

# Auxiliary notes on Optimisation

Daniel Lin; Professor: Dante Kalise

April 18, 2023

The official notes on optimisation are almost self-contained, but some details are omitted. This is an additional note for extensional study only, mainly collected from *Introduction to Non-Linear Optimisation* and lectures.

The following table of contents is clickable.

## Contents

<b>1</b>	<b>Basics</b>	<b>2</b>
1.1	<a href="#">Spectral decomposition theorem</a>	2
1.2	<a href="#">Definiteness</a>	3
1.2.1	<a href="#">Summary of criterion for definiteness and semi-definiteness</a>	4
1.3	<a href="#">Transformation</a>	5
1.4	<a href="#">Derivatives</a>	5
<b>2</b>	<b>Unconstrained optimisation</b>	<b>6</b>
2.1	<a href="#">Global Optimality</a>	6
2.2	<a href="#">Quadratic functions</a>	6
<b>3</b>	<b>Gradient Descent</b>	<b>7</b>
3.1	<a href="#">Kantorovich inequality</a>	8
3.2	<a href="#">Fermat-Weber Problem</a>	10
3.3	<a href="#">Convergence of gradient descent</a>	11
3.4	<a href="#">Newton's method</a>	13
3.5	<a href="#">Kaczmarz method</a>	13
<b>4</b>	<b>Convexity</b>	<b>15</b>
4.1	<a href="#">Visual Explanations of some Geometric Shapes</a>	15
4.2	<a href="#">Properties of Convex set</a>	15
4.3	<a href="#">Convex Functions</a>	16
4.4	<a href="#">First, second order characterisations of convex functions</a>	17
4.5	<a href="#">Properties of Convex Functions</a>	19
4.6	<a href="#">Level sets</a>	20
4.7	<a href="#">Maxima of convex functions</a>	20
<b>5</b>	<b>Convex Optimisation</b>	<b>21</b>
5.1	<a href="#">Classification using linear separator</a>	22
5.2	<a href="#">Stationary points</a>	23
5.3	<a href="#">Orthogonal Projection</a>	24
5.4	<a href="#">Gradient Projection</a>	26
5.4.1	<a href="#">Backtracking</a>	27
5.5	<a href="#">Sparsity constraint</a>	28
<b>6</b>	<b>KKT conditions</b>	<b>29</b>
6.1	<a href="#">General KKT Condition</a>	32
6.2	<a href="#">KKT for Convex Problems</a>	34
<b>7</b>	<b>Duality</b>	<b>35</b>
7.1	<a href="#">Strong Duality</a>	37

# 1 Basics

## 1.1 Spectral decomposition theorem

Recall diagonalisation of an  $n \times n$  real symmetric matrix has satisfying properties

**Theorem 1.1** (Spectral decomposition). *Given symmetric  $A \in \mathbb{R}^{n \times n}$ , there is orthogonal  $U$  and diagonal matrix  $D$  s.t.*

$$U^T A U = D$$

*elements of  $D$  are Eigenvalues of  $A$  and usually we assume WOLG they are laid in descending order (i.e.  $d_1 \geq d_2 \geq \dots \geq d_n$  where  $d_i$  are the diagonal elements of  $D$ )*

**Corollary 1.** *For symmetric matrix  $A$ ,*

$$\text{Tr}(A) = \sum_{i=1}^n \lambda_i(A), \quad \det A = \prod_{i=1}^n \lambda_i(A)$$

*where  $\lambda_i(A)$  is the  $i$ 'th Eigenvalue of  $A$ .*

*Proof.* Trace is invariant under cyclic permutation, so

$$\text{Tr}(A) = \text{Tr}(U D U^T) = \text{Tr}(U^T U D) = \text{Tr}(D) = \sum_{i=1}^n \lambda_i(A)$$

as for determinants, this function is multiplicative.

$$\det A = \det(U D U^T) = \det U \det D \det U^T = \det U \det U^{-1} \det D = \prod_{i=1}^n \lambda_i(A)$$

□

Corollary 1 gives that criteria of definiteness for  $2 \times 2$  matrices:  $\text{Tr}(A), \det A > 0 \Leftrightarrow A$  is positive definite.

Spectral decomposition also enables an upper bound for the smallest Eigenvalue  $\lambda_{\min}(A)$  and a lower bound for the largest Eigenvalue  $\lambda_{\max}(A)$

**Corollary 2.** *If  $A$  is symmetric,*

$$\lambda_{\min}(A) \leq \frac{\mathbf{x}^T A \mathbf{x}}{\|\mathbf{x}\|^2} \leq \lambda_{\max}(A)$$

*The middle term is also called Rayleigh quotient, denoted by  $R_A(\mathbf{x})$ . This quotient is closely related to the definiteness of  $A$ .*

*Proof.* Perform spectral decomposition on  $A$  and define change of variable  $\mathbf{x} = U \mathbf{y}$  so that  $\mathbf{x}^T A \mathbf{x} = \mathbf{y}^T D \mathbf{y}$ . Then

$$\max_{\mathbf{x} \neq 0} \frac{\mathbf{x}^T A \mathbf{x}}{\|\mathbf{x}\|^2} = \max_{\mathbf{y} \neq 0} \frac{\mathbf{y}^T D \mathbf{y}}{\|\mathbf{y}\|^2} = \max_{\mathbf{y} \neq 0} \frac{\sum d_i y_i^2}{\sum y_i^2}$$

the denominator  $\sum d_i y_i^2 \leq \lambda_{\max}(A) \sum (y_i^2)$  as  $\lambda_{\max}(A)$  is the maximal element of  $D$ . So

$$R_A(\mathbf{x}) \leq \max_{\mathbf{x} \neq 0} \frac{\mathbf{x}^T A \mathbf{x}}{\|\mathbf{x}\|^2} \leq \lambda_{\max}(A)$$

$R_A(\mathbf{x}) \geq \lambda_{\min}(A)$  can be deduced in a similar way.

□

## 1.2 Definiteness

We only talk about the definiteness of symmetric matrices, but you can always get a symmetric matrix by  $(A + A^T)/2$ .

**Proposition 1.2** (Properties for positive definite matrices). *If  $A$  is positive (semi-)definite,*

- *then all symmetric submatrices are positive (semi-)definite.*
- *$A^{-1}$  exists and it is also positive (semi-)definite*
- *diagonal entries of  $A$  are (non-negative) positive*
- *element of  $A$  with the largest absolute value must be on the main diagonal*
- $\det(A) > 0$
- *If  $A, B$  are positive (semi-)definite,  $A + B$  is positive (semi-)definite*

*If  $A$  is negative (semi-)definite, then  $-A$  is positive (semi-)definite so the above properties apply to  $-A$ , reversing the sign gives properties of negative (semi-)definite matrices.*

Among the three criterion in lecture notes for definiteness (principal minor, strict diagonal dominance, Eigenvalue), the strict diagonal dominance criteria is NOT necessary condition. So, for example,

$$\begin{pmatrix} 3/2 & 1 & 1 \\ 1 & 3/2 & 1 \\ 1 & 1 & 3/2 \end{pmatrix}$$

has Eigenvalues 0.5, 0.5, 3.5, it is positive definite, but diagonal dominance is not satisfied. However, this is the most convenient criterion, so you may consider trying this criteria first.

Further, without strict diagonal dominance (or even absence of diagonal dominance),  $A$  may still be positive definite. It is always worth checking other criteria or simply find  $\mathbf{x} \neq 0$  s.t.  $\mathbf{x}^T A \mathbf{x} = 0$  to prove the matrix is not positive definite.

Actually, as long as  $a \geq 1$ , the matrix  $A$  defined by

$$A_{ij} = \begin{cases} a & i = j \\ 1 & i \neq j \end{cases}$$

is positive semi-definite.

Principal minors  $D_i(A)$  used in lectures are called leading principal minors. Following is a proper definition for principal minors.

