

Non-Linear Optimization 447/550 Notes

Farid Rajkotia Zaheer

1 Mathematical Preliminaries

Theorem 1.1. (Spectral Factorization Theorem:) Let $A \in \mathbb{R}^{n \times n}$ be a $n \times n$ symmetric matrix. Then there exists an orthogonal matrix U i.e. $UU^T = U^T U = I$ and a diagonal matrix D such that

$$U^T A U = D.$$

Here the columns of U are eigenvectors of A and form an orthogonal basis. The entries of D are the eigenvalues of A .

Theorem 1.2. (Mean Value Theorem) Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable. Then, for every \mathbf{x} and $\bar{\mathbf{x}}$ in \mathbb{R}^n there exists $\mathbf{z} \in (\mathbf{x}, \bar{\mathbf{x}})$ where $(\mathbf{x}, \bar{\mathbf{x}}) := \{\mathbf{x} + s(\bar{\mathbf{x}} - \mathbf{x}) : 0 < s < 1\}$, is the line segment joining \mathbf{x} and $\bar{\mathbf{x}}$ excluding the end points such that,

$$f(\mathbf{x}) = f(\bar{\mathbf{x}}) + \nabla f(\mathbf{z})^T (\mathbf{x} - \bar{\mathbf{x}}).$$

Theorem 1.3. (Linear Approximation Theorem) Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice differentiable. Then, for every \mathbf{x} and $\bar{\mathbf{x}}$ in \mathbb{R}^n there exists $\mathbf{z} \in (\mathbf{x}, \bar{\mathbf{x}})$

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

Theorem 1.4. (Quadratic Approximation Theorem)

2 Graphical Optimization

A neat method for solving both linear and non-linear programming problems up to 3 dimensions can be done by plotting out the feasible region of the problem at hand. Explicitly, we may apply this technique if the number of variables in the objective function does not exceed 2.

The following are the steps needed to obtain a graphical solution:

1. Draw the boundaries of each of the constraint functions. The region that is obtained is the feasible set.
2. Draw several contours of the objective function and determine in which direction the function is decreasing/increasing (depending if the optimization is minimizing or maximizing).
3. Select the point that belongs to the feasible set and is the lowest/highest value of the objective function.

In order to determine in which direction the objective function is decreasing or increasing can be resolved by recalling the directional derivative from vector calculus.

Definition 2.1. (Directional Derivative) Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$, the directional derivative of f at x in the direction d is,

$$f'(x : d) = \lim_{t \rightarrow 0^+} \frac{f(x + td) - f(x)}{t}.$$

Given the limit exists and is finite.

Notice that if $f'(x : d) > 0$ then for $t > 0$ small enough we have that $f(x + td) > f(x)$. So we deduce that the function value increases in the direction d as we vary over small enough step size t .

Then we know the gradient of f in the direction of some vector d is given as,

$$\nabla f(x)^T d = \nabla f(x) \cdot d = \|\nabla f(x)\| \|d\| \cos \theta.$$

So, $\nabla f(x)^T d > 0$ iff $\theta \in [0, 2\pi)$.

Definition 2.2. (Level Set) Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$, the level set of f is the set of all points where the function takes on a constant value, c .

When $n = 2$, the level set is a level curve, when $n = 3$, the level set is a level surface etc.

Recall the well known fact from vector calculus,

Theorem 2.3. ∇f is a vector perpendicular to each point of the level set of f .

It is important to understand that it is not necessarily the case that the solution to the optimization problem is also the unconstrained optimum. Note, that the unconstrained optimum, is the minimizer or maximizer of the objective function regardless of the constraints i.e. independent of the feasible set.

We may now turn our attention to looking at the setup of a non-linear programming problem and how to go about solving it graphically.

Example 2.4. Consider the non-linear minimization problem

$$\begin{aligned} \min \quad & f(x, y) x^2 + y^2 - 4y + 4 \\ \text{s.t.} \quad & g_1(x, y) x^2 - y + 2 \leq 2 \\ & g_2(x, y) x + y - 6 \leq 0. \end{aligned} \tag{1}$$

We may plot the feasible region by plotting the boundaries of the constraint functions i.e. $x^2 - y + 2 = 2$ and $x + y - 6 = 0$.

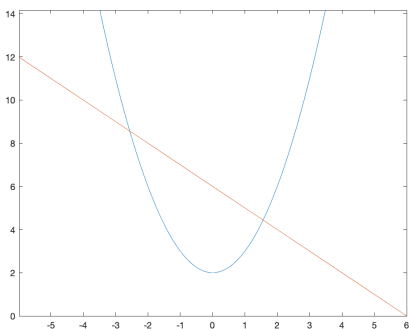


Figure 1: The feasible set is the bounded region.

By now plotting contours of the objective function we may find the optimal solution that minimizes the function.

3 Unconstrained Optimization

Unconstrained optimization refers to the set of linear or non-linear optimization problems where there are no constraints of any kind placed on the objective function. In other words, the feasible set of the problem is the entire space over which the objective function is defined in the problem. In such problems we seek to find global or local optima for the objective function.

3.1 Global and Local Optima

Consider the set $S \subset \mathbb{R}^n$ and the mapping $f : S \rightarrow \mathbb{R}$. Then,

Definition 3.1. \mathbf{x}^* is called the *global minimum point* of f over S if $f(\mathbf{x}) \geq f(\mathbf{x}^*)$ for any $\mathbf{x} \in S$.

Definition 3.2. \mathbf{x}^* is called a *strict global minimum point* of f over S if $f(\mathbf{x}) > f(\mathbf{x}^*)$ for any $\mathbf{x}^* \neq \mathbf{x} \in S$.

Definition 3.3. \mathbf{x}^* is called the *global maximum point* of f over S if $f(\mathbf{x}) \leq f(\mathbf{x}^*)$ for any $\mathbf{x} \in S$.

Definition 3.4. \mathbf{x}^* is called a *strict global maximum point* of f over S if $f(\mathbf{x}) < f(\mathbf{x}^*)$ for any $\mathbf{x}^* \neq \mathbf{x} \in S$.

The set S on which the optimization of f is performed is called the *feasible set* and any point $\mathbf{x} \in S$ is called a *feasible solution*. For unconstrained optimization problems, S will be implicitly taken to be the whole space.

A vector $\mathbf{x}^* \in S$ is called a *global optimum* if it is either a global minimum or maximum.

The maximal and minimal values of f over S are defined as the supremum and infimum respectively,

$$\begin{aligned}\max\{f(\mathbf{x}) : \mathbf{x} \in S\} &:= \sup\{f(\mathbf{x}) : \mathbf{x} \in S\} \\ \min\{f(\mathbf{x}) : \mathbf{x} \in S\} &:= \inf\{f(\mathbf{x}) : \mathbf{x} \in S\}\end{aligned}$$

Furthermore, the set of all global maximizers and minimizers of f over S is denoted respectively as,

$$\begin{aligned}\operatorname{argmax}\{f(\mathbf{x}) : \mathbf{x} \in S\} \\ \operatorname{argmin}\{f(\mathbf{x}) : \mathbf{x} \in S\}.\end{aligned}$$

Keep in mind there could be several global minimal and maximal points, however minimal and maximal values are unique. This comes from the fact that global minimal and maximal points belong to the domain of the objective function whereas minimal and maximal values belong to the range.

As a trivial example, consider any constant function, $f : \mathbb{R}^n \rightarrow \mathbb{R}$ i.e. $f(\mathbf{x}) = c$. It is clear that the set of all minimizers and maximizers are equivalent and include all points in \mathbb{R}^n . However the maximal (and in this case minimal) value is unique namely, $c \in \mathbb{R}$.

Definition 3.5. \mathbf{x}^* is called the *local minimum point* of f over S if there exists $r > 0$ such that $f(\mathbf{x}) \geq f(\mathbf{x}^*)$ for any $\mathbf{x} \in S \cap B(\mathbf{x}^*, r)$.

Definition 3.6. \mathbf{x}^* is called the *strict local minimum point* of f over S if there exists $r > 0$ such that $f(\mathbf{x}) > f(\mathbf{x}^*)$ for any $\mathbf{x}^* \neq \mathbf{x} \in S \cap B(\mathbf{x}^*, r)$.

Definition 3.7. \mathbf{x}^* is called the *local maximum point* of f over S if there exists $r > 0$ such that $f(\mathbf{x}) \leq f(\mathbf{x}^*)$ for any $\mathbf{x} \in S \cap B(\mathbf{x}^*, r)$.

Definition 3.8. \mathbf{x}^* is called the *strict local maximum point* of f over S if there exists $r > 0$ such that $f(\mathbf{x}) < f(\mathbf{x}^*)$ for any $\mathbf{x}^* \neq \mathbf{x} \in S \cap B(\mathbf{x}^*, r)$.

3.2 First Order Optimality Condition

We know that a for a differentiable one dimensional function $f(x)$ defined over some interval, if there exists some local maximum or minimum x^* , then $f'(x^*) = 0$. The multidimensional generalization states that the gradient of the function will be zero at local optima. We refer to this as the *first order optimality condition* since it relies on constraints on the first derivatives of a function.

Theorem 3.9. (First Order Necessary Optimality Condition). *Let $f : U \rightarrow \mathbb{R}$ be a function defined on the set $U \subseteq \mathbb{R}^n$. Suppose that $\mathbf{x}^* \in \text{int}(U)$ is a local optimum point and that all partial derivatives of f exist at \mathbf{x}^* . Then $\nabla f(\mathbf{x}^*) = \mathbf{0}$.*

Proof. Let $i \in \{1, 2, \dots, n\}$ and consider the one dimensional function $g(t) = f(\mathbf{x}^* + te_i)$. Notice that g is differentiable at $t = 0$ and that,

$$g'(0) = \frac{\partial f}{\partial x_i}(\mathbf{x}^*).$$

Since \mathbf{x}^* is a local optimum of f it follows that $t = 0$ is a local optimum of g which implies that $g'(0) = 0$. So we gain the equality,

$$\frac{\partial f}{\partial x_i}(\mathbf{x}^*) = 0.$$

Furthermore, this will hold for all i in the index set and the claim follows i.e. $\nabla f(\mathbf{x}^*) = \mathbf{0}$. \square

It should be noted the converse of the statement is not true i.e. there exist points that are not local optima that have gradient zero. For example $f(x) = x^3$.

We refer to points in the interior of the set whose gradient is zero as *stationary points*.

3.3 Matrix Classifications

Definition 3.10. (Positive Semidefiniteness).

1. A symmetric matrix $A \in \mathbb{R}^{n \times n}$ is called **positive definite**, denoted $A \succ 0$, if $\mathbf{x}^T A \mathbf{x} > 0$ for all $\mathbf{x} \neq \mathbf{0} \in \mathbb{R}^n$.
2. A symmetric matrix $A \in \mathbb{R}^{n \times n}$ is called **positive semidefinite**, denoted $A \succeq 0$, if $\mathbf{x}^T A \mathbf{x} \geq 0$ for all $\mathbf{x} \in \mathbb{R}^n$.

We may by symmetry derive the definition of negative definite and negative semidefinite matrices from Definition 3.10.

Definition 3.11. (Indefiniteness) A symmetric matrix $A \in \mathbb{R}^{n \times n}$ is called *indefinite* if there exist \mathbf{x} and \mathbf{y} in \mathbb{R}^n such that $\mathbf{x}^T A \mathbf{x} < 0$ and $\mathbf{y}^T A \mathbf{y} > 0$.

It is not easy in general to check whether a matrix is positive definite or positive semidefinite (same for negative). One way to check is to characterize the eigenvalues of these matrices.

Theorem 3.12. (Eigenvalue Characterization Theorem). *Let A be a symmetric $n \times n$ matrix then,*

- (a) A is positive definite if and only if all its eigenvalues are positive.
- (b) A is positive semidefinite if and only if all its eigenvalues are non-negative.
- (c) A is negative definite if and only if all its eigenvalues are negative.
- (d) A is negative semidefinite if and only if all its eigenvalues are non-positive.
- (e) A is indefinite if and only if it has at least one positive eigenvalue and at least one negative eigenvalue.

Moreover, without computer assistance it is still cumbersome to use eigenvalue characterization to determine the definiteness of a matrix as all eigenvalues need to be calculated. As one can tell, without computer assistance calculating eigenvalues for any matrix greater than 3×3 is tedious and error prone.

An alternate method is through *Leading Principal Minors*. Recall some terminology from linear algebra, namely,

- The **leading principal submatrix of order k** of an $n \times n$ matrix \mathbf{A} is obtained by deleting the last $n - k$ rows and columns of \mathbf{A} .
- The determinant of the leading principal submatrix is called the **leading principal minor** of \mathbf{A} .

Theorem 3.13. (*Leading Principal Minors Criterion*) *A matrix is*

- *positive definite if and only if all its leading principal minors are positive.*
- *negative definite if and only if all its odd principal minors are negative and even principal minors are positive.*
- *indefinite if one of its k th order leading principal minors is negative for an even k or if there are two odd leading principal minors that have different signs.*

We may articulate conditions for positive/negative semi-definiteness by changing positive/negative to non-negative/non-positive respectively.

