basic definitions about primal and dual problems.

0.5.1 Primal and Dual Problem

We refer to the above constrained convex optimization problem as the **primal** problem. Define the **Lagrangian Equation** $\mathcal{L} : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$ as the following:

$$\mathcal{L}(\mathbf{x}, \mathbf{u}, \mathbf{v}) = f_0(\mathbf{x}) + \sum_{i=1}^m u_i f_i(\mathbf{x}) + \sum_{j=1}^n v_j h_j(\mathbf{x}),$$

where $\mathbf{u} \in \mathbb{R}^m$ and $\mathbf{v} \in \mathbb{R}^n$ are called Lagrange multipliers. We restrict the domain of \mathbf{u} to be \mathbb{R}_+^m , i.e., $u_i \geq 0$ for all $1 \leq i \leq m$, in order to let inequality constraints f_i make sense.

For every $\mathbf{x} \in \mathcal{X}$, $\mathbf{u} \succeq 0$, we have $\mathcal{L}(\mathbf{x}, \mathbf{u}, \mathbf{v}) \leq f_0(\mathbf{x})$. If we pick $u_i = 0$ whenever $f_i(\mathbf{x}) < 0$, then

$$\inf_{\mathbf{x}} \sup_{\mathbf{u} \succeq 0, \mathbf{v}} \mathcal{L}(\mathbf{x}, \mathbf{u}, \mathbf{v}) = f_0(x^*)$$

Note that $\mathbf{u} \succeq 0$ means that $u_i \geq 0$ pointwisely.

Now we define the **Lagrange Dual Function** as the following:

$$\begin{aligned} g(\mathbf{u}, \mathbf{v}) &= \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{u}, \mathbf{v}) \\ &= \inf_{\mathbf{x}} f_0(\mathbf{x}) + \sum_{i=1}^m u_i f_i(\mathbf{x}) + \sum_{j=1}^n v_j h_j(\mathbf{x}) \end{aligned} \tag{49}$$

Given (58) and (49), we define the **Lagrange Dual Problem** as:

$$\begin{aligned} & \text{maximize} && g(\mathbf{u}, \mathbf{v}) \\ & \text{subject to} && u_i \geq 0, \text{ for } 1 \leq i \leq m, \end{aligned} \tag{50}$$

Let g^* denote the optimal value of the Lagrange Dual Problem. The pair (u^*, v^*) with which g^* is obtained is called dual optimal. For any fixed $x \in \mathcal{X}$, $\mathcal{L}(\mathbf{x}, \mathbf{u}, \mathbf{v})$ is an affine function of \mathbf{u} and \mathbf{v} . Then g can be viewed as the pointwise infimum of the affine function of \mathbf{u} and \mathbf{v} , thus is concave. The constraint $\mathbf{u} \succeq 0$ is an affine constraint. So the system (50) is a concave maximization problem, which is also a convex optimization problem.

Theorem 0.5.1.1. *The optimal value for the Primal problem is always greater than or equal to the optimal value for the Dual problem, i.e., $f^* \geq g^*$. This inequality is called the **Weak Duality**.*

Proof. We know that

$$\begin{aligned} f^* &= \inf_{\mathbf{x}} \sup_{\mathbf{u} \succeq 0, v} \mathcal{L}(\mathbf{x}, \mathbf{u}, v), \\ g^* &= \sup_{\mathbf{u} \succeq 0, v} \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{u}, v). \end{aligned}$$

Notice that the following inequality is necessary:

$$\inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{u}, v) \leq \sup_{\mathbf{u} \succeq 0, v} \mathcal{L}(\mathbf{x}, \mathbf{u}, v).$$

Then we first take the infimum over \mathbf{x} on both sides to get

$$\inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{u}, v) \leq \inf_{\mathbf{x}} \sup_{\mathbf{u} \succeq 0, v} \mathcal{L}(\mathbf{x}, \mathbf{u}, v),$$

and then take the supremum over $\mathbf{u} \succeq 0, v$ on both sides,

$$\sup_{\mathbf{u} \succeq 0, v} \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{u}, v) \leq \inf_{\mathbf{x}} \sup_{\mathbf{u} \succeq 0, v} \mathcal{L}(\mathbf{x}, \mathbf{u}, v) \Rightarrow f^* \geq g^*.$$

□

We define **Duality Gap** = $f^* - g^*$ and **Strong Duality** states there are x^*, u^*, v^* such that $f^* = g^*$.