**Definition 1** (Principal minor). Principal minors of  $A \in \mathbb{R}^{n \times n}$  are determinants of principal submatrices  $A$ , which are obtained by deleting any rows and the corresponding columns. If only the last  $n - k$  rows and columns are deleted, the submatrix is called the leading principal submatrix, and its determinant is called the leading principal minor, denoted as  $D_k(A)$ .

Leading principal minor criterion CANNOT be used to detect semi-definiteness, i.e. even if  $D_i(A) \geq 0 \Leftrightarrow A$  may be indefinite:

$$\begin{pmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & -1 \end{pmatrix}$$

three principal minors are 1, 0, 0, but the matrix is indefinite. (check the Eigenvalues) This counter-example is provided by Jorge Catarcha Otero Saavedra. Even if all the diagonal entries are positive, consider the following matrix

$$A := \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & a \end{pmatrix}$$

where  $0 < a < 1$ . The three leading principal minors are 1, 0, 0 (independent of  $a$ !), all non-negative. But consider  $\mathbf{x}$  s.t.  $x_1 + x_2 + x_3 = 0$

$$\mathbf{x}^T A \mathbf{x} = (x_1 + x_2 + x_3)^2 + (a - 1)z^2 = (a - 1)z^2 < 0$$

Though leading principal minors do not work, there is a criterion for semi-definiteness: If all principal minors (i.e. determinants of all possible combinations of rows and corresponding columns) are non-negative, then  $A$  is semi-definite.

There is another iff condition for positive definiteness, though not useful in practice:  $A \succ 0 \Leftrightarrow$  exists invertible  $M$  s.t.  $A = MM^T$ . You can assemble  $M$  in another order as shown in the proposition below. There is also a version for the semi-definite matrices:  $A \succeq 0 \Leftrightarrow$  there is a square matrix  $M$  s.t.  $A = MM^T$ . (dropping the requirement for  $M$  to be invertible) This is also related to the fact that any positive semi-definite matrix have a square root. The following two tricks are possible with square root:

- If  $A$  is positive semi-definite,  $\mathbf{x}^T A \mathbf{y} = \langle A^{1/2} \mathbf{x}, A^{1/2} \mathbf{y} \rangle$ .
- $\mathbf{x}^T A \mathbf{x} = \|A^{1/2} \mathbf{x}\|^2$

**Proposition 1.3.** For any matrix  $A$  (may not be square),  $A^T A$  is symmetric and positive semi-definite. If given further that  $A$  has full column rank,  $A^T A \succ 0$ .

*Proof.*

$$\mathbf{x}^T A^T A \mathbf{x} = \|A \mathbf{x}\|^2 \geq 0$$

so  $A^T A$  is always positive semi-definite. If  $A$  has full rank,  $A \mathbf{x} = 0$  iff  $\mathbf{x} = 0$ , so  $\mathbf{x}^T A^T A \mathbf{x} > 0$  for non-zero  $\mathbf{x}$ , i.e.  $A^T A \succ 0$ .  $\square$

Outer product of a vector:  $\mathbf{x} \mathbf{x}^T$  is positive semi-definite but NOT definite, because you can always find a vector  $\mathbf{y} \neq 0$  orthogonal to  $\mathbf{x}$ . So  $\mathbf{y}^T \mathbf{x} \mathbf{x}^T \mathbf{y} = 0$ . (unless you are working on one-dimensional space)

There are only two criteria for the indefiniteness of the matrix in this course:

- There is a positive and a negative Eigenvalue.
- There are  $\mathbf{x}, \mathbf{y}$  s.t.  $\mathbf{x}^T A \mathbf{x} > 0$  and  $\mathbf{y}^T A \mathbf{y} < 0$
- Or for  $2 \times 2$  matrix only, proving  $\det(A) < 0$  is enough. (use Corollary 1)

Another criterion not proved in this course is: if there is one positive and one negative element on the diagonal,  $A$  is indefinite.

**Proposition 1.4.** Covariance matrix of any vector of random variables  $\mathbf{X}$  is positive semi-definite.

*Proof.* This becomes obvious using the formulation of covariance matrix  $Q = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$  where  $n$  is the length of random vector  $\mathbf{x}$ ,  $\bar{x}$  is the mean value of  $\mathbf{x}$ .  $\square$

### 1.2.1 Summary of criterion for definiteness and semi-definiteness

- Diagonal dominance(strict):
  - for positive definiteness/semi-definiteness:

$$|A_{ii}| \geq \sum_{j \neq i} |A_{ij}|$$

and  $A_{ii} \geq 0$  for all  $i$ .

- for negative definiteness/semi-definiteness:

$$|A_{ii}| \geq \sum_{j \neq i} |A_{ij}|$$

and  $A_{ii} \leq 0$  for all  $i$ .

- Diagonal dominance(strict):

- for positive definiteness/semi-definiteness:

$$|A_{ii}| > \sum_{j \neq i} |A_{ij}|$$

and  $A_{ii} > 0$  for all  $i$ .

- for negative definiteness/semi-definiteness:  $D_k(A) < 0$  for odd  $k$ ,  $D_k(A) \geq 0$  for even  $k$

$$|A_{ii}| > \sum_{j \neq i} |A_{ij}|$$

and  $A_{ii} < 0$  for all  $i$ .

- Principal minor(non-strict):

- for positive definiteness/semi-definiteness:  $D_k(A) \geq 0$  for all  $k$  where  $D_k$  is the  $k$ th leading principal minor

- for negative definiteness/semi-definiteness:  $D_k(A) \leq 0$  for odd  $k$ ,  $D_k(A) \geq 0$  for even  $k$

- Principal minor(strict):

- for positive definiteness/semi-definiteness:  $D_k(A) > 0$  for all  $k$  where  $D_k$  is the  $k$ th leading principal minor

- for negative definiteness/semi-definiteness:  $D_k(A) < 0$  for odd  $k$ ,  $D_k(A) \geq 0$  for even  $k$

- Eigenvalue(non-strict):

- for positive definiteness/semi-definiteness: all Eigenvalues are non-negative

- for negative definiteness/semi-definiteness: all Eigenvalues are non-positive

- Eigenvalue(strict):

- for positive definiteness/semi-definiteness: all Eigenvalues are positive

- for negative definiteness/semi-definiteness: all Eigenvalues are negative

Criterion	necessary for semi-def.	sufficient for semi-def.	necessary for def.	sufficient for def.
Diagonal dominance(non-strict)		✓		
Diagonal dominance(strict)		✓		✓
Principal minor(non-strict)	✓		✓	
Principal minor(strict)		✓	✓	✓
Eigenvalue(non-strict)	✓	✓		
Eigenvalue(strict)		✓	✓	✓

### 1.3 Transformation

Given  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , you may find the transformation  $g(x) := f(Ax + b)$  where  $x, b \in \mathbb{R}^n$  useful. Gradient and Hessian of  $g$  are given below:

$$\nabla g(x) = A(\nabla f(Ax + b))^T, \quad \nabla^2 g(x) = A^T \nabla^2 f(Ax + b) A$$

these can be proved by chain rule and writing out terms element-wise.

### 1.4 Derivatives

Derivatives of linear, quadratic functions.

$f(x)$	$\nabla f(x)$
$\mathbf{a}^T \mathbf{x} + \beta$	$\mathbf{a}$
$\mathbf{x}^T A \mathbf{x}$	$(A^T + A)\mathbf{x}$

## 2 Unconstrained optimisation

Rule of thumb to remember: minimising  $f$  is equivalent to maximising  $-f$ . You can use this to prove a version of maxima for theorems on minima for free. e.g.  $f$  attains global minimum point if  $f$  is continuous and coercive. But if we know

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

(for this note, I will call this property *anti-coercive*) then  $-f$  is coercive, so  $-f$  attains global minima. i.e.  $f$  attains global maxima. Of course, in an exam, you only have the minimum version from the official note, so you have to use the theorem on  $-f$ .