Yet another check of matrix definiteness is by checking if the matrix is diagonally dominant. Such techniques can easily be found in any text on Linear Algebra. The moral of the story here is that in general to check matrix definiteness is a hard computational problem especially in practical problems that arise in applied mathematics. As such there are various techniques one can use to test matrix definiteness.

3.4 Second Order Optimality Conditions

We state the necessary and sufficient second order optimality conditions.

Theorem 3.14. (*necessary second order optimality conditions*). *Let $f : U \rightarrow \mathbb{R}$ be a function defined on the open set $U \subseteq \mathbb{R}^n$. Suppose f is twice continuously differentiable over U and that \mathbf{x}^* is a stationary point. Then the following hold:*

- If \mathbf{x}^* is a local minimum point of f over U , then $\nabla^2 f(\mathbf{x}^*) \succcurlyeq 0$.*
- If \mathbf{x}^* is a local maximum point of f over U , then $\nabla^2 f(\mathbf{x}^*) \preccurlyeq 0$.*

Proof. (a) Since \mathbf{x}^* is a local minimum point of f over U there exists a ball centred at \mathbf{x}^* with radius r i.e. $B(\mathbf{x}^*, r) \subseteq U$ for which $f(\mathbf{x}) \geq f(\mathbf{x}^*)$ for all $\mathbf{x} \in B(\mathbf{x}^*, r)$. Let \mathbf{y} be a non-zero vector in \mathbb{R}^n . For any $0 < \alpha < \frac{r}{\|\mathbf{y}\|}$, we have that,

$$\mathbf{x}_\alpha^* \equiv \mathbf{x}^* + \alpha \mathbf{y} \in B(\mathbf{x}^*, r).$$

Hence, for any such α we also have that,

$$f(\mathbf{x}_\alpha^*) \geq f(\mathbf{x}^*).$$

By the Linear Approximation Theorem, we also know there exists a non-zero vector $\mathbf{z}_\alpha \in [\mathbf{x}^*, \mathbf{x}_\alpha^*]$ such that,

$$f(\mathbf{x}_\alpha^*) - f(\mathbf{x}^*) = \nabla f(\mathbf{x}^*)^T (\mathbf{x}_\alpha^* - \mathbf{x}^*) + \frac{1}{2} (\mathbf{x}_\alpha^* - \mathbf{x}^*)^T \nabla^2 f(\mathbf{z}_\alpha) (\mathbf{x}_\alpha^* - \mathbf{x}^*).$$

Since \mathbf{x}^* is a stationary point, the equation reduces to,

$$f(\mathbf{x}_\alpha^*) - f(\mathbf{x}^*) = \frac{\alpha}{2} \mathbf{y}^T \nabla^2 f(\mathbf{z}_\alpha) \mathbf{y}.$$

It then follows for any $\alpha \in (0, \frac{r}{\|\mathbf{y}\|})$ the inequality $\mathbf{y}^T \nabla^2 f(\mathbf{z}_\alpha) \mathbf{y} \geq 0$ holds. Finally, using the fact that $\mathbf{z}_\alpha \rightarrow \mathbf{x}^*$ as $\alpha \rightarrow 0^+$ and the continuity of the Hessian, we obtain that,

$$\mathbf{y}^T \nabla^2 f(\mathbf{x}^*) \mathbf{y} \geq 0.$$

Since this will hold for any $\mathbf{y} \in \mathbb{R}^n$, the claim is established.

(b) The proof for (b) follows directly by applying the procedure above to $-f$.

□

It is important to keep in mind the distinction between the necessary and sufficient conditions for second order optimality. Theorem 3.14 is valid only when given a stationary point with the information that it is a local max or min of f over the set on which it is defined.

Note that Theorem 3.14 implies the stronger property of strict local optimality.

Definition 3.15. (Saddle Point) Let $f : U \rightarrow \mathbb{R}$ be a function defined on the open set $U \subseteq \mathbb{R}^n$. Suppose that f is continuously differentiable over U . A stationary point \mathbf{x}^* is called a saddle point if it is neither a local minimum point or local maximum point of f over U .

Theorem 3.16. (sufficient second order optimality conditions). Let $f : U \rightarrow \mathbb{R}$ be a function defined on the open set $U \subseteq \mathbb{R}^n$. Suppose f is twice continuously differentiable over U and that \mathbf{x}^* is a stationary point. Then the following hold:

- (a) If $\nabla^2 f(\mathbf{x}^*) \succ 0$, then \mathbf{x}^* is a strict global minimum point of f over U .
- (b) If $\nabla^2 f(\mathbf{x}^*) \prec 0$, then \mathbf{x}^* is a strict global maximum point of f over U .

Proof. We will prove part (a), part (b) follows dually by using the same argument for $-f$. We know that \mathbf{x}^* is a stationary point satisfying $\nabla^2 f(\mathbf{x}^*) \succ 0$. Since the Hessian is continuous it follows that there exists a ball $B(\mathbf{x}^*, r) \subseteq U$ such that $\nabla^2 f(\mathbf{x}^*) \succ 0$ for any $\mathbf{x} \in B(\mathbf{x}^*, r)$. By the Linear Approximation Theorem, we have that for any $\mathbf{x} \in B(\mathbf{x}^*, r)$ there exists a vector $\mathbf{z}_\mathbf{x} \in [\mathbf{x}^*, \mathbf{x}]$ for which,

$$f(\mathbf{x}) - f(\mathbf{x}^*) = \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^T \nabla^2 f(\mathbf{z}_\mathbf{x})(\mathbf{x} - \mathbf{x}^*).$$

Since $\nabla^2 f(\mathbf{z}_\mathbf{x}) \succ 0$, it follows that for any $\mathbf{x} \in B(\mathbf{x}^*, r)$ such that $\mathbf{x} \neq \mathbf{x}^*$ the inequality $f(\mathbf{x}) > f(\mathbf{x}^*)$ holds. The claim follows as this means that \mathbf{x}^* is a strict local minimum of f over U . □

Theorem 3.17. (Sufficient Condition for a Saddle Point) Let $f : U \rightarrow \mathbb{R}$ be a function defined on the open set $U \subseteq \mathbb{R}^n$. Suppose that f is twice continuously differentiable over U and that \mathbf{x}^* is a stationary point. If $\nabla^2 f(\mathbf{x}^*)$ is an indefinite matrix, then \mathbf{x}^* is a saddle point of f over U .

An important issue to consider is whether a function has a global minimizer or maximizer. This is an issue regarding attainment or existence. In order to develop tools to investigate this, we recall some important properties of closed and bounded sets on \mathbb{R}^n .

Theorem 3.18. (Weierstrass' Theorem) Let f be a continuous function defined over a non-empty, compact set $D \subset \mathbb{R}^n$. Then there exists a global minimum of f over D .

Recall from elementary analysis, $D \subset \mathbb{R}^n$ being closed and bounded implies D is compact. So, given f is continuous, we will always be guaranteed that there exists a global optimum in D .

Note, that Weierstrass' Theorem shows the existence of such a point, it does not however say anything about the attainment of the optima. Furthermore, when the set D is not compact, Weierstrass' theorem does not guarantee any attainment of extrema. However we may still recover information from the following definition.

Definition 3.19. (*Coerciveness*) Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuous function. The function is called coercive if,

$$\lim_{\|\mathbf{x}\| \rightarrow \infty} f(\mathbf{x}) = \infty.$$

Coercive functions exemplify the fact that the value of $f(x)$ cannot remain bounded when defined over an unbounded set $A \subseteq \mathbb{R}^n$ and f coercive.

The important and most applied property (for our purposes) of coercive functions is that they always attain their global minima when defined over closed sets. This leads to the following fact.

Theorem 3.20. (*Coercive Functions and Existence of Global Minima*) Let $f : D \subset \mathbb{R}^n \rightarrow \mathbb{R}$ such that f is coercive. Then, $f(x)$ has atleast one global minimum point.

Proof. Pick $\mathbf{x}_0 \in \mathbb{R}^n$ and set $f(\mathbf{x}_0) = M$. By the definition of coerciveness it follows that there exists $R_M \geq 0$ such that,

$$f(\mathbf{x}) \geq f(\mathbf{x}_0), \quad \forall \|\mathbf{x}\| \geq R_M.$$

Observe that the global minimum will not be in the set $\{\mathbf{x} : \|\mathbf{x}\| > R_M\}$. However, since the set $\{\mathbf{x} : \|\mathbf{x}\| \leq R_M\}$ is closed and bounded, by Weierstrass' Theorem we are guaranteed that there is a global minimum. The claim follows as we have established that the global minimum exists. \square

3.5 Global Optimality Conditions

Theorem 3.21. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a twice continuously differentiable function such that $\nabla^2 f(\mathbf{x}) \succ 0 \forall \mathbf{x} \in \mathbb{R}^n$. Let \mathbf{x}^* be a stationary point of $f(\mathbf{x})$, then \mathbf{x}^* is a global minimum of f over \mathbb{R}^n .

We may derive the definition of global maximum by a symmetric argument. Given that the Hessian is negative semi-definite $\forall \mathbf{x} \in \mathbb{R}^n$, then a stationary point \mathbf{x}^* is a global maximum.

Note here the crucial quantifier in Theorem 3.21, "for all \mathbf{x} in \mathbb{R}^n ". In other words, given that we know the Hessian of the objective function is positive semi-definite for all \mathbf{x} in \mathbb{R}^n and we have a stationary point \mathbf{x}^* , only then can we classify \mathbf{x}^* as a global minimum point of f .

3.6 Quadratic Functions

Definition 3.22. (*Quadratic Function*) A quadratic function over \mathbb{R}^n is a function of the form,

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + 2\mathbf{b} \mathbf{x}^T + c,$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric, $\mathbf{b} \in \mathbb{R}^n$ and $c \in \mathbb{R}$.

We also gain the helpful identities that are straightforward to prove,

$$\begin{aligned} \nabla f(\mathbf{x}) &= 2\mathbf{A} \mathbf{x} + 2\mathbf{b} \\ \nabla^2 f(\mathbf{x}) &= 2\mathbf{A}. \end{aligned}$$

We may also deduce the following important properties of quadratic functions.

Lemma 3.23. Let $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + 2\mathbf{b} \mathbf{x}^T + c$ such that $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric, $\mathbf{b} \in \mathbb{R}^n$ and $c \in \mathbb{R}$. Then,

- (a) \mathbf{x} is a stationary point of f if and only if $\mathbf{A} \mathbf{x} = -\mathbf{b} \implies \mathbf{x} = -\mathbf{A}^{-1} \mathbf{b}$.
- (b) If $\mathbf{A} \succcurlyeq 0$, then \mathbf{x} is a global minimum point of f if and only if $\mathbf{A} \mathbf{x} = -\mathbf{b}$.
- (c) If $\mathbf{A} \succ 0$, then $\mathbf{x} = -\mathbf{A}^{-1} \mathbf{b}$ is a strict global minimum point of f .

Following from (c), we may also gain the minimal value of f by $f(\mathbf{x}) = c - \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b}$.

Lemma 3.23 follows directly as a consequence of the Global Optimality Conditions.

Lemma 3.24. *Let $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + 2\mathbf{b}^T \mathbf{x} + c$ such that $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric, $\mathbf{b} \in \mathbb{R}^n$ and $c \in \mathbb{R}$. Then f is coercive if and only if $\mathbf{A} \succ 0$.*

4 Least Squares

4.1 Solution to Over Determined Systems

Consider the linear system $\mathbf{A} \mathbf{x} \approx \mathbf{b}$, such that $\mathbf{A} \in \mathbb{R}^{m \times n}$. We assume that \mathbf{A} has a full column rank.

When $m > n$ i.e. there are more equations than unknowns, the system is inconsistent. The classic approach to solving such kinds of problems is through **Least Squares**. Solutions of the form,

$$\mathbf{x}_{LS} = \min \|\mathbf{A} \mathbf{x} - \mathbf{b}\|^2.$$

Alternatively, the least squares problem is the same as,

$$\min_{\mathbf{x} \in \mathbb{R}^n} \{f(\mathbf{x}) \equiv \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} - 2\mathbf{b}^T \mathbf{A} \mathbf{x} + \|\mathbf{b}\|^2\}.$$

Recall that $\nabla^2 = 2\mathbf{A}^T \mathbf{A} \prec 0$.