Theorem 0.5.1.2. Slater's Condition

Suppose we have a convex optimization problem. If there exists at least one strict feasible $\tilde{x} \in \mathbb{R}^n$, then the strong duality holds.

The very last thing in this section is the Karush-Kuhn-Tucker conditions or KKT conditions. Given (58), the KKT conditions are

- $\nabla f(x) + \sum_{i=1}^m u_i \nabla f_i(x) + \sum_{j=1}^n v_j \nabla h_j(x) = 0$ (Stationarity)
- $u_i f_i(x) = 0$ for $1 \leq i \leq m$ (Complementary slackness)
- $f_i(x) \leq 0, h_j(x) = 0$ for $1 \leq i \leq m$ (Primal feasibility)
- $u_i \geq 0$ for $1 \leq i \leq m$ (Dual feasibility)

Theorem 0.5.1.3. *(x^*, u^*, v^*) are primal and dual solutions such that the strong duality holds if and only if (x^*, u^*, v^*) satisfy the KKT conditions.*

Proof. (\Rightarrow) Since the strong duality holds with (x^*, u^*, v^*) , we know that

$$\begin{aligned} f(x^*) &= g(u^*, v^*) \\ &= \inf_{\mathbf{x}} f_0(\mathbf{x}) + \sum_{i=1}^m u_i^* f_i(\mathbf{x}) + \sum_{j=1}^n v_j^* h_j(\mathbf{x}) \\ &\leq f_0(x^*) + \sum_{i=1}^m u_i^* f_i(x^*) + \sum_{j=1}^n v_j^* h_j(x^*) \\ &\leq f(x^*). \end{aligned}$$

We have the first inequality because plugging in any value for \mathbf{x} would give a value greater than the infimum over all \mathbf{x} . The second inequality comes from the fact that $f_i(x^*) \leq 0$, $u_i^* \geq 0$, and $h_i(x^*) = 0$. Since we cannot have $f(x^*) < f(x^*)$, all inequalities should be equality. Therefore, we get x^* gives the infimum of $\mathcal{L}(\mathbf{x}, \mathbf{u}, \mathbf{v}) = f_0(\mathbf{x}) + \sum_{i=1}^m u_i^* f_i(\mathbf{x}) + \sum_{j=1}^n v_j^* h_j(\mathbf{x})$, which give the stationarity condition. On top of that, the last equality tells us that $\sum_{i=1}^m u_i^* f_i(x^*) = 0$, which is the complementary slackness condition. The primal the dual feasibility are obviously true.

(\Leftarrow) By integrate the stationarity condition, we get the following:

$$\begin{aligned} g(u^*, v^*) &= f_0(x^*) + \sum_{i=1}^m u_i^* f_i(x^*) + \sum_{j=1}^n v_j^* h_j(x^*) \\ &= f(x^*). \end{aligned}$$

We get the last equality from complementary slackness and primal feasibility conditions. Hence, we get that (x^*, u^*, v^*) satisfy that $f(x^*) = g(u^*, v^*)$, which gives the strong duality. \square

0.5.2 Newton's Method with Equality constraints

In this section, we talk about the extension of Newton's Method to optimization problems with equality constraints. This method will play an important role for our later discussion about the Interior-point Method. Consider the following convex quadratic function with equality constraints:

$$\begin{aligned} &\text{minimize} \quad \frac{1}{2}x^T A x + b^T x + c \\ &\text{subject to} \quad Qx = q, \end{aligned} \tag{51}$$

where A is a $\mathbb{R}^{n \times n}$ symmetric positive semidefinite matrix and $Q \in \mathbb{R}^{p \times n}$, $\text{rank}(Q) = p < n$. The assumption on the dimension of Q says that there are fewer equality constraints than variables.