Regarding second-order optimality conditions, they all assume the following two conditions

- $f$  is twice continuously differentiable on the required domain. So you should be careful with functions involving  $\|\mathbf{x}\|$ .
- $\mathbf{x}^*$  is a stationary point.

Tricky cases of second order optimality conditions are  $\nabla^2 f(\mathbf{x}^*) \succeq 0, \nabla^2 f(\mathbf{x}^*) \preceq 0$ .  $\mathbf{x}^*$  may be local extrema or saddle point. You must use definitions in these cases. i.e. find an open ball around  $\mathbf{x}^*$  s.t.  $f(\mathbf{x}^*)$  is the smallest/largest in that ball, or find two trajectories  $\mathbf{x} = \mathbf{x}(t), y = y(t)$  starting at  $\mathbf{x}^*$  s.t.  $\mathbf{x}^*$  is minima on one trajectory but maxima on the other to prove  $\mathbf{x}^*$  is saddle point.

Optimising on balls: When the domain is  $S = \{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\| \leq r\}$ , i.e. a ball, we need to play tricks on functions. Either differentiate for the interior of  $S$  and consider the boundary separately or write the function as inner product form  $\mathbf{a}^T \mathbf{x}$ . Then maxima is obtained at  $\mathbf{x} = r\mathbf{a}/\|\mathbf{a}\|$  and maximal value is  $r\|\mathbf{a}\|$ . If you have an oval  $\{\mathbf{x} = (x, y) \in \mathbb{R}^2 \mid x^2/a + y^2/b = 1 \leq r\}$  instead, change of variable  $u = x/\sqrt{a}, v = y/\sqrt{b}$  should be performed first.

### 2.1 Global Optimality

Pushing local extrema to global: Assume  $S$  is a non-empty closed set

- If  $f$  is coercive, and you only found one stationary point on  $S$ , it must be global minima (on  $S$ ).
- If  $f$  (defined) is anti-coercive and you only found one stationary point on  $S$ , it must be global maxima (on  $S$ ).

If  $f$  is coercive/anti-coercive, but you cannot find any stationary point for the interior of  $S$ , the extrema must be on the boundary of  $S$ .

Proving local extrema are not global:

- If  $\mathbf{x}^*$  is local minima, either prove  $f$  is not bounded below or find a smaller local minimum to show  $\mathbf{x}^*$  is not global minima.
- If  $\mathbf{x}^*$  is the local maxima, either prove  $f$  is not bounded above or find a larger local maximum to show  $\mathbf{x}^*$  is not global maxima.

### 2.2 Quadratic functions

A quadratic function is in the following form

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

where  $A$  is usually assumed to be symmetric.

Quadratic functions are important as, for example, the linear least square problem of finding  $\mathbf{x}$  s.t.  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$  is minimised and is essentially a quadratic after expanding it.

$$\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 = (\mathbf{A}\mathbf{x} - \mathbf{b})^T (\mathbf{A}\mathbf{x} - \mathbf{b}) = \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{A} \mathbf{x} - \mathbf{x}^T \mathbf{A}^T \mathbf{b} + \mathbf{b}^T \mathbf{b} = \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} - 2\mathbf{b}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{b}$$

So gradient, Hessian, and optimality conditions for quadratic functions are studied in depth.

Roles of  $A, \mathbf{b}, c$  in quadratic functions

- $A$  determines the main behaviours of the function, including whether maxima/minima exists
- $b$  determines location of extrema together with  $A$ , namely  $x$  s.t.  $Ax = -b$
- $c$  only determines  $f$ 's value at extrema.

The Coerciveness of quadratic functions is also determined by  $A$

**Proposition 2.1.** *Quadratic function is coercive iff  $A \succ 0$ .*

*Proof.* By corollary 2,  $x^T Ax \geq \lambda_{\min}(A)\|x\|^2$ , so LHS blows to  $\infty$  if  $\|x\| \rightarrow \infty$ . And since  $x^T Ax$  is the leading term of quadratic function,  $f$  blows to  $\infty$ . (You need to use Cauchy-Schwarz inequality or equivalent to prove this rigorously)

The converse is proved by contradiction, assume  $A$  has negative eigenvalue  $\lambda$  i.e.  $Av = \lambda v$  for some vector  $v$ .

$$f(\alpha v) = \lambda \|\alpha v\|^2 + 2(b^T v)\alpha + c \rightarrow -\infty$$

as  $\alpha \rightarrow -\infty$ . Again you need some inequalities to prove this rigorously, but intuitively  $\lambda \|\alpha v\|^2 < 0$  is the leading term. You also need contradictions for  $A$  having 0 Eigenvalue, left as an exercise.  $\square$

### 3 Gradient Descent

After having a scratch of gradient descent algorithm, or any other iterative algorithm, there is always something to consider

- **Initial values:** this will affect the convergence and even where the iteration arrives. Especially for functions with multiple local minima, different initial values lead to different minima.
- **Details of each iteration:** this depends on the algorithm considered. For gradient descent, the details are the descent direction (given by  $-\nabla f$  as the gradient is essentially the steepest ascent. See the video [Khan Academy](#) for explanation), and step size/learning rate
- **Convergence criteria:** convergence of which variable would you use as criteria for convergence? Criterion is usually in the form  $\|\dots\| < \epsilon$  where  $\epsilon$  is called tolerance. How small tolerance should be is another topic. For gradient descent, criterion can be either  $\|f(x_{k+1}) - f(x_k)\|$  or  $\|\nabla f(x_{k+1})\|$ . The second criterion is more promising as later we'll prove  $\nabla f(x_k) \rightarrow 0$ .

There are three ways to choose a step size

1. **constant step size:** choosing descent direction  $d_k = -\nabla f(x_k)$  only ensures you will descent in this direction in a SMALL neighbourhood of  $x_k$ . So picking a large step size cause failure of decrease, but a small step size will give you slow convergence. Overall, the constant step size is difficult to choose.
2. **Exact line search:** each iteration gradient descent is only focused on one direction  $d_k$ , so it is essentially a 1-dimensional problem. Therefore optimising on this line and finding the step size giving the greatest descent is acceptable.
3. But since greatest descent is not required for convergence, a numerical method called **Backtracking line search**(BLS) can be used instead to find step size giving a "relatively" large descent. The idea is to first pick a larger step-size  $s$ , so you get the point  $x_k + sd_k$  on the line  $x_k + td_k$  ( $t > 0$ ). Let  $t_0 = s$ . Then you gradually shrink this step size by multiplication of a constant  $\beta \in (0, 1)$ . In Armijo's original paper (Armijo, 1966), he used  $\beta = 1/2$ . The value of step sizes while shrinking are  $t_j = \beta^j t_0$ . Descent rate is measured by  $f(x_k) - f(x_k + td_k)$  (Note:  $k$  is considered as fixed when shrinking, we are working in one iteration of gradient descent) Using Maclaurin expansion,

$$f(x_k) - f(x_k + td_k) = -t \nabla f(x_k)^T d_k - o(t\|d_k\|)$$

$-t \nabla f(x_k)^T d_k > 0$  as  $d_k$  is descent direction.  $o(t\|d_k\|) > 0$ , so LHS  $< -t \nabla f(x_k)^T d_k$ . But we want inequality in the other direction, i.e. lower bound for the descent rate. Multiplying constant  $\alpha \in (0, 1)$  to  $t \nabla f(x_k)^T d_k$  gives a milder condition.

$$f(x_k) - f(x_k + td_k) = -\alpha t \nabla f(x_k)^T d_k - (1 - \alpha)t \nabla f(x_k)^T d_k - o(t\|d_k\|)$$

Now using limit

$$\lim_{t \downarrow 0} \frac{(1 - \alpha)t \nabla f(x_k)^T d_k + o(\|d_k\|t)}{t} = (1 - \alpha) \nabla f(x_k)^T d_k < 0$$

therefore, for small enough  $t$ ,

$$f(x_k) - f(x_k + t d_k) \geq -\alpha t \nabla f(x_k)^T d_k$$

This will be used as a criterion to stop shrinking  $t_j$ , i.e. when  $t_j$  satisfies

$$f(x_k) - f(x_k + t_j d_k) \geq -\alpha t_j \nabla f(x_k)^T d_k$$

when  $d_k := -\nabla f(x_k)$ , this is equivalent to

$$f(x_k) - f(x_k + t_j d_k) \geq \alpha t_j \|\nabla f(x_k)\|^2$$

If your initial choice  $s$  already satisfies this, maybe your  $s$  is too small.