So the unique, optimal solution \mathbf{x}_{LS} is the solution $\nabla f(\mathbf{x}) = \mathbf{0}$,

$$\left(\mathbf{A}^T \mathbf{A}\right) \mathbf{x}_{LS} = \mathbf{A}^T \mathbf{b} \implies \mathbf{x}_{LS} = \left(\mathbf{A}^T \mathbf{A}\right)^{-1} \mathbf{A}^T \mathbf{b}.$$

4.2 Data Fitting

An important application of least squares is in linear data fitting. Consider data points (\mathbf{s}_i, t_i) where $i = 1, 2, 3, \dots, m$ where $\mathbf{s}_i \in \mathbb{R}^n$ and $t_i \in \mathbb{R}$. Assume a linear relation of the form,

$$t_i = \mathbf{s}_i^T \mathbf{x},$$

holds in an approximate sense. In least squares we want to find the parameters vector $\mathbf{x} \in \mathbb{R}^n$ that solves,

$$\min_{\mathbf{x} \in \mathbb{R}^n} \sum_{i=1}^m (\mathbf{s}_i^T \mathbf{x} - t_i)^2 = \min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{S} \mathbf{x} - \mathbf{t}\|^2.$$

4.2.1 Least Squares in Non-linear Data Fitting

The least squares approach can also be used in non-linear fitting. Suppose we are given a set of points in \mathbb{R}^2 of the form (x_i, y_i) , where $i = 1, 2, \dots, m$. Furthermore, we know apriori that these points are approximately related via a polynomial of degree at most d ,

$$\sum_{i=0}^d a_j u_i^j \approx y_i.$$

The least squares approach here seeks a_0, a_1, \dots, a_d that are the least squares solution to the linear system,

$$\underbrace{\begin{pmatrix} 1 & u_1 & u_1^2 & \cdots & u_1^d \\ 1 & u_2 & u_2^2 & \cdots & u_2^d \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & u_m & u_m^2 & \cdots & u_m^d \end{pmatrix}}_{\mathbf{U}} \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_d \end{pmatrix} = \begin{pmatrix} y_0 \\ y_1 \\ \vdots \\ y_m \end{pmatrix}$$

Notice then, the solution is well-defined if the $m \times (d + 1)$ matrix \mathbf{U} is of full column rank. The matrix \mathbf{U} consists of the first $d + 1$ columns of the vandermonde matrix.

4.3 Regularized Least Squares

There are many instances where the least squares solution does not give rise to a good enough approximate solution. For example in systems where there are more variables than equations.

In such cases we must modify the least squares technique to incorporate a *regularization function* $R(\cdot)$. The problem takes the form.

$$\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|^2 + \lambda R(\mathbf{x}).$$

$\lambda > 0$ and is called the *regularization parameter*. As λ gets larger, more weight is given to the regularization function.

In many cases $R(\mathbf{x})$ is taken to be quadratic i.e. $R(\mathbf{x}) = \|\mathbf{Dx}\|^2$ such that $\mathbf{D} \in \mathbb{R}^{p \times n}$ is a given matrix. The aim of the regularization function is to then control the norm of \mathbf{Dx} .

$$\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|^2 + \lambda \|\mathbf{Dx}\|^2.$$

Expressing the problem equivalently as,

$$\min_{\mathbf{x}} \{f_{RLS}(\mathbf{x}) \equiv \mathbf{x}^T (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{D}^T \mathbf{D}) \mathbf{x} - 2\mathbf{b}^T \mathbf{Ax} + \|\mathbf{b}\|^2\}.$$

Since f_{RLS} is quadratic the Hessian is,

$$\nabla^2 f_{RLS}(\mathbf{x}) = 2 (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{D}^T \mathbf{D}) \succcurlyeq 0.$$

So, an stationary point is a global minimum. Furthermore, the stationary points satisfy $\nabla f_{RLS} = \mathbf{0}$ i.e.

$$(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{D}^T \mathbf{D}) \mathbf{x} = \mathbf{A}^T \mathbf{b}.$$

Hence, if $\mathbf{A}^T \mathbf{A} + \lambda \mathbf{D}^T \mathbf{D} \succ 0$, then the RLS solution is given by,

$$\mathbf{x}_{RLS} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{D}^T \mathbf{D})^{-1} \mathbf{A}^T \mathbf{b}.$$

4.4 Denoising

A common area where RLS is implemented is in denoising data. Suppose that a a noisy signal of a measurement is given as,

$$\mathbf{b} = \mathbf{x} + \mathbf{w}.$$

Here $\mathbf{x} \in \mathbb{R}^n$ is an unknown signal and $\mathbf{w} \in \mathbb{R}^n$ is noise. Given \mathbf{b} , we want to find a good estimate for \mathbf{x} . However, the problem,

$$\min \|\mathbf{b} - \mathbf{x}\|^2,$$

is meaningless as the optimal solution is clearly $\mathbf{b} = \mathbf{x}$.

In order to articulate a more relevant problem, we must add a regularization term. Supposing we have some a-priori information regarding the signal. For example, we know that is it smooth in some sense. Then, our regularization function can be quadratic. Specifically, the sum of squares of consecutive differences i.e.

$$R(\mathbf{x}) = \sum_{i=1}^{n-1} (x_i - x_{i+1})^2.$$

We may also write $R(\mathbf{x})$ in matrix form as $\|\mathbf{Lx}\|^2$, where $\mathbf{L} \in \mathbb{R}^{(n-1) \times n}$ is given as,

$$\mathbf{L} = \begin{pmatrix} 1 & -1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 1 & -1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & -1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 1 & -1 \end{pmatrix}$$

Therefore the resulting RLS problem is,

$$\min \|\mathbf{x} - \mathbf{b}\|^2 + \lambda \|\mathbf{L}\mathbf{x}\|^2.$$

Furthermore, the optimal solution is given as,

$$\mathbf{x}_{RLS}(\lambda) = (\mathbf{I} + \lambda \mathbf{L}^T \mathbf{L})^{-1} \mathbf{b}.$$

4.5 Non-Linear Least Squares

So far the least squares methods we have outlined are for solution to approximate linear equalities. As we know however, non-linear systems arise ubiquitously in practice.

$$f_i(\mathbf{x}) = c_i.$$

In such a case, we want to formulate a non-linear least squares problem of the form,

$$\min \sum_{i=1}^n (f_i(\mathbf{x}) - c_i)^2.$$

As opposed to linear least squares, there is no general method for solving non-linear least squares problems. The Gauss-Newton Method (which will be discussed later) can be used to solve problems of such form, however this does not guarantee the optimal solution. Rather, only a solution that is stationary.

4.6 Circle Fitting

Suppose we are given m points $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m \in \mathbb{R}^n$. The circle fitting problem seeks to find a circle,

$$C(\mathbf{x}, r) = \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{x}\| = r\},$$

that best fits the m points. Note, we are using the term circle, however our procedure is being constructed in general on \mathbb{R}^n . Furthermore, observe that $C(\mathbf{x}, r)$ is the boundary of the corresponding ball $B(\mathbf{x}, r)$.

The non-linear approximate equations associated with the problem are given as,

$$\|\mathbf{x} - \mathbf{a}_i\| \approx r, \quad i = 1, 2, \dots, m.$$

We wish to deal with differentiable functions and the norm function is not differentiable. Hence,

$$\|\mathbf{x} - \mathbf{a}\|^2 \approx r^2.$$

The NLS problem here can be formulated as follows,

$$\min_{\mathbf{x} \in \mathbb{R}^n, r \in \mathbb{R}} \sum_{i=1}^m (\|\mathbf{x} - \mathbf{a}_i\|^2 - r^2)^2.$$

The trick here is to realize that this NLS problem can be shown to be equivalent to a linear LS problem. It would follow that the optimal solution can be easily obtained.

Note, the problem is the same as,

$$\min_{\mathbf{x}, r} \left\{ \sum_{i=1}^m \left(-2\mathbf{a}_i^T \mathbf{x} + \|\mathbf{x}\|^2 - r^2 + \|\mathbf{a}_i\|^2 \right)^2 : \mathbf{x} \in \mathbb{R}^n, r \in \mathbb{R} \right\}.$$

Where we set the regularization function, $R(\mathbf{x}) = \|\mathbf{x}\|^2 - r^2$. Then re-write the above problem as,

$$\min_{\mathbf{x} \in \mathbb{R}^n, R \in \mathbb{R}} \left\{ f(\mathbf{x}, R) \equiv \sum_{i=1}^m \left(-2\mathbf{a}_i^T \mathbf{x} + R + \|\mathbf{a}_i\|^2 \right)^2 : \|\mathbf{x}\|^2 \geq R \right\}.$$

Notice, that the change of variable introduced a constraint namely, $\|\mathbf{x}\|^2 \geq R$. We want to show that this constraint can be dropped and the problem is equivalent to the LS,

$$\min_{\mathbf{x}, R} \left\{ \sum_{i=1}^m \left(-2\mathbf{a}_i^T \mathbf{x} + R + \|\mathbf{a}_i\|^2 \right)^2 : \mathbf{x} \in \mathbb{R}^n, R \in \mathbb{R} \right\}.$$

5 Gradient Methods

5.1 Descent Direction Methods

We now turn to looking at the general unconstrained minimization problem,

$$\min\{f(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^n\}.$$

We have already seen in the previous section that one procedure to solve such problems is by computing the stationary points i.e. solutions $\nabla f(\mathbf{x}) = \mathbf{0}$ and then finding the stationary point with lowest function value.

In general however this will not work. This is due to several reasons, here are a couple:

- (i) In a majority of cases it is very difficult to solve the system of mainly non-linear $\nabla f(\mathbf{x}) = \mathbf{0}$.
- (ii) We may have infinitely many stationary points and finding the one that corresponds to the minimal value may be a difficult optimization problem in and of itself. This could be even more difficult than the original problem!

It is easy to see that finding exact analytical solutions will be (in general) a very difficult task.

Hence using some iterative numerical technique may be an easy way to find a numerical/computational solution. The iterative algorithm we will consider takes the form,

$$\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k, \quad k = 0, 1, 2, \dots$$

Here $\mathbf{d}_k \in \mathbb{R}^n$ is the direction and t_k is the stepsize. In this method, we restrict ourselves, by convention, to descent directions.

Definition 5.1. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuously differentiable function over \mathbb{R}^n . Let a non-zero vector $\mathbf{d} \in \mathbb{R}^n$ be called a descent direction of f at \mathbf{x} if the directional derivative is negative. Mathematically,

$$f'(\mathbf{x}; \mathbf{d}) = \nabla f(\mathbf{x})^T \cdot \mathbf{d} < 0.$$

It is straightforward to see that taking small enough step size, moving in the direction of descent leads to a decrease in objective function value. We may state this more rigorously as a Lemma.

Lemma 5.2. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuously differentiable function and let $\mathbf{x} \in \mathbb{R}^n$. Suppose that \mathbf{d} is a descent direction of f at \mathbf{x} . Then there exists $\epsilon > 0$ such that

$$f(\mathbf{x} + t\mathbf{d}) < f(\mathbf{x}), \quad \forall t \in (0, \epsilon].$$

The proof of Lemma 5.2 follows directly from the definition of the directional derivative.

Here is a schematic on how the Gradient method works within a computation.

Initialization: Pick $\mathbf{x}_0 \in \mathbb{R}^n$ arbitrarily.

General Step: For any $k = 0, 1, 2, \dots$

- (a) Pick a descent direction \mathbf{d}_k .
- (b) Find a step size t_k such that $f(\mathbf{x}_k + t_k \mathbf{d}_k) < f(\mathbf{x}_k)$.
- (c) Set $\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k$.
- (d) If stopping criterion is satisfied, then STOP and the output is \mathbf{x}_{k+1} .

There are still many details missing, hence the schematic and by extension the method remains conceptual.

We must still have rigorous answers to the questions what is the stopping criteria? What stepsize should be taken? What is the starting point? How does one choose a descent direction?

We begin by analyzing stopping criteria. A popular stopping criterion is given by,

$$\|\nabla f(\mathbf{x}_{k+1})\| \leq \epsilon.$$

We assume the stepsize is chosen in such a way that $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$. Notice then we have assumed a descent direction as the function value decreases iteration to iteration.

5.1.1 Step-Size Selection

The process of finding the step size is called **line search**. This can be regarded as a minimization problem in and of itself as we try and minimize the one dimensional function,

$$g(t) = f(\mathbf{x}_k + t_k \mathbf{d}_k).$$

There are many choices for step size selection. We will describe three major ones now.

- **Constant Step-Size:** $t_k = t$ for all k .
- **Exact Line Search:** t_k is a minimizer of f along the ray $\mathbf{x}_k + t_k \mathbf{d}_k$.

$$t_k \in \operatorname{argmin}_{t \geq 0} f(\mathbf{x}_k + t_k \mathbf{d}_k).$$

- **Backtracking:** This method requires three parameters, namely $s > 0$, $\alpha \in (0, 1)$ and $\beta \in (0, 1)$. The choice of t_k is taken by the following procedure. First, t_k is set to equal the initial s . Then, while,

$$f(\mathbf{x}_k) - f(\mathbf{x}_k + t_k \mathbf{d}_k) < -\alpha t_k \nabla f(\mathbf{x}_k)^T \mathbf{d}_k$$

we set $\beta t_k \rightarrow t_k$. In words, we want to find the smallest step-size such that the condition above is satisfied.

We can begin to see that choosing a step-size is a non-trivial task and in practice will be more of an art than an exact method while solving gradient descent problems.