Since there is no inequality constraints, the KKT conditions can be simplified into the following:

- $Ax^* + b + Q^T v^* = 0$
- $Qx^* = q$

From the previous section, we know that x^* is the primal solution and v^* is the dual solution if and only if they satisfy the KKT conditions. The KKT conditions can be rewritten in the matrix form:

$$\begin{bmatrix} A & Q^T \\ Q & 0 \end{bmatrix} \begin{bmatrix} x^* \\ v^* \end{bmatrix} = \begin{bmatrix} -b \\ q \end{bmatrix}, \tag{52}$$

which is a set of $(n + p)$ linear equations. The coefficient matrix is called the KKT matrix. We can solve for x^*, v^* if the KKT matrix is nonsingular. To answer this question, we introduce the following theorem:

Theorem 0.5.2.1. *Suppose A is a $\mathbb{R}^{n \times n}$ symmetric positive semidefinite matrix and $Q \in \mathbb{R}^{p \times n}$, $\text{rank}(A) = p < n$. Then the following are equivalent to KKT matrix being nonsingular:*

- $\mathcal{N}(A) \cap \mathcal{N}(Q) = \{\mathbf{0}\}$, i.e., the only vector that satisfies $Ax = Qx = 0$ is the zero vector.

- $Qx = 0, x \neq 0$ implies $x^T Ax > 0$.
- $F^T AF \succ 0$, where $F \in \mathbb{R}^{n \times (n-p)}$ is a matrix such that $\mathcal{R}(F) = \mathcal{N}(Q)$.

Proof. (1 \Rightarrow 2) Choose $x \neq 0 \in \mathcal{N}(Q)$. Since $\mathcal{N}(A) \cap \mathcal{N}(Q) = \{0\}$, we know that $x \notin \mathcal{N}(A)$, which means that $x^T Ax > 0$. Proved 2.

(2 \Rightarrow 3) Choose any $x \in \mathbb{R}^{n-p}$. Then $Fx = z$ for some $z \in \mathcal{N}(Q)$. $x^T F^T AFx = z^T Az$. By the Rank-Nullity Theorem, $\text{rank}(A) + \dim(\mathcal{N}(A)) = n$, i.e., $\dim(\mathcal{R}(F)) = n - p$. So $\dim(\mathcal{N}(F)) = 0$, which means that $\mathcal{N}(F) = \{0\}$. If $z \neq 0$, then we have $x \neq 0$ and $x^T F^T AFx = z^T Az > 0$. Therefore, for all $x \neq 0$, $x^T F^T AFx > 0$, which means that $F^T AF \succ 0$.

(3 \Rightarrow 1) Choose any $x \neq 0$. $\dim(\mathcal{N}(F)) = 0$ implies that $z = Fx \neq 0$. $\mathcal{R}(F) = \mathcal{N}(Q)$ implies that $z \in \mathcal{N}(A)$. Since $F^T AF \succ 0$, we get that $z^T Az > 0$ for all $z \in \mathcal{N}(A)$. This means that $\mathcal{N}(A) \cap \mathcal{N}(Q) = \{0\}$. \square

In conclusion, if A is instead a symmetric positive definite matrix, then the KKT matrix is always nonsingular.