The choice of  $\alpha, \beta$  for backtracking line search is another topic to study, out of the scope of this course.

Exact line search for quadratic functions is relative easy to find: assume  $A$  is positive definite,  $f(x) = x^T A x + 2b^T x + c$  where  $b \in \mathbb{R}^n, c \in \mathbb{R}$ . Then the mission is find  $\lim_{t \geq 0} f(x + t d)$  where  $d = -\nabla f = -(2Ax + b)$ . After long and boring expansion of  $f(x + t d)$ ,

$$g(t) := f(x + t d) = (d^T A d)t^2 + 2(d^T A x + d^T b)t + f(x)$$

now  $g'(t) = 2(d^T A d)t + 2d^T (Ax + b)$  so stationary point is

$$t = -\frac{d^T 2(Ax + b)}{2d^T A d} = -\frac{d^T \nabla f(x)}{2d^T A d}$$

$g''(t) = 2(d^T A d) > 0$  for all  $t$  as  $A$  is positive definite. So the answer above gives minima, and  $t \geq 0$ .

If your descent direction happens to be  $d = -\nabla f(x)$ ,

$$t = \frac{d^T d}{2d^T A d}$$

but you need to be careful in cases like Gauss-Newton, Newton's method and scaled gradient descent where the descent direction is NOT  $-\nabla f(x)$ .

### 3.1 Kantorovich inequality

We can define a partial order for symmetric matrices:

**Definition 2** (Loewner Order). For symmetric matrices  $A, B$ ,  $A \prec B$  if  $B - A \succ 0$ , and  $A \preceq B$  if  $B - A \succeq 0$ . Note positive definiteness  $A \succ 0$ , and semi-definite  $A \succeq 0$  are special cases of this partial order.

**Proposition 3.1** (Equivalent condition).  $A \prec B$  iff all Eigenvalues of  $A$  are smaller than that of  $B$  i.e.  $\lambda_i(A) < \lambda_j(B)$  for all  $i, j$ , and  $A \preceq B$  iff  $\lambda_i(A) \leq \lambda_j(B)$  for all  $i, j$ .

**Proposition 3.2.** If  $A \prec B$ , then  $x^T A x \leq x^T B x$  for any vector  $x$ .

*Proof.* Left as an easy exercise. □

We will use matrix order later. For now, consider optimisation of  $f(x) = x^T A x$  where  $A \succ 0$ , which is quadratic, using the step size we discussed in the above session,

$$t_k = \frac{d_k^T d_k}{2d_k^T A d_k}$$

where  $d_k = -2Ax_k$ . Let us find a relation between  $f(x_{k+1})$  and  $f(x_k)$  to study convergence,

$$f(x_{k+1}) = (x_k + t_k d_k)^T A (x_k + t_k d_k)^T$$



after expansion, you get  $x_k^T Ax_k - t_k d_k^T d_k + t_k^2 d_k^T A d_k$ . Then plug in  $t_k$  above gives you

$$x_k^T Ax_k \left( 1 - \frac{1}{4} \frac{(d_k^T d_k)^2}{4(d_k^T A d_k)(x_k^T Ax_k)} \right)$$

but  $x_k^T Ax_k = x_k^T A A^{-1} A x_k = (d_k^T A^{-1} d_k)/4$ , so we can remove  $x_k$  from the above equation

$$= f(x_k) \left( 1 - \frac{1}{4} \frac{(d_k^T d_k)^2}{(d_k^T A d_k)(d_k^T A^{-1} d_k)} \right)$$

To study convergence, we need an upper bound for the ratio 1 – the crazy fraction, i.e. lower bound for the fraction given by Kantorovich inequality.

**Lemma 3.3.** For  $x \neq 0$  and matrix  $A \succ 0$ ,

$$\frac{(x^T x)^2}{(x^T A x)(x^T A^{-1} x)} \geq \frac{4Mm}{(M+m)^2}$$

where  $M := \lambda_{\max}(A)$ ,  $m = \lambda_{\min}(A)$

*Proof.* Eigenvalues for  $A + MmA^{-1}$  are considered, and the spectrum is exactly

$$\left\{ \lambda_i(A) + \frac{Mm}{\lambda_i(A)} : \lambda_i(A) \in \text{Spec}(A) \right\}$$

using calculus,  $\phi(t) = t + Mm/t$  defined on  $[m, M]$  has maximum attained at  $t = M, m$  with value  $M + m$ , so Eigenvalues of  $A + MmA^{-1}$  are smaller than  $M + m$ , i.e.

$$A + MmA^{-1} \preceq (M + m)I$$

so by Proposition 3.2,

$$x^T Ax + Mm(x^T A^{-1} x) \leq (M + m)(x^T x)$$

using the fact that  $\alpha\beta \leq (\alpha + \beta)^2/4$ ,

$$x^T Ax [Mm(x^T A^{-1} x)] \leq \frac{1}{4} [x^T Ax + Mm(x^T A^{-1} x) \leq (M + m)(x^T x)]^2 \leq \frac{(M + m)^2}{4} (x^T x)^2$$

rearranging yields the required inequality.  $\square$

Then you can prove the theorem on the notes that strongly suggest the behaviour of gradient descent on this quadratic function:

**Theorem 3.4.**

$$f(x_{k+1}) \leq \left( \frac{M - m}{M + m} \right)^2 f(x_k)$$

Note dividing by  $m$ , the ratio becomes

$$\left( \frac{\kappa(A) - 1}{\kappa(A) + 1} \right)^2$$

where  $\kappa(A) := M/m = \lambda_{\max}(A)/\lambda_{\min}(A)$  is called condition number. This upper bound becomes small if  $\kappa(A) \approx 1$ , which leads to faster convergence. But larger condition numbers (ill-conditioned) give slower convergence. Similar behaviours appears for general quadratic function as  $A$  is the Hessian matrix for any  $f(x) = x^T A x + 2b^T x + c$ .

A modified gradient descent is to do transformation  $x = S_k y$  so  $g(y) := f(S_k y)$  becomes a new function to optimise (new  $S_k$  can be chosen for each iteration). The advantage is that Hessian of  $g$ , according to transformation rules in section 1.3,

$$\nabla^2 g(y) = S_k \nabla^2 f(S_k y) S_k$$

can be adjusted to have a better condition number. The new descent direction is  $-D_k \nabla f(x_k)$  where  $D_k := S_k S_k^T$  is called *scaling matrix*. Choosing  $D_k = (\nabla^2 f(x_k))^{-1}$  makes sure  $\nabla^2 g(y_k) = S_k \nabla^2 f(S_k y) S_k = I$  as  $S_k = D_k^{1/2}$ .  $\kappa(I) = 1$ , the best conditioned matrix ever. But matrix inversion is very time-consuming. Simpler  $D_k$ , such as a diagonal matrix, works the best. Note with,  $D_{ii} := (\nabla^2 f(x_k))_{ii}^{-1}$ , diagonal elements of  $\nabla^2 g(y_k)$  are 1, so the Hessian will also be well-conditioned.