For constant step size, we know that its advantage is its simplicity, however the question of how to choose a constant remains ambiguous. It then seems as though line search becomes the perfect solution. However, in practice choosing the minimizer of f along the ray may cause long convergence times and computing inefficiencies. Also it may not always be possible to find an exact minimizer.

In a sense then, backtracking becomes the mathematical compromise between the two other methods in step size selection. It does not perform an exact line search but does generate a "good enough" step size. Rigorously, good-enough here means the condition as outlined above is satisfied.

Lemma 5.3. Let f be a continuously differentiable function on \mathbb{R}^n and let $\mathbf{x} \in \mathbb{R}^n$. Suppose that $\mathbf{d} \neq \mathbf{0} \in \mathbb{R}^n$ is a descent direction of f at \mathbf{x} and let $\alpha \in (0, 1)$. Then there exists $\epsilon > 0$ such the inequality,

$$f(\mathbf{x}) - f(\mathbf{x} + t\mathbf{d}) \geq -\alpha t \nabla f(\mathbf{x})^T \mathbf{d},$$

Holds for all $t \in [0, \epsilon]$.

We now show an example to find the analytical formula for the exact line search when the objective function is quadratic.

Example 5.4. Let $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + 2\mathbf{b}^T \mathbf{x} + c$ where $\mathbf{A} \in \mathbb{R}^{n \times n}$ and positive definite, $\mathbf{b} \in \mathbb{R}^n$ and $c \in \mathbb{R}$. Let $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{d} \in \mathbb{R}^n$ be a descent direction of f at \mathbf{x} .

We will find the exact analytical formula for finding the step size as generated by line search. Explicitly, we will find a solution for the problem,

$$\min_{t \geq 0} \{f(\mathbf{x} + t\mathbf{d})\}.$$

So, we have that,

$$\begin{aligned} g(t) = f(\mathbf{x} + t\mathbf{d}) &= (\mathbf{x} + t\mathbf{d})^T \mathbf{A} (\mathbf{x} + t\mathbf{d}) + 2\mathbf{b}^T (\mathbf{x} + t\mathbf{d}) + c \\ &= (\mathbf{d}^T \mathbf{A} \mathbf{d}) t^2 + 2(\mathbf{d}^T \mathbf{A} \mathbf{x} + \mathbf{d}^T \mathbf{b}) t + \mathbf{x}^T \mathbf{A} \mathbf{x} + 2\mathbf{b}^T \mathbf{x} + c \\ &= (\mathbf{d}^T \mathbf{A} \mathbf{d}) t^2 + 2(\mathbf{d}^T \mathbf{A} \mathbf{x} + \mathbf{d}^T \mathbf{b}) t + f(\mathbf{x}). \end{aligned}$$

Note, in the second step we use the fact that $\mathbf{A} \succ 0$ and hence symmetric. Now, we know that,

$$\begin{aligned} g'(t) &= 2(\mathbf{d}^T \mathbf{A} \mathbf{d}) t + 2(\mathbf{d}^T \mathbf{A} \mathbf{x} + \mathbf{d}^T \mathbf{b}) \\ \nabla f(\mathbf{x}) &= 2(\mathbf{A} \mathbf{x} + \mathbf{b}). \end{aligned}$$

It then must be the case that $g'(t) = 0$ if and only if,

$$t = \bar{t} \equiv -\frac{\mathbf{d}^T \nabla f(\mathbf{x})}{2\mathbf{d}^T \mathbf{A} \mathbf{d}}.$$

Since \mathbf{d} is a descent direction of f at \mathbf{x} , it follows that $\mathbf{d}^T \nabla f(\mathbf{x}) < 0$ and hence $\bar{t} > 0$. This implies that the step size dictated by the exact line search is,

$$\bar{t} = -\frac{\mathbf{d}^T \nabla f(\mathbf{x})}{2\mathbf{d}^T \mathbf{A} \mathbf{d}}.$$

Note, we have implicitly used the fact that the second derivative with respect to $g(t)$ is always positive.

5.2 Gradient Descent Method

We now turn towards outlining the gradient descent method or just gradient descent. Building from the theory thus far, the gradient descent method takes the descent direction to be the negative gradient of the objective function i.e $\mathbf{d}_k = -\nabla f(\mathbf{x}_k)$. By definition, this will give us the direction of steepest descent.

As done earlier, we outline the schematic of the gradient descent method program.

Input: $\epsilon > 0$, tolerance parameter.

Initialization: Pick $\mathbf{x}_0 \in \mathbb{R}^n$ arbitrarily.

General Step: For $k : 0, 1, 2, \dots$

- (a) Pick step size via line search on the function $g(t) = f(\mathbf{x}_k - t\nabla f(\mathbf{x}_k))$.
- (b) Set $\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \nabla f(\mathbf{x}_k)$.
- (c) If $\|\nabla f(\mathbf{x}_{k+1})\| < \epsilon$, STOP and \mathbf{x}_{k+1} is the output.

We now turn to looking at an example of an exact line search within the gradient descent method as implemented on MATLAB.

Example 5.5. Consider the two dimensional minimization problem,

$$\min_{x,y} x^2 + 2y^2.$$

It is easy to see that there is only one minimizer and it is global namely $(0,0)$, where the function value is 0. To invoke the gradient method with the step size chosen by exact line search we write a MATLAB script as provided below.

The program finds the exact step size via exact line search to problems of the form,

$$\min_{\mathbf{x} \in \mathbb{R}^n} \{ \mathbf{x}^T \mathbf{A} \mathbf{x} + 2\mathbf{b}^T \mathbf{x} \}.$$

Where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is $\mathbf{A} \succ 0$ and $\mathbf{b} \in \mathbb{R}^n$. By Example 5.5, we know the exact line search formula for quadratic minimization problem. Hence, in this case we know the step size is given by the k th iteration,

$$t_k = \frac{\|\nabla f(\mathbf{x}_k)\|^2}{2\nabla f(\mathbf{x}_k)^T \mathbf{A} \nabla f(\mathbf{x}_k)}.$$

Keep in mind this is a gradient method and hence the direction of descent is taken to be the negative gradient of the objective function.

We now present the example of applying the gradient method via backtracking in order to select the step size.

5.3 The Condition Number

Consider the quadratic minimization problem,

$$\min_{\mathbf{x} \in \mathbb{R}^n} \{ f(\mathbf{x}) \equiv \mathbf{x}^T \mathbf{A} \mathbf{x} \}.$$

Such that $\mathbf{A} \succ 0$. The optimal solution is obviously $\mathbf{x}^* = \mathbf{0}$. We know that the gradient method with exact line search takes the form,

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t_k \mathbf{d}_k,$$

where $\mathbf{d}_k = \nabla f(\mathbf{x}_k) = 2\mathbf{A}\mathbf{x}_k$. Moreover, we know the step size is given by,

$$t_k = \frac{\mathbf{d}_k^T \mathbf{d}_k}{2\mathbf{d}_k^T \mathbf{A} \mathbf{d}_k}.$$

Then we have that

$$\begin{aligned} f(\mathbf{x}_{k+1}) &= \mathbf{x}_{k+1}^T \mathbf{A} \mathbf{x}_{k+1} \\ &= (\mathbf{x}_k - t_k \mathbf{d}_k)^T \mathbf{A} (\mathbf{x}_k - t_k \mathbf{d}_k) \\ &= \mathbf{x}_k^T \mathbf{A} \mathbf{x}_k - 2t_k \mathbf{d}_k^T \mathbf{A} \mathbf{x}_k + t_k^2 \mathbf{d}_k^T \mathbf{A} \mathbf{d}_k \\ &= \mathbf{x}_k^T \mathbf{A} \mathbf{x}_k - t_k \mathbf{d}_k^T \mathbf{d}_k + t_k^2 \mathbf{d}_k^T \mathbf{A} \mathbf{d}_k \end{aligned}$$

Now, substituting the expression for t_k into the last step,

$$\begin{aligned} f(\mathbf{x}_{k+1}) &= \mathbf{x}_k^T \mathbf{A} \mathbf{x}_k - \frac{1}{4} \frac{(\mathbf{d}_k^T \mathbf{d}_k)^2}{\mathbf{d}_k^T \mathbf{A} \mathbf{d}_k} = \mathbf{x}_k^T \mathbf{A} \mathbf{x}_k \left(1 - \frac{1}{4} \frac{(\mathbf{d}_k^T \mathbf{d}_k)^2}{(\mathbf{d}_k^T \mathbf{A} \mathbf{d}_k) (\mathbf{x}_k^T \mathbf{A} \mathbf{A}^{-1} \mathbf{A} \mathbf{x}_k)} \right) \\ &= \left(1 - \frac{(\mathbf{d}_k^T \mathbf{d}_k)^2}{(\mathbf{d}_k^T \mathbf{A} \mathbf{d}_k) (\mathbf{d}_k^T \mathbf{A}^{-1} \mathbf{d}_k)} \right) f(\mathbf{x}_k) \end{aligned}$$

We now turn to a well known result from matrix theory called Kantorovich's inequality.

Lemma 5.6. Let \mathbf{A} be a real, square, positive definite matrix. Then for any $\mathbf{x} \neq \mathbf{0} \in \mathbb{R}^n$, the inequality,

$$\frac{(\mathbf{x}^T \mathbf{x})^2}{(\mathbf{x}^T \mathbf{A} \mathbf{x})(\mathbf{x}^T \mathbf{A}^{-1} \mathbf{x})} \geq \frac{4\lambda_{\max}(\mathbf{A})\lambda_{\min}(\mathbf{A})}{(\lambda_{\max}(\mathbf{A}) + \lambda_{\min}(\mathbf{A}))^2},$$

holds.

We want to analyse the convergence rate of the gradient method on the quadratic minimization problem. By using the Kantorovich inequality we have that,

$$f(\mathbf{x}_{k+1}) \leq \left(1 - \frac{4Mm}{(M+m)^2}\right) f(\mathbf{x}_k) = \left(\frac{M-m}{M+m}\right)^2 f(\mathbf{x}_k).$$

Where $M = \lambda_{\max}(\mathbf{A})$ and $m = \lambda_{\min}(\mathbf{A})$. We may summarize the ideas here in the following lemma.

Lemma 5.7. Let $\{\mathbf{x}_k\}_{k \geq 0}$ be the sequence generated by the gradient descent method via exact line search for solving the quadratic minimization problem. Then for any k ,

$$f(\mathbf{x}_{k+1}) \leq \left(\frac{M-m}{M+m}\right)^2 f(\mathbf{x}_k).$$

This inequality implies that, $f(\mathbf{x}_k) \leq c^k f(\mathbf{x}_0)$, where $c = \left(\frac{M-m}{M+m}\right)^2$.

We have that the sequence of function values is bounded above by a decreasing geometric sequence. In this case we say that the sequence of a function values converges at a linear rate to the optimal value.

There is great dependence on c , as such we may deduce convergence rates by analysing c . **As c becomes larger, convergence speed becomes slower.** We may also write,

$$c = \left(\frac{x-1}{x+1}\right)^2.$$

Where $x = \frac{M}{m}$ i.e. the ratio between the largest eigenvalue of \mathbf{A} and the smallest eigenvalue of \mathbf{A} .

Definition 5.8. Let \mathbf{A} be an $n \times n$ positive definite matrix. Then the **condition number** of \mathbf{A} is defined by,

$$\chi(\mathbf{A}) = \frac{\lambda_{\max}(\mathbf{A})}{\lambda_{\min}(\mathbf{A})}.$$

We may in general define the condition number for any square matrix, however we will restrict ourselves to the positive definite case for our purposes with respect to optimization.

One will notice that when gradient method take a relatively small number of iterations to converge to the optimal solution, the condition number is small.

Matrices with a small condition number are said to be **well-conditioned**, whereas matrices with large condition number are said to be **ill-conditioned**.

We have restricted ourselves to a class of problems where the objective function is quadratic, hence the Hessian is constant and therefore we gain well defined condition numbers. However, the same analysis will work for non-quadratic objective functions as well.

Convergence Analysis of the Gradient Method

We now turn to providing a more rigorous treatment regarding convergence of the gradient method.

Lipschitz Property of the Gradient

Consider the unconstrained optimization problem,

$$\min \{f(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^n\}.$$

We assume that the objective function is continuously differentiable and that the gradient is Lipschitz continuous i.e.

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \text{ for any } \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$$

Note that if $\nabla f(\mathbf{x})$ is Lipschitz continuous with constant L , then it is also Lipschitz continuous with constant $\tilde{L} \geq L$.

The class of functions with Lipschitz gradient with constant L is denoted $C_L^{1,1}(\mathbb{R}^n)$ or just $C_L^{1,1}$. Furthermore, linear functions and quadratic functions fall in the class of $C^{1,1}$.

Theorem 5.9. *Let f be a twice continuously differentiable function over \mathbb{R}^n . Then the following two claims are equivalent:*

- (a) $f \in C_L^{1,1}(\mathbb{R}^n)$.
- (b) $\|\nabla^2 f(\mathbf{x})\| \leq L$ for any $\mathbf{x} \in \mathbb{R}^n$.