Next, we discuss how Newton's Method can be applied. The previous discussion about convex quadratic function is useful because we need to use the second-order Taylor approximation. Suppose we are interested in finding the minimum of a twice continuously differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ whose Hessian matrix is positive definite. Let Q be the same matrix that denotes the equality constraints. Assume that we start at a feasible point x_0 such that $Qx_0 = q$. We approximate $f(x_0)$ by its second-order Taylor expansion $\hat{f}(x_0 + v)$ and get the following minimization problem:

$$\begin{aligned} & \text{minimize} && \frac{1}{2}v^T \nabla^2 f(x_0)v + \nabla f(x_0)^T v + f(x_0) \\ & \text{subject to} && Q(x_0 + v) = q, \text{ equivalently } Qv = 0. \end{aligned} \tag{53}$$

Here, the variable is v which is a descent direction that decreases the function value. Define the Newton step Δx to be the previous optimal solution x^* and v^* as before to be the dual solution. The KKT conditions for this minimization problem become:

- $\nabla^2 f(x_0)\Delta x + \nabla f(x_0) + Q^T v^* = 0$
- $Q\Delta x = 0$

To write this in the matrix form, we get

$$\begin{bmatrix} \nabla^2 f(x_0) & Q^T \\ Q & 0 \end{bmatrix} \begin{bmatrix} \Delta x \\ v^* \end{bmatrix} = \begin{bmatrix} -\nabla f(x_0) \\ 0 \end{bmatrix}. \tag{54}$$

Since $\nabla^2 f(x_0)$ is positive definite, we know that the KKT matrix must be nonsingular. So we can always invert the matrix and solve for Δx and v^* . Hence, we have computed the Newton step Δx which is the optimal solution that satisfies the equality constraints and minimize the second-order approximation function. However, we still get some confusions that haven't been verified yet. How do we know that Δx is a descent direction for $f(x)$, i.e., $f(x + t\Delta x) < f(x)$? Similar to the classic Newton's Method, is Newton decrement $\lambda(x)$ still a good estimate of the distance between $f(x)$ and $\inf_v \hat{f}(x + v)$?

Newton Decrement

We define the Newton's Decrement for the equality constrained minimization problem to be:

$$\lambda(x) = (\Delta x^T \nabla^2 f(x) \Delta x)^{\frac{1}{2}}.$$

This is actually the same definition as in chapter 3. Previously, we defined $\Delta x = -\nabla^2 f(x)^{-1} \nabla f(x)$, then $\lambda(x) = (\Delta x^T \nabla^2 f(x) \Delta x)^{\frac{1}{2}} = (\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x))^{\frac{1}{2}}$, which is our previous definition for $\lambda(x)$ for the classic Newton's Method. Now, we discover the relation between $\lambda(x)$ and $f(x) - \inf_v f(x + v)$

Theorem 0.5.2.2. *The difference between f and its second-order Taylor approximation satisfies $f(x) - \inf\{\hat{f}(x + v) \mid Qv = 0\} = \frac{\lambda(x)^2}{2}$.*

Proof. $\hat{f}(x + \Delta x) = \inf\{\hat{f}(x + v) \mid Qv = 0\}$. The Newton step is defined by the KKT conditions:

$$\begin{bmatrix} \nabla^2 f(x) & Q^T \\ Q & 0 \end{bmatrix} \begin{bmatrix} \Delta x \\ v^* \end{bmatrix} = \begin{bmatrix} -\nabla f(x) \\ 0 \end{bmatrix}. \quad (55)$$

From the first row of the KKT matrix, we get that

$$\nabla^2 f(x) \Delta x + Q^T v^* = -\nabla f(x).$$

We multiply Δx^T to both sides of the equation:

$$\Delta x^T \nabla^2 f(x) \Delta x = -\Delta x^T \nabla f(x), \quad (56)$$

because $(Q\Delta x)^T v^* = 0$. Then

$$\begin{aligned} \hat{f}(x + \Delta x) &= \frac{1}{2} \Delta x^T \nabla^2 f(x) \Delta x + \nabla f(x)^T \Delta x + f(x) \\ &= \frac{1}{2} \Delta x^T \nabla^2 f(x) \Delta x - \Delta x^T \nabla^2 f(x) \Delta x + f(x) \\ &= -\frac{1}{2} \lambda(x)^2 + f(x). \end{aligned} \quad (57)$$