### 3.2 Fermat-Weber Problem

This problem is called the "facility location problem" imagine you are a land agent trying to plan the location of a new airport or train station  $\mathbf{x}$ , and there are many cities/towns  $\mathbf{a}_i$  around. You seek to ensure the facility is not distanced too far from someone, i.e. minimise

$$f(\mathbf{x}) = \sum_{i=1}^m \|\mathbf{x} - \mathbf{a}_i\|$$

But maybe people in some cities are wealthier, or the population is larger; you want to maximise your benefit by adding weights  $\omega_i$  to the distances.

$$f(\mathbf{x}) = \sum_{i=1}^m \omega_i \|\mathbf{x} - \mathbf{a}_i\|$$

From the official notes, you already know an iterative process (Weiszfeld's method) can be used to find the stationary point of  $f$  i.e.

$$\mathbf{x}_{k+1} = T(\mathbf{x}_k), \quad \text{where } T \text{ is operator } T(\mathbf{x}) = \frac{1}{\sum_i \omega_i / \|\mathbf{x} - \mathbf{a}_i\|} \sum_i \frac{\omega_i \mathbf{a}_i}{\|\mathbf{x} - \mathbf{a}_i\|}$$

But there are several things to consider here:

- The assumption  $\mathbf{x} \neq \mathbf{a}_i$  for any  $i$  is made, how to ensure  $\mathbf{x}_k$  does not become one of  $\mathbf{a}_i$  during iteration.
- Will this iterative process converge?

For the rest of this section, we will study these problems. Though these are not required for exams, the proofs involve useful techniques and provide you with the basic routes of proving convergence.

The first thing is to show  $f(\mathbf{x}_k)$  is non-increasing. We can start by using the fact that the value of  $T(\mathbf{x})$  can be viewed as an optimisation problem of the auxiliary function

$$h(\mathbf{y}, \mathbf{x}) = \sum_{i=1}^m \omega_i \frac{\|\mathbf{y} - \mathbf{a}_i\|^2}{\|\mathbf{x} - \mathbf{a}_i\|}$$

**Lemma 3.5.**  $T(\mathbf{x}) = \text{Argmin}_{\mathbf{y}} \{h(\mathbf{y}, \mathbf{x})\}$

$h(\mathbf{y}, \mathbf{x})$  function has some sweet properties connecting to  $f(\mathbf{x})$

**Lemma 3.6.**  $h(\mathbf{x}, \mathbf{x}) = f(\mathbf{x})$  and  $h(\mathbf{y}, \mathbf{x}) \geq 2f(\mathbf{y}) - f(\mathbf{x})$

*Proof.* Left as an exercise. (Hint: for the target inequality, use the inequality  $a^2/b \geq 2a - b$  for any  $a \geq 0, b > 0$ )  $\square$

Let's prove Lemma 3.5,

*Proof.* Fix  $\mathbf{x}$ ,  $h(\cdot, \mathbf{x})$  can be viewed as quadratic function. After expanding, you would find the leading term is

$$\mathbf{y}^T \underbrace{\left( \sum_i \frac{\omega_i}{\|\mathbf{x} - \mathbf{a}_i\|} \right)}_{=: A} I \mathbf{y}$$

note matrix  $A \succ 0$ . Hence, by proposition on page 27 of lecture notes,  $\mathbf{y}^* := -A^{-1}\mathbf{b}$ , is the unique strict global minimum point. So

$$\nabla_{\mathbf{y}} h(\mathbf{y}^*, \mathbf{x}) = 2 \sum_i \omega_i \frac{\mathbf{y}^* - \mathbf{a}_i}{\|\mathbf{x} - \mathbf{a}_i\|} = 0$$

extracting  $\mathbf{y}^*$  gives  $\mathbf{y}^* = T(\mathbf{x})$   $\square$

**Lemma 3.7.** Let  $\mathbf{x}_k$  be the sequence taken from Weiszfeld's method and  $\mathbf{x}_k \neq \mathbf{a}_i$ . Then

- sequence  $\{f(\mathbf{x}_k)\}$  is non-increasing.

- $f(\mathbf{x}_k) = f(\mathbf{x}_{k+1})$  iff  $\nabla f(\mathbf{x}_k) = 0$ . If decrement is stopped, we have found a stationary point of  $f$ .

*Proof.* For the first statement, just have to prove  $f(T(\mathbf{x})) \leq f(\mathbf{x})$ . By Lemma 3.5,

$$h(T(\mathbf{x}), \mathbf{x}) \leq h(\mathbf{x}, \mathbf{x}) = f(\mathbf{x})$$

On the other hand, by the inequality in Lemma 3.6,

$$h(T(\mathbf{x}), \mathbf{x}) \geq 2f(T(\mathbf{x})) - f(T(\mathbf{x}))$$

so

$$2f(T(\mathbf{x})) - f(T(\mathbf{x})) \leq h(T(\mathbf{x}), \mathbf{x}) \leq f(\mathbf{x}) \quad (*)$$

which yields the required inequality.

For the second statement,  $\nabla f(\mathbf{x}) = 0 \Leftrightarrow \mathbf{x} = T(\mathbf{x})$  so we have to prove  $\mathbf{x} = T(\mathbf{x}) \Leftrightarrow f(\mathbf{x}) = f(T(\mathbf{x}))$ . By (\*), if  $f(\mathbf{x}) = f(T(\mathbf{x}))$ , then

$$h(T(\mathbf{x}), \mathbf{x}) = h(\mathbf{x}, \mathbf{x}) = f(\mathbf{x})$$

by uniqueness of the minimiser  $T(\mathbf{x})$  for  $h(\cdot, \mathbf{x})$ ,  $\mathbf{x} = T(\mathbf{x})$ . □

With the above lemma, we can pick  $\mathbf{x}_0$  s.t.

$$f(\mathbf{x}_0) < \min_i \{f(\mathbf{a}_i)\}$$

then any  $\mathbf{x}_k$  will not be in the set  $\{\mathbf{a}_i\}$ . With this assumption, we can state the main theorem:

**Theorem 3.8.** *Let  $\{\mathbf{x}_k\}$  be the sequence taken from Weizfeld's method and  $\mathbf{x}_0$  satisfies*

$$f(\mathbf{x}_0) < \min_i \{f(\mathbf{a}_i)\}$$

*then any limit point of  $\{\mathbf{x}_k\}$  is stationary point of  $f$ .*

*Proof.* If subsequence  $\{\mathbf{x}_{k_n}\}$  converges to  $\mathbf{x}^*$ , want to show  $\nabla f(\mathbf{x}^*) = 0$ . By monotonicity of sequence  $\{f(\mathbf{x}_k)\}$ ,

$$f(\mathbf{x}^*) \leq f(\mathbf{x}_0) < f(\mathbf{a}_i) \forall i$$

so  $\mathbf{x}^* \neq \mathbf{a}_i$  for any  $i$ ,  $\nabla f(\mathbf{x}^*)$  is well-defined.

Note  $\mathbf{x}_{k_n+1} = T(\mathbf{x}_{k_n}) \rightarrow T(\mathbf{x}^*)$ . Sequence  $\{f(\mathbf{x}_k)\}$  is decreasing and bounded below by 0 so it converges to some  $f^*$ . So all subsequences, including  $\{f(\mathbf{x}_{k_n})\}, \{f(\mathbf{x}_{k_n+1})\}$ , converges to  $f^*$ . By continuity of  $f$ ,  $f(T(\mathbf{x}^*)) = f(\mathbf{x}^*) = f^*$ . Then since  $f(T(\mathbf{x})) = f(\mathbf{x}) \Rightarrow T(\mathbf{x}) = \mathbf{x} \Rightarrow \nabla f(\mathbf{x}) = 0$ ,  $\mathbf{x}^*$  is the stationary point. □

The above theorem can be improved: all limit points of  $\{\mathbf{x}_k\}$  are the global minimum and gives the same value  $f(\mathbf{x}^*)$  where  $\mathbf{x}^*$  is the limit point. (Proof beyond scope)