In words, Theorem 5.9 says that a twice continuously differentiable function has Lipschitz gradient if and only if it is equivalent to saying that the norm on its Hessian is bounded above by the Lipschitz constant.

Also note that the norm of the Hessian is an induced matrix norm and will generally (for our purposes) be the spectral norm.

The Descent Lemma

An important result for all $C^{1,1}$ functions is that they can be bounded above by quadratic functions over the entire space. This result is known as the **Descent Lemma** and is of paramount importance when doing analysis regarding convergence of gradient based methods.

Lemma 5.10. *Let $f \in C_L^{1,1}(\mathbb{R}^n)$. Then for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,*

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

Proof.

□

Convex Sets

We now begin the theoretical discussion of convex sets and convex analysis. The theoretical foundations we build will then be applied within the context of a large sub-field of optimization called convex optimization.

Definition 5.11. Convex Sets: *Let $S \subseteq \mathbb{R}^n$, S is said to be convex if for any $\mathbf{x}, \mathbf{y} \in S$,*

$$\lambda\mathbf{x} + (1 - \lambda)\mathbf{y} \in S,$$

for any $\lambda \in [0, 1]$.

Geometrically, this definition states that a set is convex only when for every pair of points within the set, those two points can be connect by a line segment and the entire line segment stays within the set.

As a trivial example, a line segment is a convex set. Consider the following lemma on more important sets that turn out to be convex.

Lemma 5.12. *Let $\mathbf{a} \in \mathbb{R}^N \setminus \{\mathbf{0}\}$ and $b \in \mathbb{R}^n$. Then the following sets are convex,*

(a) *The hyper-plane $H = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^T \mathbf{x} = b\}$.*

(b) *The half-space $H^- = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^T \mathbf{x} \leq b\}$.*

(c) *The open-half-space $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^T \mathbf{x} < b\}$.*

Proof. Let us prove the convexity of the half-space, the proof for the other two can be extended trivially. Following from the definition of convexity, let $\mathbf{x}, \mathbf{y} \in H^-$ and $\lambda \in [0, 1]$. We want to show that $\lambda\mathbf{x} + (1-\lambda)\mathbf{y} \in H^-$. Take $\mathbf{z} \in H^-$ such that,

$$\mathbf{a}^T \mathbf{z} = \mathbf{a}^T (\lambda\mathbf{x} + (1-\lambda)\mathbf{y}) = \lambda\mathbf{a}^T \mathbf{x} + (1-\lambda)\mathbf{a}^T \mathbf{y} \leq \lambda b + (1-\lambda)b.$$

Where the inequality follows from the definition of the half space i.e. $\mathbf{a}^T \mathbf{x} \leq b$ and $\mathbf{a}^T \mathbf{y} \leq b$. Hence the claim follows. \square

To establish convexity of more topological sets, we can also show that open and closed balls are convex.

Lemma 5.13. *Let $\mathbf{c} \in \mathbb{R}^n$ and $r > 0$ and let $\|\cdot\|$ be some norm on \mathbb{R}^n . Then, the closed and open balls,*

$$B(\mathbf{c}, r) = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{c} - \mathbf{x}\| < r\}, B[\mathbf{c}, r] = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{c} - \mathbf{x}\| \leq r\},$$

are convex.

Proof. We prove the convexity of the closed ball, the proof can be trivially extended to open balls. Let $\mathbf{x}, \mathbf{y} \in B[\mathbf{c}, r]$ and $\lambda \in [0, 1]$, it follows that,

$$\|\mathbf{c} - \mathbf{x}\| \leq r, \|\mathbf{c} - \mathbf{y}\| \leq r.$$

Now let $\mathbf{z} = \lambda\mathbf{x} + (1-\lambda)\mathbf{y}$, we want to show that $\mathbf{z} \in B[\mathbf{c}, r]$. Take,

$$\begin{aligned} \|\mathbf{z} - \mathbf{c}\| &= \|\lambda\mathbf{x} + (1-\lambda)\mathbf{y} - \mathbf{c}\| \\ &= \|\lambda(\mathbf{x} - \mathbf{c}) + (1-\lambda)(\mathbf{y} - \mathbf{c})\| \\ &\leq \|\lambda(\mathbf{x} - \mathbf{c})\| + \|(1-\lambda)(\mathbf{y} - \mathbf{c})\| \\ &= \lambda\|\mathbf{x} - \mathbf{c}\| + (1-\lambda)\|\mathbf{y} - \mathbf{c}\| \\ &\leq \lambda r + (1-\lambda)r = r. \end{aligned}$$

Note, the first inequality comes from the triangle inequality. The claim follows. \square

Note, that these results will work for all norms on \mathbb{R}^n .

Algebraic Operations on Convex Sets

An important property of convexity is that it is preserved under intersections of sets.

Lemma 5.14. *Let $C_i \in \mathbb{R}^n$ where i is in some index set I . If each C_i is convex then,*

$$\bigcap_i C_i,$$

is also convex.

Proof. Let $\mathbf{x}, \mathbf{y} \in \bigcap_i C_i$ and let $\lambda \in [0, 1]$ and assume that C_i is convex for any i . It then follows from the definition of convexity that $\lambda\mathbf{x} + (1 - \lambda)\mathbf{y} \in C_i$ for any i . By the definition of intersection of sets, it must be the case that,

$$\lambda\mathbf{x} + (1 - \lambda)\mathbf{y} \in \bigcap_i C_i.$$

The claim follows. □

In a similar manner we can prove that convexity is preserved under addition, cartesian product, linear mappings and inverse linear mappings. These results also show that convexity is preserved under translations, rotations and other symmetries.

Note, that convexity is preserved under set intersections however this is not the case for set unions. For example, take two distinct closed balls in \mathbb{R}^2 . The union of these two sets is clearly not convex as there is not line segment that connect one point from the first ball to the another point in the second ball while remaining within the set union.

Definition 5.15. Convex Combinations *Given k vectors $\mathbf{x}_1, \dots, \mathbf{x}_k$, the convex combination of these vectors is,*

$$\sum_{i=1}^k \lambda_i \mathbf{x}_i,$$

such that $\sum_{i=1}^k \lambda_i = 1$ and each λ_i is non-negative.

We can see how then how the definition of the convex set originates. Moreover, we may re-define a convex set to be one, where the convex combinations of vectors within some set remain within the set.

In fact for any m vectors in a convex set, we can prove via induction that the convex combination of these m vectors remains within the set.

Theorem 5.16. (Countable Convex Combinations) *Let C be a convex set and let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m \in C$. Then for any $\lambda \in \Delta_m$,*

$$\sum_{i=1}^m \lambda_i \mathbf{x}_i \in C.$$

Proof. We will show this by induction on m . For the base case $m = 1$, the result is trivial as $\mathbf{x}_1 \in C$.

The induction hypothesis is that $\sum_{i=1}^m \lambda_i \mathbf{x}_i \in C$.

We now want to show that $\mathbf{z} := \sum_{i=1}^{m+1} \lambda_i \mathbf{x}_i \in C$. If $\lambda_{m+1} = 1$ we have that $\mathbf{z} = \mathbf{x}_{m+1} \in C$ and the claim follows. If $\lambda_{m+1} < 1$, then,

$$\begin{aligned} \mathbf{z} &= \sum_{i=1}^{m+1} \lambda_i \mathbf{x}_i \\ &= \sum_{i=1}^m \lambda_i \mathbf{x}_i + \lambda_{m+1} \mathbf{x}_{m+1}. \end{aligned}$$

If we now multiply the first term of the second step by one in the form $\frac{1-\lambda_{m+1}}{1-\lambda_{m+1}} = 1$, we may rewrite \mathbf{z} as,

$$\mathbf{z} = (1 - \lambda_{m+1}) \sum_{i=1}^m \frac{\lambda_i}{1 - \lambda_{m+1}} \mathbf{x}_i + \lambda_{m+1} \mathbf{x}_{m+1}.$$

Let,

$$\mathbf{v} := \sum_{i=1}^m \frac{\lambda_i}{1 - \lambda_{m+1}} \mathbf{x}_i.$$

Then, \mathbf{v} is a convex combination of the first m terms and by the induction hypothesis $\mathbf{v} \in C$. Furthermore, by Definition ??,

$$\mathbf{z} = (1 - \lambda_{m+1}) \mathbf{v} + \lambda_{m+1} \mathbf{x}_{m+1} \in C.$$

The induction is complete and the claim is established. \square

Definition 5.17. Convex Hulls Let $S \subseteq \mathbb{R}^n$, then the convex hull of S , denoted $\text{conv}(S)$, is the set comprising all convex combinations of S ,

$$\text{conv}(S) \equiv \left\{ \sum_{i=1}^k \lambda_i \mathbf{x}_i : \mathbf{x}_1, \dots, \mathbf{x}_k \in S, \lambda \in \Delta_k, k \in \mathbb{N} \right\}.$$

We can say that $\text{conv}(S)$ is the smallest convex set containing S . This fact requires a proof.

Lemma 5.18. Let $S \subseteq \mathbb{R}^n$, if $S \subseteq T$ for some convex set T , then $\text{conv}(S) \subseteq T$.

Proof. Suppose that $S \subseteq T$ and take $\mathbf{z} \in \text{conv}(S)$. By the definition of the convex hull, \mathbf{z} can be represented as the convex combination,

$$\mathbf{z} = \sum_{i=1}^k \lambda_i \mathbf{x}_i,$$

of vectors $\mathbf{x}_i \in \text{conv}(S)$. But then, by the convexity of T and the fact that $\mathbf{x}_i \in T$, it must be the case that $\mathbf{z} \in T$. The claim follows. \square

We now turn to presenting a central idea from Caratheodory.

Theorem 5.19. Caratheodory's Theorem Let $S \subset \mathbb{R}^n$ and let $\mathbf{x} \in \text{conv}(S)$. Then there exists $\mathbf{x}_1, \dots, \mathbf{x}_n$ such that $\mathbf{x} \in \text{conv}(\{\mathbf{x}_1, \dots, \mathbf{x}_{n+1}\})$. That is there exists $\lambda \in \Delta_{n+1}$ such that,

$$\mathbf{x} = \sum_{i=1}^{n+1} \lambda_i \mathbf{x}_i.$$

In words, Caratheodory's Theorem says that any element of the convex hull of a subset $S \subseteq \mathbb{R}^n$ can be represented as a convex combination of no more than $n + 1$ vectors of S .

Proof. $\mathbf{x} \in \text{conv}(S)$. By Definition of the convex hull, we know that there exists $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k \in S$ and $\lambda \in \Delta_k$ such that,

$$\mathbf{x} = \sum_{i=1}^k \lambda_i \mathbf{x}_i.$$

We may assume without loss of generality that $\lambda_i > 0$ otherwise the term would simply be omitted from the convex combination. Also, keep in mind $\sum_{i=1}^k \lambda_i = 1$. If $k \leq n + 1$ the claim is proven.

Assume $k > n + 1$, then the vectors, $\mathbf{x}_2 - \mathbf{x}_1, \mathbf{x}_3 - \mathbf{x}_1, \dots, \mathbf{x}_k - \mathbf{x}_1$, which are more than n vectors in \mathbb{R}^n are necessarily linearly dependent. So, there exists $\mu_2, \mu_3, \dots, \mu_k$ which are not all zeros such that,

$$\sum_{i=2}^k \mu_i (\mathbf{x}_i - \mathbf{x}_1) = 0.$$

If we define,

$$\mu_1 := - \sum_{i=2}^k \mu_i,$$

We obtain that,

$$\sum_{i=1}^k \mu_i \mathbf{x}_i = 0.$$

Where not all coefficients μ_i are zero and they satisfy $\sum_{i=1}^k \mu_i = 0$. □

5.4 Convex Cones

A set S is called a cone if it satisfies the following property, for any $\mathbf{x} \in S$ and $\lambda \geq 0$ the inclusion $\lambda \mathbf{x} \in S$ is satisfied.

We now propose a simple and elegant characterization of convex cones.

Lemma 5.20. *A set is a convex cone if and only if the following properties hold,*

- (a) $\mathbf{x}, \mathbf{y} \in S \implies \mathbf{x} + \mathbf{y} \in S$.
- (b) $\mathbf{x} \in S \lambda \geq 0 \implies \lambda \mathbf{x} \in S$.

Proof. (convex cone \implies (a),(b)). Assume S is a convex cone, property (b) follows from the definition. To prove (a), assumed $\mathbf{x}, \mathbf{y} \in S$. By the convexity of S we have that $\frac{1}{2}(\mathbf{x} + \mathbf{y}) \in S$. Then, since S is a cone, we have that $\mathbf{x} + \mathbf{y} = 2 \cdot \frac{1}{2}(\mathbf{x} + \mathbf{y}) \in S$.