So $f(x) - \hat{f}(x + \Delta x) = \frac{1}{2} \lambda(x)^2$ finishes the proof. \square

This theorem means that similar to the classic Newton's Method, $\frac{1}{2} \lambda(x)^2$ is a good indicator of accuracy and can serve as the stopping criterion. Now, we explain why Δx is a descent direction for $f(x)$. We just need to check that the directional derivative of f in the direction Δx is negative:

$$\left. \frac{d}{dt} f(x + t\Delta x) \right|_{t=0} = \nabla f(x)^T \Delta x = -\lambda(x)^2 < 0.$$

So the Newton step is a descent direction for $f(x)$.

Algorithm for Equality Constrained Newton's Method

The following algorithm uses the Backtracking Line Search to determine each iteration step size.

Algorithm 8: Newton's Method with Equality Constraints

```

NewtonMethod ( $f, Q, x_0, \epsilon, \alpha, \gamma$ );
Compute the initial  $\lambda^2$  and  $\Delta x$  for  $x_0$ ;
while  $\frac{\lambda^2}{2} \geq \epsilon$  do
     $x_{i+1} = x_i + \Delta x$ ;
    while  $f(x_i) - f(x_{i+1}) < -\alpha t \nabla f(x_i)^T \Delta x$ , do
         $t = \gamma t$ ;
         $x_{i+1} = x_i + t \Delta x$ ;
     $x_i = x_{i+1}$ ;
     $\begin{bmatrix} \Delta x \\ v^* \end{bmatrix} = \begin{bmatrix} \nabla^2 f(x) & Q^T \\ Q & 0 \end{bmatrix}^{-1} \begin{bmatrix} -\nabla f(x) \\ 0 \end{bmatrix}$ ;
     $\lambda^2 = \Delta x^T \nabla^2 f(x) \Delta x$ ;
return  $x_i$ ;

```

0.5.3 Barrier Method and Logarithmic Barrier Function

Our goal now is to use equality constrained Newton's method to help solve the minimization problem with both inequality and equality constraints. Here is the statement of the problem:

$$\begin{aligned}
 & \text{minimize} && f_0(x) \\
 & \text{subject to} && f_i(x) \leq 0, \text{ for } 1 \leq i \leq m, \\
 & && Ax = b,
 \end{aligned} \tag{58}$$

where $f_0, f_1, \dots, f_m : \mathbb{R}^n \rightarrow \mathbb{R}$ are convex and twice continuously differentiable functions and $A \in \mathbb{R}^{p \times n}$, $\text{rank}(A) = p < n$. Let \mathcal{D} be the domain of this optimization problem. How can we transform the problem to be a minimization problem with only equality constraints? The clever idea is that we use a differentiable function to approximate the inequality constraint.

One immediate answer would be to use an indicator function to denote the inequality constraint, that is

$$\mathbb{1}_{\{f_i(x) > 0\}} = \begin{cases} 0 & \text{if } f_i(x) \leq 0 \\ 1 & \text{if } f_i(x) > 0 \end{cases}$$

Then the original minimization problem can be changed into the following:

$$\begin{aligned}
 & \text{minimize} && f_0(x) + \sum_{i=1}^m \mathbb{1}_{\{f_i(x) > 0\}} \cdot \infty \\
 & \text{subject to} && Ax = b.
 \end{aligned}$$

We get rid of the inequality constraints but the objective function is no longer differentiable. So using the indicator function is not a good choice. Our next step is to approximate this indicator function.

We use the following log function to approximate the indicator function:

$$I_t(x) = -\frac{1}{t} \log(-x), \text{dom}(I_t) = -\mathbb{R}_+$$

where $t > 0$ is a variable that determines the accuracy of the approximation. The following Figure 12 shows the approximation. The red curve has a much bigger value of t compared to the blue curve. As a result, the red curve is a better approximation to the indicator function than the blue curve. From the picture, we know that as t gets larger, $I_t(x)$ approximates the indicator function better. Unlike the indicator function, $I_t(x)$ is convex, differentiable and it increases to positive infinity as x goes to 0.