### 3.3 Convergence of gradient descent

The concept of Lipschitz continuity is fundamental. The official notes have already shown that the existence of a Lipschitz constant  $L$  of the gradient is equivalent to an upper bound of Hessian. A small remark on the notation  $C^{1,1}$ :  $C^1$  is the space of a continuously differentiable function, and it guarantees the existence of a gradient. The extra 1 means the gradient is Lipschitz continuous. Functions in  $C_L^{1,1}$  have another vital property assisting theorems on convergences:

**Lemma 3.9** (Descent lemma).

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

*Proof.* By the fundamental theorem of calculus,

$$f(y) - f(x) = \int_0^1 \langle \nabla f(x + t(y - x)), (y - x) \rangle dt$$

then by the trick of adding 0,

$$f(y) - f(x) = \langle \nabla f(x), (y - x) \rangle + \int_0^1 \langle \nabla f(x + t(y - x)) - \nabla f(x), (y - x) \rangle dt$$

Only have to show

$$\left| \int_0^1 \langle \nabla f(x + t(y - x)) - \nabla f(x), (y - x) \rangle dt \right| \leq \frac{L}{2} \|x - y\|^2$$

This follows a similar process to the proof of the theorem (Equivalence to Boundedness of the Hessian) in official notes.  $\square$

Descent lemma can be used on gradient descent, i.e. substitute  $y = x - t\nabla f(x)$ , we get the following result (*sufficient decrease inequality*),

$$f(x) - f(x - t\nabla f(x)) \geq t \left(1 - \frac{Lt}{2}\right) \|\nabla f(x)\|^2$$

So for constant step size  $t_k \equiv t$ , the descent rate is sufficient by the following inequality

$$f(x_k) - f(x_{k+1}) \geq t \left(1 - \frac{Lt}{2}\right) \|\nabla f(x_k)\|^2$$

Of course, we require  $t \in (0, 2/L)$  for this lower bound to be positive. The maximiser for this lower bound is  $t^* = 1/L$ . And the descent rate bound becomes

$$f(x_k) - f(x_{k+1}) \geq \frac{1}{2L} \|\nabla f(x_k)\|^2$$

which is the bound for exact linear search as exact linear search looks for the fastest descent.

For backtracking, recall an initial step size  $t_0 = s$  is picked and then  $t_k$  is shrunk by factor  $\beta$  each time, until reaching criterion

$$f(x_k) - f(x_k - t_k \nabla f(x_k)) \geq \alpha t_k \|\nabla f(x_k)\|^2$$

Let's find lower bound for  $t_k$ , assuming that it is the first  $t_i$  satisfying above criterion, so  $t_{k-1} = t_k/\beta$  does not satisfy the criterion, i.e.

$$f(x_k) - f(x_k - \frac{t_k}{\beta} \nabla f(x_k)) < \alpha \frac{t_k}{\beta} \|\nabla f(x_k)\|^2$$

but according to sufficient decrease inequality,

$$f(x_k) - f\left(x_k - \frac{t_k}{\beta} \nabla f(x_k)\right) \geq \frac{t_k}{\beta} \left(1 - \frac{Lt_k}{2\beta}\right) \|\nabla f(x_k)\|^2$$

So

$$\begin{aligned} \frac{t_k}{\beta} \left(1 - \frac{Lt_k}{2\beta}\right) &< \alpha \frac{t_k}{\beta} \\ \Rightarrow t_k &> \frac{2(1 - \alpha)\beta}{L} \end{aligned}$$

considering  $t_k \leq s$ , we can get a general lower bound for the descent rate shown below

$$f(x_k) - f(x_k - t_k \nabla f(x_k)) \geq \alpha \min \left\{ s, \frac{2(1 - \alpha)\beta}{L} \right\} \|\nabla f(x_k)\|^2$$

Collecting the results above together:

**Lemma 3.10** (Sufficient Decrease Rate). *The descent rate for three step-size choices is sufficient in the sense of inequality below,*

$$f(x_k) - f(x_k - t_k \nabla f(x_k)) \geq M \|\nabla f(x_k)\|^2$$

where

$$M = \begin{cases} t \left(1 - \frac{tL}{2}\right) & \text{constant size} \\ \frac{1}{2L} & \text{exact line search} \\ \alpha \min \left\{ s, \frac{2(1-\alpha)\beta}{L} \right\} & \text{backtracking} \end{cases}$$

You can see the importance of the Lipschitz constant  $L$  again; it almost determines the convergence rate. Conclusions like  $f(x_{k+1}) = f(x_k)$  iff  $\nabla f(x_k) = 0$  and  $\nabla f(x_k) \rightarrow 0$  follows easily from above lemma. Further, the convergence rate of  $\nabla f(x_k)$  is given by the theorem below

**Theorem 3.11.** *If  $f(x_k) \rightarrow f^*$ , then for any  $n$ ,*

$$\min_{0 \leq k \leq n} \|\nabla f(x_k)\| \leq \sqrt{\frac{f(x_0) - f^*}{M(n+1)}}$$

where  $M$  is the same as that in lemma 3.10.

*Proof.* Left as an exercise. □

### 3.4 Newton's method

We have seen that matrix  $A$  is essential for a quadratic function  $f(x) = x^T A x + 2b^T x + c$ , and when  $A \succ 0$ ,  $x^* = -A^{-1}b$  which is extremely convenient. What about general functions? We can use quadratic approximation! Assume we start with  $x_0$

$$f(x) \approx f(x_0) + \nabla f(x_0)^T (x - x_0) + \frac{1}{2} (x - x_0)^T \nabla^2 f(x_0) (x - x_0)$$

then approximately the minimum satisfies (need to use a change of variable  $y = x - x_0$ )  $\nabla^2 f(x_0)y = -\nabla f(x_0)$  according to properties of quadratic functions.

If luckily  $\nabla^2 f(x) \succ 0$  for all  $x$  (such functions are said to be *concave up*), then the following equation leads  $x_0$  to  $x_1$ , a point closer to the minimum:

$$x_1 = x_0 - (\nabla^2 f(x_0))^{-1} \nabla f(x_0)$$

iterating this process gives pure Newton's method, i.e.

$$x_{n+1} = x_n - (\nabla^2 f(x_n))^{-1} \nabla f(x_n)$$

Note Newton's method for finding the zero of a function  $g(x)$  is

$$x_{k+1} = x_k - \frac{g(x_k)}{g'(x_k)}$$

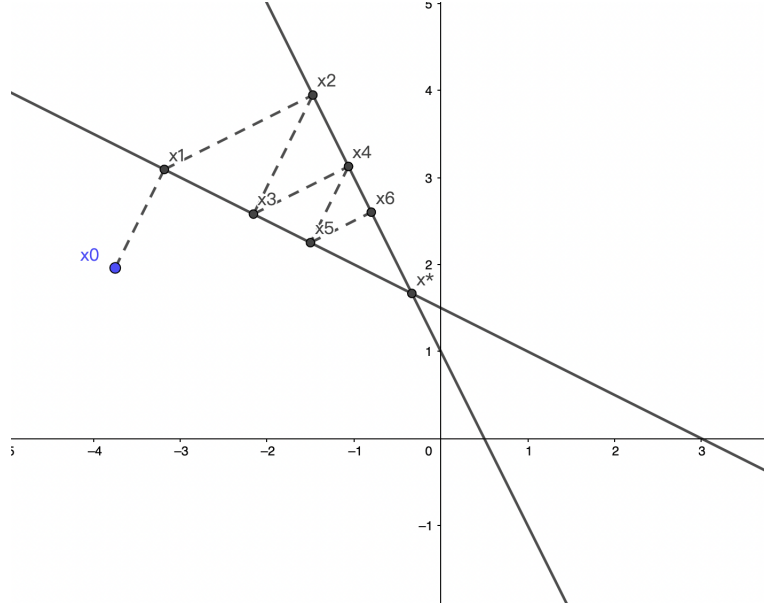
if  $g(x) = \nabla f(x)$ , this is exactly the iteration we derived above.