((a),(b) \implies convex cone). A set S satisfies properties (a) and (b). By property (b), S is a cone. To show convexity, let $\mathbf{x}, \mathbf{y} \in S$ and let $\lambda \in [0, 1]$. Then, $\lambda \mathbf{x}, (1 - \lambda)\mathbf{y} \in S$ by property (b). Thus, by property (a), $\lambda \mathbf{x} + (1 - \lambda)\mathbf{y} \in S$ thereby establishing convexity of S . □

Definition 5.21. (Conic Combinations) *Given k points such that $\mathbf{x}_1, \dots, \mathbf{x}_k \in \mathbb{R}^n$, a conic combination of these k vectors, is another vector of the form,*

$$\sum_{i=1}^k \lambda_i \mathbf{x}_i,$$

such that each $\lambda \in \mathbb{R}_+^k$.

Now we have the elementary albeit important result.

Lemma 5.22. *Let C be a convex cone, and let $\mathbf{x}_1, \dots, \mathbf{x}_k \in C$ and $\lambda_1, \dots, \lambda_k \geq 0$. Then the convex combination,*

$$\sum_{i=1}^k \lambda_i \mathbf{x}_i \in C.$$

In words, the conic combination remains within the convex cone.

Proof. By the lemma on convex cones, it must be the case that,

$$\lambda_1 \mathbf{x}_1 + \lambda_2 \mathbf{x}_2 + \dots + \lambda_k \mathbf{x}_k = \sum_{i=1}^k \lambda_i \mathbf{x}_i \in C.$$

Where, $\lambda_k \in \mathbb{R}_+^k$. The claim follows as this is the conic combination of the k vectors in C . □

Just as for convex sets and convex hulls, we may naturally propose a definition of the conic hull.

Definition 5.23. (Conic Hull) *Let $S \subset \mathbb{R}^n$. The conic hull of S , denoted $\text{cone}(S)$, is the set comprising all conic combinations of vectors within S ,*

$$\text{cone}(S) := \left\{ \sum_{i=1}^k \lambda_i \mathbf{x}_i : \mathbf{x}_i \in S, \lambda \in \mathbb{R}_+^k, k \in \mathbb{N} \right\}.$$

As we stated for convex hulls, the conic hull of S is the smallest convex cone containing S .

Now the natural question is, whether we can establish an analog of Caratheodory's theorem for convex cones. Interestingly enough, we may establish an even stronger result.

Theorem 5.24. (Conic Representation Theorem) *Let $S \subset \mathbb{R}^n$ and let $\mathbf{x} \in \text{cone}(S)$. Then there exists k linearly independent vectors, $\mathbf{x}_1, \dots, \mathbf{x}_k \in S$ such that $\mathbf{x} \in \text{cone}(\{\mathbf{x}_1, \dots, \mathbf{x}_k\})$. That is, there exists $\lambda \in \mathbb{R}_+^k$ such that,*

$$\mathbf{x} = \sum_{i=1}^k \lambda_i \mathbf{x}_i.$$

In addition, $k \leq n$.

Definition 5.25. (Basic Feasible Solutions) *Let $P = \{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} = \mathbf{b}, \mathbf{x} \geq 0\}$ where $A \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. Suppose that the rows of A are linearly independent. Then $\tilde{\mathbf{x}}$ is a basic feasible solution (bfs) of P if the columns of A corresponding to the indices of the positive values of $\tilde{\mathbf{x}}$ are linearly independent.*

Definition 5.26. (Extreme Point) *Let $S \subseteq \mathbb{R}^n$. A point $\mathbf{x} \in S$ is called an extreme point of S if there do not exist $\mathbf{x}_1, \mathbf{x}_2 \in S$ where $\mathbf{x}_1 \neq \mathbf{x}_2$ and $\lambda \in (0, 1)$ such that $\mathbf{x} = \lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2$.*

That is, an extreme point of a set is a point that cannot be expressed as a trivial convex combination of any other distinct points in the set. Furthermore, the set of all extreme points of S is denoted $\text{ext}(S)$. Geometrically, the set of extreme points of a convex polytope consists of all its vertices.

Convex Functions

We now turn towards the ideas and concepts related to convex functions. These play an essential part in the theory of convex optimization and more generally in non-linear optimization and obviously convex analysis.

Definition 5.27. (Convex Function) A function $f : C \subseteq \mathbb{R}^n \rightarrow R$ where C is convex, is called convex if,

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

$\forall \mathbf{x}, \mathbf{y} \in C$ and $\lambda \in [0, 1]$.

The definition is also known as the *fundamental inequality*. Furthermore, if we do not allow equality in the fundamental inequality, then f is known as a *strictly convex* function.

Another important and related concept is that of *concavity*. We say that a function is concave if $-f$ is convex. More formally,

Definition 5.28. (Concave Function) A function $f : C \subseteq \mathbb{R}^n \rightarrow R$ where C is convex, is called concave if,

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

$\forall \mathbf{x}, \mathbf{y} \in C$ and $\lambda \in [0, 1]$.

Theorem 5.29. (Jensen's Inequality) Let $f : C \rightarrow R$ be a convex function, where $C \subseteq \mathbb{R}^n$ is a convex set. Then for any $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k \in C$ and $\lambda \in \Delta_k$ the following inequality holds,

$$f\left(\sum_{i=1}^k \lambda_i \mathbf{x}_i\right) \leq \sum_{i=1}^k \lambda_i f(\mathbf{x}_i).$$

Proof. We will prove the claim by induction on k . For the base case, $k = 1$, the statement is trivial and follows from the definition of the convex function.

The induction hypothesis is that,

$$f\left(\sum_{i=1}^k \lambda_i \mathbf{x}_i\right) \leq \sum_{i=1}^k \lambda_i f(\mathbf{x}_i). \quad (2)$$

We want to show that the result will hold for $k + 1$. Suppose that $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{k+1} \in C$ and $\lambda \in \Delta_{k+1}$. Let,

$$\mathbf{z} := \sum_{i=1}^{k+1} \lambda_i \mathbf{x}_i.$$

We want to establish that,

$$f(\mathbf{z}) \leq \sum_{i=1}^k \lambda_i f(\mathbf{x}_i).$$

If $\lambda_{k+1} = 1$, we have that $\mathbf{z} = \mathbf{x}_{k+1}$ and the claim follows straightforwardly. If $\lambda_{k+1} < 1$, then

$$\begin{aligned} f(\mathbf{z}) &= f\left(\sum_{i=1}^{k+1} \lambda_i \mathbf{x}_i\right) \\ &= f\left(\sum_{i=1}^k \lambda_i \mathbf{x}_i + \lambda_{k+1} \mathbf{x}_{k+1}\right) \\ &= f\left((1 - \lambda_{k+1}) \sum_{i=1}^k \frac{\lambda_i}{1 - \lambda_{k+1}} \mathbf{x}_i + \lambda_{k+1} \mathbf{x}_{k+1}\right) \\ &\leq (1 - \lambda_{k+1}) f\left(\sum_{i=1}^k \frac{\lambda_i}{1 - \lambda_{k+1}} \mathbf{x}_i\right) + \lambda_{k+1} f(\mathbf{x}_{k+1}). \end{aligned}$$

Let,

$$\mathbf{v} := \sum_{i=1}^k \frac{\lambda_i}{1 - \lambda_{k+1}} \mathbf{x}_i.$$

It follows that \mathbf{v} is a convex combination of k points in the convex set C . By the induction hypothesis, $f(\mathbf{v}) \leq \sum_{i=1}^k \frac{\lambda_i}{1 - \lambda_{k+1}} f(\mathbf{x}_i)$. The induction is thereby complete and the claim follows. \square

First Order Characterization of Convex Functions

Convex functions are not necessarily differentiable. When they are however, we can replace terms in Jensen's inequality with the gradient of the convex function. We will see through the *Gradient Inequality*, tangent hyperplanes of convex functions i.e. gradients of said functions are underestimates.

As one can see, this will play a vital role in optimization problems and finding minima of functions that are convex.

Theorem 5.30. (The Gradient Inequality) *Let $f : C \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuously differentiable function defined over a convex set C . Then f is convex if and only if,*

$$f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) \leq f(\mathbf{y}),$$

for any $\mathbf{x}, \mathbf{y} \in C$.

Proof. Suppose f is a convex function. If $\mathbf{x} = \mathbf{y}$, the claim follows trivially. Assume $\mathbf{x} \neq \mathbf{y}$. By the definition of f being convex, we have that,

$$\begin{aligned} f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) &\leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) \\ \frac{f(\mathbf{x} + \lambda(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\lambda} &\leq f(\mathbf{y}) - f(\mathbf{x}). \end{aligned}$$

Notice, letting $\lambda \rightarrow 0^+$, the quotient on the left hand side will converge to the directional derivative of f at \mathbf{x} in the $\mathbf{y} - \mathbf{x}$ direction.

$$f'(\mathbf{x}; \mathbf{y} - \mathbf{x}) \leq f(\mathbf{y}) - f(\mathbf{x}).$$

By assumption, f is continuously differentiable, so it follows that,

$$f'(\mathbf{x}; \mathbf{y} - \mathbf{x}) \leq \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}).$$

To prove the reverse, assume the gradient inequality holds. Let $\mathbf{z}, \mathbf{v} \in C$ and let $\lambda \in (0, 1)$. We want to show that,

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{v}) \leq \lambda f(\mathbf{z}) + (1 - \lambda) f(\mathbf{v}).$$

Let $\mathbf{u} = \lambda \mathbf{z} + (1 - \lambda) \mathbf{v} \in C$. Then,

$$\mathbf{z} - \mathbf{u} = \frac{\mathbf{u} - (1 - \lambda) \mathbf{v}}{\lambda} - \mathbf{u} = -\frac{1 - \lambda}{\lambda} (\mathbf{v} - \mathbf{u}).$$

We now employ the Gradient Inequality on the pairs $\{\mathbf{z}, \mathbf{u}\}$ and $\{\mathbf{v}, \mathbf{u}\}$.

$$\begin{aligned} f(\mathbf{u}) + \nabla f(\mathbf{u})^T (\mathbf{z} - \mathbf{u}) &\leq f(\mathbf{z}) \\ f(\mathbf{u}) - \frac{\lambda}{1 - \lambda} \nabla f(\mathbf{u})^T (\mathbf{z} - \mathbf{u}) &\leq f(\mathbf{v}), \end{aligned}$$

Multiplying the first inequality by $\frac{\lambda}{1 - \lambda}$ and adding to the second, we get,

$$\frac{1}{1 - \lambda} f(\mathbf{u}) \leq \frac{\lambda}{1 - \lambda} f(\mathbf{z}) + f(\mathbf{v}).$$

Which can be re-arranged,

$$f(\mathbf{u}) \leq \lambda f(\mathbf{z}) + (1 - \lambda) f(\mathbf{v}).$$

Hence, the claim follows, f is convex function. \square

Proposition 5.31. (Sufficiency of Stationarity under Convexity) Let f be a continuously differentiable convex function over the convex set $C \subseteq \mathbb{R}^n$. Suppose the $\nabla f(\mathbf{x}^*) = \mathbf{0}$ then \mathbf{x}^* is a global minimizer of f over C .

Proof. Let $\mathbf{z} \in C$, applying the gradient inequality and using the fact that $\nabla f(\mathbf{x}^*) = \mathbf{0}$, we get that,

$$f(\mathbf{x}^*) \leq f(\mathbf{z}), \forall \mathbf{z} \in C.$$

So, \mathbf{x}^* is a global minimizer of f over C . □

Note that proposition becomes a necessary condition when $C = \mathbb{R}^n$.

Theorem 5.32. (Necessity and Sufficiency of Stationarity under Convexity) Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuously differentiable convex function. Then $\nabla f(\mathbf{x}^*) = \mathbf{0}$ if and only if \mathbf{x}^* is a global minimum point of f over \mathbb{R}^n .

Now we turn to establishing strict convexity of quadratic convex functions.

Theorem 5.33. (Convexity and strict convexity of quadratic functions with positive semidefinite matrices) Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be the quadratic function of the form $\mathbf{x}^T A \mathbf{x} + 2\mathbf{b}^T \mathbf{x} + c$ where $A \in \mathbb{R}^{n \times n}$ is a symmetric matrix, $\mathbf{b} \in \mathbb{R}^n$ and $c \in \mathbb{R}$. Then, f is convex if $A \succcurlyeq 0$ and strictly convex if $A \succ 0$.

Proof. We start by proving convexity. We know from the theorem on the gradient inequality, the convexity of f is equivalent to the validity of the gradient inequality. So, we want to show that,

$$\mathbf{x}^T A \mathbf{x} + 2\mathbf{b}^T \mathbf{x} + c + 2(\mathbf{A}\mathbf{x} + \mathbf{b})^T (\mathbf{y} - \mathbf{x}) \leq \mathbf{y}^T A \mathbf{y} + 2\mathbf{b}^T \mathbf{y} + c.$$

Where the inequality holds for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. With some rearrangement we may write the inequality as,

$$(\mathbf{y} - \mathbf{x})^T A (\mathbf{y} - \mathbf{x}) \geq 0.$$

Let $\mathbf{d} = \mathbf{y} - \mathbf{x}$, then we have that $\mathbf{d}^T A \mathbf{d} \geq 0$ for any $\mathbf{d} \in \mathbb{R}^n$, therefore we have that $A \succcurlyeq 0$.