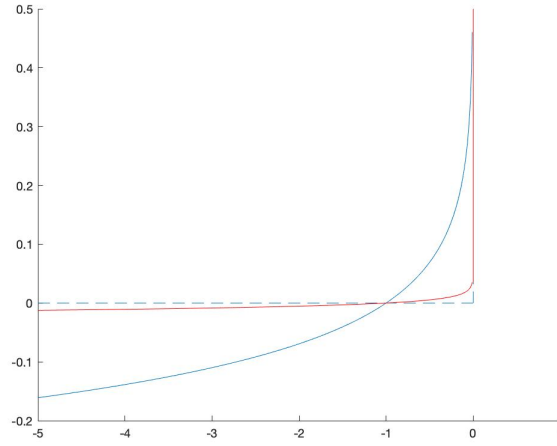


Figure 12: The log function that approximates the indicator function

We can update our minimization problem by replacing the indicator function with this log function:

$$\begin{aligned} \text{minimize} \quad & f_0(x) + \sum_{i=1}^m I_t(f_i(x)) = f_0(x) + \sum_{i=1}^m -\frac{1}{t} \log(-f_i(x)) \\ \text{subject to} \quad & Ax = b. \end{aligned}$$

We define $\phi(x) = -\sum_{i=1}^m \log(-f_i(x))$ with $\text{dom}(\phi(x)) = \{x \in \mathcal{D} | f_i(x) < 0, i = 1, 2, \dots, m\}$ to be the **logarithmic barrier function**. This new objective function is convex and twice continuously differentiable. To simplify the problem, we multiply t to the objective function and get the following equivalent problem:

$$\begin{aligned} \text{minimize} \quad & t f_0(x) + \phi(x) \\ \text{subject to} \quad & Ax = b. \end{aligned} \tag{59}$$

This new problem has a different minimum function value but the minimizer x^* is the same as in the previous problem. For now, we assume that the above minimization problem can be solved by the equality constrained Newton's Method. Here we compute the gradient and Hessian matrix of

the objective function that will be useful when using Newton's Method:

$$\begin{aligned}\nabla(t f_0(x) + \phi(x)) &= t \nabla f_0(x) + \nabla \phi(x) = t \nabla f_0(x) - \sum_{i=1}^m \frac{\nabla f_i(x)}{f_i(x)}, \\ \nabla^2(t f_0(x) + \phi(x)) &= t \nabla^2 f_0(x) + \nabla^2 \phi(x) \\ &= t \nabla^2 f_0(x) + \sum_{i=1}^m \frac{1}{f_i(x)^2} \nabla f_i(x)^T \nabla f_i(x) + \sum_{i=1}^m \frac{1}{-f_i(x)} \nabla^2 f_i(x).\end{aligned}$$

For each $t > 0$, we have a corresponding minimization problem (59), from which we get a unique solution $x^*(t)$. We define the **central path** to be the set of points $x^*(t)$, $t > 0$. For future reference, x^*, u^*, v^* all depend on t . Different t gives different value for these three variables. To simplify the notation, we just use x^*, u^*, v^* to mean $x^*(t), u^*(t), v^*(t)$. By the KKT conditions, we know that points on the central path x^* satisfies

- $t \nabla f_0(x^*) + \nabla \phi(x^*) + A^T \hat{v} = 0$, for some $v^* \in \mathbb{R}^p$
- $f_i(x^*) < 0$, for $i = 1, 2, \dots, m$, $Ax^* = b$