### 3.5 Kaczmarz method

Suppose given a system of linear equations,  $Ax = b$ , where  $A \in \mathbb{R}^{m \times n}$  i.e. given equations  $a_i^T x = b_i$  where  $a_i^T$  is the  $i$ 'th row of  $A$ . Kaczmarz gives an iterative method which moves  $x$  along one row  $a_i^T$  at a time (Kaczmarz, 1937). Geometrically, his idea is a sequence of projections. Taking  $\mathbb{R}^2$  as an example, we have two lines  $a_1^T x = b_1$ ,  $a_2^T x = b_2$  and want to find the intersection  $x^*$  as shown in figure 1.

Pick  $x_0$  arbitrarily and project it to one line getting  $x_1$ . Continue this process until convergence. From the graph, you can see the convergence is quite fast.

In general, such lines are defined as hyperplanes in the form  $\mathcal{H} = \{x \in \mathbb{R}^n : \langle a, x \rangle = b\}$  where  $a \in \mathbb{R}^n, b \in \mathbb{R}$ . We must find a projection from arbitrary  $x \in \mathbb{R}^n$  to the hyperplane. Note if  $b = 0$ ,  $a$  is the orthogonal direction of the hyperplane, so the projection point should be in the form  $x^* := x - sa/\|a\|$  (here,  $a$  is normalised) for some



**Figure 1: Image representation of Kaczmarz method**

constant  $s \in \mathbb{R}$ . This coefficient  $s$  can be found by the inner product  $\langle a, x \rangle$ . This inner product is the projection length from  $x$  to  $a$ , but the scaled by  $\|a\|$ . So we need to divide  $\|a\|$  here.

$$x^* = x - \frac{\langle a, x \rangle}{\|a\|} \frac{a}{\|a\|}$$

Now if  $b \neq 0$ , simply pick any  $x_0 \in \mathcal{H}$ , i.e.  $\langle a, x_0 \rangle = b$ . That means  $\langle a, x - x_0 \rangle = 0$ . So using change of variable  $u = x - x_0$  (now  $u^* = x^* - x_0$ )

$$u^* = u - \frac{\langle a, u \rangle}{\|a\|^2} a$$

which is equivalent to

$$x^* = x - \frac{\langle a, x \rangle - b}{\|a\|^2} a$$

using this to build iterations:

$$x_{k+1} = x_k - \frac{\langle a, x_k \rangle - b}{\|a\|^2} a$$

will be the mathematical realisation of the algorithm shown in figure 1.

If you let  $i$  go through 1 to  $m$ , i.e. take rows of  $A$  one by one, this algorithm would work fine if  $m$  is small. But what if you are given a matrix of size  $1e6 \times 1e5$  with many 0 entries?

A simple solution is to use the stochastic method, i.e. pick  $i$  randomly according to the uniform distribution on  $\{i\}_{i=1,2,\dots,m}$ . Stochastic methods are beneficial for big data analysis as you can always randomly pick a data point or small data group for each iteration. This would improve efficiency and, surprisingly, usually converges to the desired solution!

Of course, you don't have to stick to the uniform distribution. In the case of the Kaczmarz method, weights can be given to each row according to  $\|a_i\|^2$  as a larger norm means a more significant impact on the solution.

Similarly in gradient descent, if you have an objective function

$$\frac{1}{m} \sum_{i=1}^m Q_i(x)$$

where  $1/m$  is to ensure  $g(x) = E_i[Q_i(x)]$  if uniform distribution is used. If  $m$  is extremely large, you can randomly select  $i$  and let the descent direction be  $-\nabla Q_i(x_k)$ . Or you can select small group  $K \subseteq \{1, \dots, m\}$  and let descent direction be  $-\sum_{i \in K} \nabla Q_i(x_k)$ .

## 4 Convexity

### 4.1 Visual Explanations of some Geometric Shapes

**Hyperplane:**

$$H := \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^T \mathbf{x} = b\} \quad b \in \mathbb{R}, \mathbf{a} \in \mathbb{R}^n$$

it is also mentioned in the Kaczmarz method that you can always pick any point  $x_0$  on the hyperplane and transform the equation into  $\mathbf{a}^T(\mathbf{x} - \mathbf{x}_0) = 0$ , which means vector  $\mathbf{a}$  is the normal direction to the plane.

Using any hyperplane, the space  $\mathbb{R}^n$  can be divided into two half-spaces:

$$H^- := \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^T \mathbf{x} \leq b\}, \quad H^+ := \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^T \mathbf{x} \geq b\}$$

It is easy to prove that  $H, H^-, H^+$  are all convex.

**Ellipsoid:** Our most familiar form of the ellipsoid is

$$\frac{x_1^2}{a_1^2} + \dots + \frac{x_n^2}{a_n^2} = 1$$

you can change any  $x_i^2$  to  $(x_i - c)^2$  to translate the ellipsoid. All ellipsoids can be written in the matrix form

$$\{\mathbf{x} \in \mathbb{R}^n : \mathbf{x}^T Q \mathbf{x} + 2\mathbf{b}^T \mathbf{x} + c \leq 0\}$$

where  $Q$  must be positive semi-definite, because by Theorem 1.1, any symmetric matrix  $Q$  can be decomposed into  $UDU^T$  where  $D$  is diagonal. So  $\mathbf{x}^T Q \mathbf{x} = \mathbf{x}^T U D U^T \mathbf{x} = \mathbf{y}^T D \mathbf{y}$  where  $\mathbf{y} := U^T \mathbf{x}$ . But the orthogonal matrix  $U^T$  is a rigid transformation, i.e. shape of objects will not be changed. So the shape of  $\mathbf{x}^T Q \mathbf{x} - 1 = 0$  will be the same as  $\mathbf{y}^T D \mathbf{y} - 1 = 0$  which can be written as  $\sum_i d_i y_i^2 = 1$ . Hence, if the object is ellipsoid, then  $d_i \geq 0$ , i.e.  $Q$  is positive semi-definite. The term  $2\mathbf{b}^T \mathbf{x}$  will not affect the shape but will move the centre of the ellipsoid away from the origin.

### 4.2 Properties of Convex set

The intersection of arbitrary convex sets is convex by the nature of the definition. But the union of convex sets may not be convex. The easiest example is the union of two distinct lines.

The image and pre-image of linear transformation on a convex set are still convex, i.e. if  $C$  is convex,  $A$  is any matrix, the following two sets are convex:

$$A(C) := \{A\mathbf{x} + \mathbf{b} : \mathbf{x} \in C\}, \quad A^{-1}(C) := \{\mathbf{y} : A\mathbf{y} + \mathbf{b} \in C\}$$

And any linear combination and Cartesian product of convex sets are convex:

$$\mu_1 C_1 + \dots + \mu_k C_k, \quad C_1 \times \dots \times C_k$$

are convex.

The Interior and closure of a convex set are convex, as formulated in the following theorems

**Theorem 4.1** (Preservation under Closure). *If  $C$  is convex set, closure  $\overline{C}$  is convex*

*Proof.* Easy proof using the sequential definition of closure. i.e. if  $x \in \overline{C}$ , exists  $(x_k)_{k \geq 1} \subseteq C$  s.t.  $x_k \rightarrow x$  as  $k \rightarrow \infty$ .  $\square$

**Theorem 4.2** (Preservation under Interior). *If  $C$  is convex, then interior  $\text{int}(C)$  is convex.*

*Proof.* In the case,  $\text{int}(C) = \emptyset$ , the theorem is trivially true.

If  $\text{int}(C) \neq \emptyset$ , it is easier to prove another version: if  $x \in \text{int}(C)$ ,  $y \in \overline{C}$ , then  $(1-\lambda)x + \lambda y \in \text{int}(C)$  for all  $\lambda \in (0, 1)$ . This is called *line segment principle*, see proof in (Beck 2014, p.109).  $\square$

Using the line segment principle, it can be proved that for the convex set,  $\overline{\text{int}(C)} = \overline{C}$  and  $\text{int}(\overline{C}) = \text{int}(C)$ .