We now want to show strict convexity. This will amount to using the strict gradient inequality,

$$\mathbf{x}^T A \mathbf{x} + 2\mathbf{b}^T \mathbf{x} + c + 2(\mathbf{A}\mathbf{x} + \mathbf{b})^T (\mathbf{y} - \mathbf{x}) < \mathbf{y}^T A \mathbf{y} + 2\mathbf{b}^T \mathbf{y} + c.$$

By the analogous argument we get that $\mathbf{d}^T A \mathbf{d} > 0$, which is equivalent to saying that $A \succ 0$ and the claim follows. □

Another type of first order characterization is the monotonicity of the gradient. In the one dimensional case, this characterization states that the derivative is non-decreasing. However, we will need to generalize for n -dimensions.

Theorem 5.34. (Monotonicity of the Gradient) Suppose that $f : C \subset \mathbb{R}^n \rightarrow \mathbb{R}$ is a continuously differentiable function and C is a convex set. Then, f is convex over C if and only if,

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

for any $\mathbf{x}, \mathbf{y} \in C$.

Second Order Characterizations of Convex Functions

When the function is twice continuously differentiable, the convexity of the function is characterized by the positive semidefiniteness of the Hessian matrix.

Theorem 5.35. (Second order characterization of convexity) Let $f : C \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$ be a twice continuously differentiable function defined over the convex set C . Then f is convex if and only if the Hessian is positive semidefinite i.e. $\nabla^2 f(x) \succeq 0$ for any $x \in C$.

To prove the theorem it suffices to show that the gradient inequality holds true. For the opposite direction, we take f to be convex and then show that the Hessian is positive semi-definite.

We can also state conditions for strict convexity.

Theorem 5.36. (Sufficient second order condition for strict convexity) Let f be a twice continuously differentiable function over a convex set $C \subset \mathbb{R}^n$. Suppose $\nabla^2 f \succ 0$, then f is strictly convex over C .

The proof of the theorem follows the same form in one direction for the general second order characterization for convexity.

6 Convex Optimization

We now begin a central concept in non-linear optimization, namely convex optimization. Our efforts here will be to minimize a convex function $f : C \subset \mathbb{R}^n \rightarrow \mathbb{R}$, over the closed, convex set C .

We will consider explicit convex optimization problems of the form,

$$\min f(\mathbf{x}), \text{ s.t. } g_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) \leq 0.$$

Where i, j belong to some finite index set. For the problem to be convex, inequality constraints must be themselves convex functions and equality constraints must be affine. We now note a hall mark of convex problems. The implications of the following theorem are obvious.

Theorem 6.1. (Local min=Global min) Let $f : C \rightarrow \mathbb{R}$ be a convex function over the convex set C . Let $\mathbf{x}^* \in C$ be a local minimum of f over C . Then \mathbf{x}^* is a global minimum of f over C .

Proof.

□

It follows by only a slight modification, that for a strictly convex function over a convex set, a local minimum is indeed a strict global minimum.

Definition 6.2. (Optimal Set) The optimal set of a convex problem, is the set of all minimizers. Denoted,

$$\operatorname{argmin}\{f(\mathbf{x}) : \mathbf{x} \in C\} \quad (3)$$

We may extend this notion to any general optimization problem. Furthermore, **we may show that the optimal set of a convex problem is itself convex.**

We may also define problems such as minimizing concave functions over concave sets as convex problems by a symmetry argument. Just multiply the objective function by -1 .

All linear programming problems are convex problems.

6.1 Orthogonal Projection Operator

Given a non-empty closed, convex set C . We define the orthogonal projection operator as $P_C : \mathbb{R}^n \rightarrow C$ as,

$$P_C(\mathbf{x}) = \operatorname{argmin}\{\|\mathbf{y} - \mathbf{x}\|^2 : \mathbf{y} \in C\}.$$

Here \mathbf{x} is fixed and not necessarily in C .

Geometrically, the orthogonal projection operator returns a point in C that is closest to the point $\mathbf{x} \in \mathbb{R}^n$. Obviously, if $\mathbf{x} \in C$, then $P_C(\mathbf{x}) = \mathbf{x}$.

We also notice that the OPO is defined via a convex optimization problem. Namely, minimizing the quadratic, square of the norm of the difference over the convex set C .

We now show that P_C is in fact well defined.

Theorem 6.3. (First Projection Theorem) *Let C be a non-empty closed convex set. Then, the problem,*

$$P_C(\mathbf{x}) = \operatorname{argmin}\{\|\mathbf{y} - \mathbf{x}\|^2 : \mathbf{y} \in C\},$$

has a unique optimal solution.

Proof. Since the objective function is quadratic with positive definite matrix we know that the function must be coercive. It follows that it must attain at least one optimal solution. Furthermore, we may deduce that the function is strictly convex, we can conclude there exists only one optimal solution and it is attained by coerciveness. \square

We know that the norm induces a metric,

$$d(\mathbf{x}, C) = \min_{\mathbf{y} \in C} \|\mathbf{x} - \mathbf{y}\|.$$

We may also express this metric in terms of the OPO,

$$d(\mathbf{x}, C) = \|\mathbf{x} - P_C(\mathbf{x})\|.$$

7 Optimization over Convex Sets

We now consider a class of optimization problems, where f is a continuously differentiable function over a closed convex set C .

Instead now of considering stationary points of a function within an unconstrained problem. We must now consider stationary points of a problem.

Definition 7.1. (Stationary points of a constrained problem) *Let f be a continuously differentiable function over a closed convex set C . Then $\mathbf{x}^* \in C$ is called a stationary point of the problem if,*

$$\nabla f(\mathbf{x}^*)^T (\mathbf{x} - \mathbf{x}^*) \geq 0,$$

for any $\mathbf{x} \in C$.

We can see, that geometrically, stationarity here means that there are no feasible descent directions at the optimal \mathbf{x}^* . Furthermore, we can deduce that stationarity is a necessary condition for a local minimum of the problem.

Theorem 7.2. (Stationarity as a necessary optimality condition) *Let f be a continuously differentiable function over a closed convex set C and let \mathbf{x}^* be a local minimum of the problem. Then \mathbf{x}^* is a stationary point of the problem.*

7.1 Stationarity in Convex Problems

Stationarity is a necessary optimality condition for local optimality. However, when the objective function is also convex, stationarity is a necessary and sufficient condition for optimality.

Theorem 7.3. *Let f be a continuously differentiable convex function over the closed and convex set C . Then \mathbf{x}^* is a stationary point of the constrained optimization problem if and only if \mathbf{x}^* is an optimal solution of the problem.*

7.2 Orthogonal Projection Revisited

We now revisit the idea of the orthogonal projection operator to put forward the second projection theorem.

To give a general idea, given a closed convex set C , $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{y} \in C$. The second projection theorem asserts that the angle between $\mathbf{x} - P_C(\mathbf{x})$ and $\mathbf{y} - P_C(\mathbf{x})$ is obtuse i.e. greater than or equal to 90 degrees. More formally we have the following.

Theorem 7.4. (Second projection theorem) *Let C be a closed convex set and let $\mathbf{x} \in \mathbb{R}^n$. Then $\mathbf{z} = P_C(\mathbf{x})$ if and only if,*

$$(\mathbf{x} - P_C(\mathbf{x}))^T(\mathbf{y} - P_C(\mathbf{x})) \leq 0,$$

for any $\mathbf{y} \in C$.

8 Optimality Conditions for Linearly Constrained Problems

In the last section we discussed the concept of stationarity within optimization problems where the objective function is differentiable and the feasible set is closed and convex. We saw that the condition is quite tricky to use apart from some simpler feasible sets.

In this section we will look at constructing a stronger optimality condition (at least from linearly constrained problems) that will be easier to use while performing optimizations. These will be the so called Karush-Kuhn-Tucker conditions for linearly constrained problems. Later on in these notes we will also extend the KKT conditions to other types of optimization problems.

8.1 Separation and Alternative Theorems

We begin by stating a somewhat straight forward fact about convex sets that will provide a strong foundation for what is to follow.

Given a set $S \subset \mathbb{R}^n$, a hyperplane $H = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^T \mathbf{x} = b\}$ is said to strictly separate a point $\mathbf{y} \notin S$ if $\mathbf{a}^T \mathbf{y} > b$ and $\mathbf{a}^T \mathbf{x} \leq b$ for all $\mathbf{x} \in S$ and $\mathbf{y} \notin S$. Given this simple idea, consider the following theorem.

Theorem 8.1. (Strict Separation Theorem) *Let $C \subset \mathbb{R}^n$ be a closed and convex set and let $\mathbf{y} \notin C$. Then there exists non-zero \mathbf{p} and $\alpha \in \mathbb{R}$ such that,*

$$\mathbf{p}^T \mathbf{y} > \alpha$$

and,

$$\mathbf{p}^T \mathbf{x} \leq \alpha,$$

for all $\mathbf{x} \in C$.

We now use the theorem to prove Farka's lemma, which is regarded as an alternative theorem. Reason being, it states that exactly one of the two systems is feasible.

Theorem 8.2. (Farka's Lemma) Let $C \in \mathbb{R}^n$ and $\mathbf{A} \in \mathbb{R}^{m \times n}$. Then, exactly one of the following systems has a solution.

(i) $\mathbf{Ax} \leq \mathbf{0}$, $\mathbf{c}^T \mathbf{x} > 0$.

(ii) $\mathbf{A}^T \mathbf{y} = \mathbf{c}$, $\mathbf{y} \geq \mathbf{0}$.

Before carrying on, consider the following illustration. Let $\mathbf{c}^T = (-1, 9)$ and

$$\mathbf{A} = \begin{pmatrix} 1 & 5 \\ -1 & 2 \end{pmatrix}.$$

Stating the system (i) is infeasible means that the system $\mathbf{Ax} \leq \mathbf{0}$ implies that the inequality $\mathbf{c}^T \mathbf{x} \leq 0$. This is true as adding twice the second inequality to the first establishes the result. It follows that \mathbf{c}^T can be written as a conic combination of rows of \mathbf{A} .

Now, in generality, the question becomes whether a set of linear inequalities in a base system can be expressed as linear inequalities in some new system if and only if the new inequalities can be written as conic combinations of inequalities in the base system. Amazingly, Farka's lemma says yes!

With this as a motivation, we may re-formulate Farka's Lemma to make it easier to prove.

Theorem 8.3. (Farka's Lemma, second formulation) Let $\mathbf{c} \in \mathbb{R}^n$ and $\mathbf{A} \in \mathbb{R}^{m \times n}$. Then the following are equivalent,

(A) The implication $\mathbf{Ax} \leq \mathbf{0} \implies \mathbf{c}^T \mathbf{x} \leq 0$ holds true.

(B) There exists $\mathbf{y} \in \mathbb{R}_+^m$ such that $\mathbf{A}^T \mathbf{y} = \mathbf{c}$.

Another alternative theorem to consider is Gordon's Theorem.

Theorem 8.4. (Gordon's Alternative Theorem) Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ then exactly one of the following systems has a solution.

(A) $\mathbf{Ax} < \mathbf{0}$.

(B) $\mathbf{p} \neq \mathbf{0}$, $\mathbf{A}^T \mathbf{p} = \mathbf{0}$, $\mathbf{p} \geq \mathbf{0}$.

9 KKT Conditions

We now use Gordon's theorem to state a special case of the KKT conditions. These conditions in turn extend the idea of stationarity as discussed in the previous section.

Theorem 9.1. (KKT Conditions for Linearly Constrained Problems) Consider the minimization problem,

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{a}_i^T \mathbf{x} \leq b_i, \quad i = 1, 2, \dots, m \end{aligned}$$

Where f is continuously differentiable over \mathbb{R}^n , $\mathbf{a}_1, \dots, \mathbf{a}_m \in \mathbb{R}^n$ and $b_1, \dots, b_m \in \mathbb{R}$ and let \mathbf{x}^* denote the local minimum point of the problem.

Then there exist Lagrange multipliers $\lambda_1, \dots, \lambda_m \geq 0$ such that,

$$\nabla f(\mathbf{x}^*) + \sum_1^m \lambda_i \mathbf{a}_i = \mathbf{0},$$

and

$$\lambda_i (\mathbf{a}_i^T \mathbf{x}^* - b_i) = 0, \quad i = 1, 2, \dots, m.$$

Proof. Since \mathbf{x}^* is a local minimum point of the problem, it follows from the stationary condition of differentiable functions over convex sets that,

$$\nabla f(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) \geq 0,$$

for all $\mathbf{x} \in \mathbb{R}^n$ such that \mathbf{x} satisfies the constraints i.e. $\mathbf{a}_i^T \mathbf{x} \leq b_i$.