From the above conditions, we can derive a dual solution (u^*, v^*) that satisfy

$$\nabla f(x^*) + \sum_{i=1}^m u_i^* \nabla f_i(x^*) + A^T v^* = 0. \quad (60)$$

Define u^* and v^* by the following formula,

$$u_i^* = \frac{-1}{t f_i(x^*)} \text{ for } i = 1, 2, \dots, m, \text{ and } v^* = \frac{\hat{v}}{t} \quad (61)$$

We see that $u_i^* > 0$ because $f_i(x^*) < 0$ and $t > 0$. In addition,

$$\begin{aligned}0 &= t \nabla f_0(x^*) + \nabla \phi(x^*) + A^T \hat{v} \\ &= t \left(\nabla f_0(x^*) - \sum_{i=1}^m \frac{\nabla f_i(x)}{t f_i(x)} + A^T v^* \right) \\ &= t \left(\nabla f_0(x^*) - \sum_{i=1}^m u_i^* \nabla f_i(x) + A^T v^* \right)\end{aligned}$$

Therefore, we find the dual solution such that (60) is true. This means that x^* minimizes the lagrangian equation $\mathcal{L}(x, u^*, v^*) = f_0(x) + \sum_{i=1}^m u_i^* f_i(x) + (Ax - b)^T v^*$. Since the dual function $g(u^*, v^*)$ is defined to be the infimum of $\mathcal{L}(x, u^*, v^*)$ over all x , we get that

$$\begin{aligned}g(u^*, v^*) &= f_0(x^*) + \sum_{i=1}^m u_i^* f_i(x^*) + (Ax^* - b)^T v^* \\ &= f_0(x^*) + \sum_{i=1}^m \frac{-1}{t f_i(x^*)} f_i(x^*) + (Ax^* - b)^T v^* \\ &= f_0(x^*) - \frac{m}{t}.\end{aligned} \quad (62)$$

The last term $(Ax^* - b)^T v^*$ disappears because $Ax^* = b$. By Weak Duality, we know that

$$\frac{m}{t} = f_0(x^*) - g(u^*, v^*) \geq f_0(x^*) - f_0^*,$$

where we proved the important fact that as $t \rightarrow \infty$, $f_0(x^*(t))$ approaches to the infimum of f_0 . where we proved the important fact that as $t \rightarrow \infty$, $f_0(x^*(t))$ approaches to the infimum of f_0 . As said before, x^* is a variable that depends on t . In order to see the result clearly, we change x^* into $x^*(t)$.

Algorithm for Barrier Method

The following algorithm uses nested while loop. The outer while loop generates a new objective function $tf_0 + \phi$ and updates t and x^* . The inner while loop uses Newton's Method to compute the minimizer x^* given the objective function.

Algorithm 9: Newton's Method with Equality Constraints

BarrierMethod ($f_0, f_i, A, x_0, \epsilon, \alpha, \gamma, t_0 > 0, \beta > 1$);

Compute the $\phi(x)$, KKT matrix and Δx for x_0 ;

$t = t_0$;

while $m/t \geq \epsilon$ **do**

1. Use equality constrained Newton's Method to minimize $tf_0 + \phi$ with x_0 as the starting point and solve for x^* ;
2. $x_0 = x^*$;
3. $t = \beta t$;

return x_0 ;

The iteration step for Newton's Method is the following:

while $\frac{\lambda^2}{2} \geq \epsilon$ **do**

$x_{i+1} = x_i + \Delta x$;

while $f(x_i) - f(x_{i+1}) < -\alpha t \nabla f(x_i)^T \Delta x$, **do**

$t = \gamma t$;

$x_{i+1} = x_i + t \Delta x$;

$x_i = x_{i+1}$;

$$\begin{bmatrix} t \nabla^2 f_0(x) + \nabla^2 \phi(x) & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} x^* \\ v^* \end{bmatrix} = \begin{bmatrix} -t \nabla f(x) - \nabla \phi(x) \\ 0 \end{bmatrix} \cdot;$$

$\lambda^2 = \Delta x^T \nabla^2 f(x) \Delta x$;

0.6 Finite Element Method

[?]