Let the set of active constraints be denoted,

$$I(\mathbf{x}^*) := \{i : \mathbf{a}_i^T \mathbf{x}^* = b_i\}.$$

Furthermore, we introduce the change of variable $\mathbf{y} = (\mathbf{x} - \mathbf{x}^*)$, making the stationarity condition $\nabla f(\mathbf{x}^*)^T \mathbf{y} \geq 0$ for any \mathbf{y} satisfying $\mathbf{a}_i^T (\mathbf{x}^* + \mathbf{y}) \leq b_i$. That is for any $\mathbf{y} \in \mathbb{R}^n$ such that,

$$\begin{aligned} \mathbf{a}_i^T \mathbf{y} &\leq 0, \quad i \in I(\mathbf{x}^*) \\ \mathbf{a}_i^T \mathbf{y} &\leq b_i - \mathbf{a}_i^T \mathbf{x}^*, \quad i \notin I(\mathbf{x}^*). \end{aligned}$$

We now want to show that the latter inequalities can be removed. That is $\mathbf{a}_i^T \mathbf{y} \leq 0 \implies \nabla f(\mathbf{x}^*)^T \mathbf{y} \geq 0$.

Suppose then that \mathbf{y} satisfies, $\mathbf{a}_i^T \mathbf{y} \leq 0$, $i \in I(\mathbf{x}^*)$. Since $b_i - \mathbf{a}_i^T \mathbf{x}^* > 0$ for all $i \notin I(\mathbf{x}^*)$, it follows there exists a small enough $\alpha > 0$ for which $\mathbf{a}_i^T \mathbf{y} \leq b_i - \mathbf{a}_i^T \mathbf{x}^*$. Since in addition $\mathbf{a}_i^T \alpha \mathbf{y} \leq 0$, it follows by stationarity that $\nabla f(\mathbf{x}^*)^T (\alpha \mathbf{y}) \geq 0$ and hence $\nabla f(\mathbf{x}^*)^T \mathbf{y} \geq 0$. We have thus shown that,

$$\mathbf{a}_i^T \mathbf{y} \leq 0 \text{ for all } i \in I(\mathbf{x}^*) \implies \nabla f(\mathbf{x}^*)^T \mathbf{y} \geq 0.$$

Then, by employing Farka's lemma, there exists a $\lambda_i \geq 0$ such that

$$-\nabla f(\mathbf{x}^*) = \sum_{i \in I(\mathbf{x}^*)} \lambda_i \mathbf{a}_i.$$

Lastly, defining $\lambda_i = 0$ for $i \notin I(\mathbf{x}^*)$ we get that $\lambda_i (\mathbf{a}_i^T \mathbf{x}^* - b_i) = 0$. So,

$$\nabla f(\mathbf{x}^*) + \sum_i^m \lambda_i \mathbf{a}_i = 0.$$

□

These KKT conditions are a necessary optimality condition, but when the objective function is convex, then KKT are both necessary and sufficient global optimality conditions.

We refer to the condition $\lambda_i (\mathbf{a}_i^T \mathbf{x} - b_i) = 0$ as *complementary slackness*.

Furthermore, we may generalize the theorem to encapsulate linear equality constraints. The proof of this variant relies on the fact that a linear equality can be decomposed into two linear inequalities. The proof then follows the same structure as before. We present the theorem below.

Theorem 9.2. (KKT conditions for linearly constrained problems) Consider the minimization problem,

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{a}_i^T \mathbf{x} \leq b_i, \quad i = 1, 2, \dots, m . \\ & \mathbf{c}_j^T \mathbf{x} = d_j, \quad j = 1, 2, \dots, p \end{aligned}$$

f is continuously differentiable over \mathbb{R}^n , $\mathbf{a}_1, \dots, \mathbf{a}_m, \mathbf{c}_1, \dots, \mathbf{c}_p \in \mathbb{R}^n$ and $b_1, \dots, b_m, d_1, \dots, d_p \in \mathbb{R}$. Then, we have the following,

(a) If \mathbf{x}^* is a local minimum point of the problem, then there exist $\lambda_1, \dots, \lambda_m \geq 0$ and $\mu_1, \dots, \mu_p \in \mathbb{R}$ such that

$$\nabla f(\mathbf{x}^*) + \sum_1^m \lambda_i \mathbf{a}_i + \sum_1^p \mu_j \mathbf{c}_j = \mathbf{0},$$

and

$$\lambda_i (\mathbf{a}_i^T \mathbf{x}^* - b_i) = 0.$$

(b) If in addition, f is convex over \mathbb{R}^n and \mathbf{x}^* is a feasible solution of the problem and the existence of the Lagrange multipliers is as above, then \mathbf{x}^* is an optimal solution of the problem.

Observe that the Lagrange multipliers for linear equality constraints need not be non-negative.

Furthermore, \mathbf{x}^* is called a KKT point of the problem if (a) is satisfied.

9.0.1 The Lagrangian

A popular and compact representation of the KKT conditions is via the Lagrangian function. We present this idea here within the setting of general non-linear programming problems. Consider the problem,

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & g_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, m . \\ & h_j(\mathbf{x}) = 0, \quad j = 1, 2, \dots, p \end{aligned}$$

Here $f, g_1, \dots, g_m, h_1, \dots, h_p$ are continuously differentiable over \mathbb{R}^n . The associated Lagrangian takes the form,

$$L(\mathbf{x}, \lambda, \mu) = f(\mathbf{x}) + \sum_1^m \lambda_i g_i(\mathbf{x}) + \sum_1^p \mu_j h_j(\mathbf{x}).$$

Notice then, we gain back the stationarity condition by taking the gradient of the Lagrangian with respect to \mathbf{x} ,

$$\nabla_{\mathbf{x}} L(\mathbf{x}, \lambda, \mu) = \nabla f(\mathbf{x}) + \sum_1^m \lambda_i \nabla g_i(\mathbf{x}) + \sum_1^p \mu_j \nabla h_j(\mathbf{x}) = \mathbf{0}.$$

Now considering the problem as in the KKT Theorem, if we define matrices \mathbf{A}, \mathbf{C} and vectors $\mathbf{b} \in \mathbb{R}^m$ $\mathbf{d} \in \mathbb{R}^p$ such that,

$$\mathbf{A} = \begin{pmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_m^T \end{pmatrix}, \quad \mathbf{C} = \begin{pmatrix} \mathbf{c}_1^T \\ \vdots \\ \mathbf{c}_p^T \end{pmatrix}.$$

Considering this, the constraints can be written as $\mathbf{A}\mathbf{x} = \mathbf{b}$ and $\mathbf{C}\mathbf{x} = \mathbf{d}$. Then we may re-write the Lagrangian as,

$$L(\mathbf{x}, \lambda, \mu) = f(\mathbf{x}) + \lambda^T (\mathbf{A}\mathbf{x} - \mathbf{b}) + \mu^T (\mathbf{C}\mathbf{x} - \mathbf{d}).$$

In turn making the stationary condition, expressed in terms of the Lagrangian as,

$$\nabla_{\mathbf{x}} L(\mathbf{x}, \lambda, \mu) = \nabla f(\mathbf{x}) + \mathbf{A}^T \lambda + \mathbf{C}^T \mu = \mathbf{0}.$$

Duality

In general, we know convex problems are easier to solve than non-convex problems. Furthermore, outside convexity, we know that unconstrained problems are easier to solve than constrained ones.

We want to put forward a theory that somehow can use an unconstrained version of a constrained optimization problem to provide an estimate of the original constrained problem.

This notion sparks the idea of duality within mathematical optimization. Which we now construct.

Consider a non-linear minimization problem P . We find that in general, the unconstrained version of the problem P_0 , has optimal solution such that,

$$\text{val}(P_0) \leq \text{val}(P).$$

In words, the optimal value of the unconstrained problem acts as a lower bound on the optimal value of the constrained problem.

If we now optimize over $\text{val}(P_0)$ i.e. find the best lower bound, we may formulate another optimization problem D such that,

$$\max\{\text{val}(P_0)\}.$$

We call this problem the **dual problem** and the original problem P , the primal problem.

We find that,

$$\text{val}(D) \leq \text{val}(P).$$

Deriving the dual problem from the primal problem is not difficult but is non-trivial. Given the primal problem P . We denote the Lagrangian as $L(\mathbf{x}, \lambda, \mu)$, assuming P is in general non-linear and consists of both equality and inequality constraints.

Then, the **dual objective function** is defined as,

$$q(\lambda, \mu) = \inf_{\mathbf{x} \in X} L(\mathbf{x}, \lambda, \mu).$$

When the minimum is attained, the infimum can be written explicitly as a minimum.

One may notice that $q(\lambda, \mu) = -\infty$ for some values λ and μ . We thereby define the domain of q to be,

$$\text{dom}(q) = \{(\lambda, \mu) \in \mathbb{R}_+^m \times \mathbb{R}^p : q(\lambda, \mu) > -\infty\}.$$

We may then define the dual problem explicitly as,

$$q^* = \max q(\lambda, \mu), \text{ s.t. } (\lambda, \mu) \in \text{dom}(q).$$

We then gain the following powerful theorem.

Theorem 9.3. (Convexity of the Dual Problem) Consider the general constrained problem P . Let the constraint functions be finite valued defined on the set $X \subset \mathbb{R}^n$. Let q be the function as defined before, then,

- (a) $\text{dom}(q)$ is a convex set.
- (b) q is concave over $\text{dom}(q)$, hence $-q$ is convex over $\text{dom}(q)$.

Proof.

□

Exploiting this fact we can formulate the so called duality theorems.

Theorem 9.4. (Weak Duality) Consider the primal and dual problem. Then,

$$\text{val}(D) \leq \text{val}(P).$$

Straightforwardly, the weak duality theorem states that the optimal value of the dual problem is a lower bound on the optimal value of the primal problem.

Strong Duality in the Convex Case

We now turn to establishing in the convex case i.e. when the primal problem is convex, **strong duality** holds. As the name suggests, the strong duality is precisely when $\text{val}(D) = \text{val}(P)$.

To do this rigorously, we will need an idea of separating a point in a convex set from another point not necessarily within the convex set via hyperplanes.

Theorem 9.5. (Supporting Hyperplane Theorem) Let $C \subset \mathbb{R}^n$ be a convex set and let $\mathbf{y} \notin C$. Then there exists $\mathbf{0} \neq \mathbf{p} \in \mathbb{R}^n$ such that,

$$\mathbf{p}^T \mathbf{x} \leq \mathbf{p}^T \mathbf{y},$$

for any $\mathbf{x} \in C$.

From here we may also deduce a separation theorem between disjoint convex sets.

Theorem 9.6. (Separation of two convex sets) Let $C_1, C_2 \subset \mathbb{R}^n$ be two non-empty, disjoint convex sets. Then there exists $\mathbf{0} \neq \mathbf{p} \in \mathbb{R}^n$ such that,

$$\mathbf{p}^T \mathbf{x} \leq \mathbf{p}^T \mathbf{y},$$

for any $\mathbf{x} \in C_1$ and $\mathbf{y} \in C_2$.

This in turn will help in developing a non-linear version of Farkas Lemma, which will then be the key in proving the strong duality result.

Theorem 9.7. (Non-Linear Farkas Lemma) Let $X \subset \mathbb{R}^n$ be a convex set and let f, g_1, g_2, \dots, g_m be convex functions over X . Assume that there exists a $\hat{\mathbf{x}} \in X$ such that, $g_1(\hat{\mathbf{x}}) < 0, \dots, g_m(\hat{\mathbf{x}}) < 0$. Let $c \in \mathbb{R}$. Then the following two claims are equivalent.

(a) The following implication holds:

$$\mathbf{x} \in X, g_i(\mathbf{x}) \leq 0, \implies f(\mathbf{x}) \geq c.$$

(b) There exists Lagrange multipliers $\lambda_1, \dots, \lambda_m$ such that,

$$\min_{\mathbf{x} \in X} \left\{ f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) \right\} \geq c.$$

We see that the non-linear Farkas lemma shows that the convex function f is bounded below by c and this is equivalent to the minimum over $\mathbf{x} \in X$ of the Lagrangian being bounded below by c .

We can then turn to establishing a strong duality result with a Slater-type condition.

Theorem 9.8. (Strong duality for convex problems with inequality constraints) Consider the optimization problem,

$$f^* = \min f(\mathbf{x}), \text{ s.t. } g_i(\mathbf{x}) \leq 0, \mathbf{x} \in X.$$

Here X is a convex set and f, g_1, \dots, g_m are convex functions over X . Suppose that there exists $\hat{\mathbf{x}} \in X$ such that $g_i(\hat{\mathbf{x}}) < 0$. Suppose also that the problem has a finite optimal value. Then the optimal value of the dual problem,

$$q^* = \max\{q(\lambda) = \min_{\mathbf{x} \in X} L(\mathbf{x}, \lambda) : \lambda \in \text{dom}(q)\}.$$

Then the optimal value of the primal and dual problem are,

$$q^* = f^*.$$