Recall Carathéodory theorem says that elements in a convex hull in  $\mathbb{R}^n$  can be written as a linear combination of at most  $n+1$  vectors in  $S$ . The additional one dimension is used to amend for requirement on the coefficients:  $\lambda \in \Delta_{n+1}$  (this means the sum of entries of  $\lambda$  is 1). In fact, on the conic hull, where this requirement on coefficients is dropped and they can be any positive numbers, any element in the conic hull of  $S$  can be written as a linear combination of at most  $n$  linearly independent vectors  $x_i$  in  $S$ . This is called the *conic representation theorem*.

The convex hull of the closed set may not be closed, e.g.

$$S = \{(0, 0)\} \cup \{(x, y) : xy \geq 1, x \geq 0, y \geq 0\}$$

$S$  is closed, but its hull is not closed:

$$\text{conv}(S) = \{(0, 0)\} \cup \{(x, y) : x > 0, y > 0\}$$

But compactness can be preserved

**Theorem 4.3** (Preservation of compactness). *If  $S$  is compact, then  $\text{conv}(S)$  is compact.*

*Proof.* .

**Step 1. Boundedness**

$S$  is compact so there is  $M > 0$  s.t.  $\|x\| \leq M \forall x \in S$ . We want to find a bound for  $\text{conv}(S)$ , and Carathéodory theorem becomes useful now: for all  $y \in \text{conv}(S)$ , exists  $x_1, \dots, x_{n+1} \in S$  and  $\lambda \in \Delta_{n+1}$  (i.e.  $\sum_{i=1}^{n+1} \lambda_i = 1$ ) s.t.  $y = \sum_{i=1}^{n+1} \lambda_i x_i$ . Therefore,

$$\|y\| \leq \sum_{i=1}^{n+1} \lambda_i \|x_i\| \leq M \sum_{i=1}^{n+1} \lambda_i = M$$

it turns out that  $\text{conv}(S)$  is bounded by the same bound of  $S$ .

**Step 2. Closedness**

Given sequence  $(y^k)_{k \geq 1} \subseteq \text{conv}(S)$  converging to  $y \in \mathbb{R}^n$ , the target is to prove  $y \in \text{conv}(S)$ . Again by Carathéodory theorem, for any  $k$ , there are  $x_1^k, \dots, x_{n+1}^k \in S$  and  $\lambda^k \in \Delta_{n+1}$  s.t.

$$y^k = \sum_{i=1}^{n+1} \lambda_i^k x_i^k$$

$S$  is compact so  $((\lambda^k, x_1^k, \dots, x_{n+1}^k))_{k \geq 1}$  have convergent subsequence, say the limit is  $(\lambda, x_1, \dots, x_{n+1})$ . Taking limit of the equation  $y^k = \sum_{i=1}^{n+1} \lambda_i^k x_i^k$  on the convergent subsequence (i.e. letting  $k = k_j$  where  $k_j$  are indices of the subsequence, and send  $j \rightarrow \infty$ ) yields

$$y = \sum_{i=1}^{n+1} \lambda_i x_i$$

$\square$

## 4.3 Convex Functions

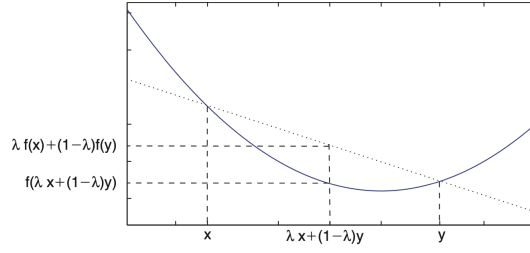
Convex functions are defined by the *Fundamental inequality*

$$f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y), \quad \forall x, y \in C, \lambda \in [0, 1]$$

where  $C$  is a convex set. This can be visually understood in figure 2 taken from (page 117, Beck, 2014)

If the fundamental inequality becomes a strict inequality in the above definition, the function is called a *strictly convex* function.





**Figure 2: Graph representation of the fundamental inequality**

An important property of convex functions, Jensen's inequality, basically says that this inequality generalises to any linear combination of vectors in  $C$  as long as coefficients add up to 1. The proof is done by induction on the number of vectors in linear combination and proceeds in a similar way to proving if  $C$  is convex, then any linear combination of  $m$  vectors in  $C$  is still in  $C$  if  $\lambda \in \Delta_m$ .

**Exercise 1.** Using the convex function  $f(x) = -\ln x$  defined on  $\{x \in \mathbb{R} : x > 0\}$  (it is convex because second derivative always positive) and Jensen's inequality, prove the AGM inequality

$$\frac{\sum_i x_i}{n} \geq \sqrt[n]{\prod_i x_i}$$

#### 4.4 First, second order characterisations of convex functions

In this section, equivalent conditions (characterisations) for the convexity of a continuously differentiable function are considered, based on the first and second derivatives. Continuously differentiability ensures the existence of derivatives.

Staring at figure 2, you may notice that if we draw the tangent line (or tangent hyperplane in higher dimension), then the tangent line is always below the function.

**Theorem 4.4** (Gradient inequality). *If  $f$  is continuously differentiable function, then  $f$  is convex over  $C$*

$$\Leftrightarrow f(x) + \nabla f(x)^T(y - x) \leq f(y) \quad \forall x, y \in C$$

$f(x) + \nabla f(x)^T(y - x)$  is exactly the tangent hyperplane taken at  $f(x)$ .

There is another characterisation for convexity of continuously differentiable functions, from Figure 2, it can be seen that for convex functions on  $\mathbb{R}$ , the derivative is always non-decreasing. In a higher dimension, this characterisation is stated below

**Theorem 4.5** (Monotonicity of gradient). *If  $f : C \rightarrow \mathbb{R}$  is continuously differentiable function defined on convex set  $C \subseteq \mathbb{R}^n$ , then  $f$  is convex iff*

$$(\nabla f(x) - \nabla f(y))^T(x - y) \geq 0 \quad \forall x, y \in C$$

*Proof.* ( $\Rightarrow$ ) First assume  $f$  is convex, then by gradient inequality,

$$f(x) \geq f(y) + \nabla f(y)^T(x - y)$$

$$f(y) \geq f(x) + \nabla f(x)^T(y - x)$$

And adding two inequalities yields the required inequality.

( $\Leftarrow$ ) Assume  $f$  satisfies the inequality (i.e. gradient is monotonic), we aim to prove the gradient inequality. Construct function  $g(t) := f(x + t(y - x))$  which are the values of  $f$  on the line joining  $x, y$ . Then

$$f(y) = g(1) = g(0) + \int_0^1 g'(t) dt$$

then by chain rule,  $g'(t) = (y - x)^T \nabla f(x + t(y - x))$ , so

$$\begin{aligned} f(y) &= f(x) + \int_0^1 (y - x)^T \nabla f(x + t(y - x)) dt \\ &= f(x) + \nabla f(x)^T (y - x) + \int_0^1 (y - x)^T (\nabla f(x + t(y - x)) - \nabla f(x)) dt \end{aligned}$$

The integrand is positive as it can be rearranged in the following way

$$(y - x)^T (\nabla f(x + t(y - x)) - \nabla f(x)) = \frac{1}{t} (\nabla f(x + t(y - x)) - \nabla f(x))^T (x + t(y - x) - x) \geq 0$$

that yields

$$f(y) \geq f(x) + \nabla f(x)^T (y - x)$$

□

For continuously differentiable convex functions defined over a convex set, any stationary point is a global minimiser. (the example  $f(x) = x^2$  illustrates this point)

**Proposition 4.6** (Sufficient global optimality condition of convex function). *If  $f$  defined on  $C \subseteq \mathbb{R}^n$  is convex and continuously differentiable. Assume  $\nabla f(x^*) = 0$ , then  $x^*$  is global minimiser of  $f$  over  $C$ .*

*Proof.* Given  $z \in C$ , using the gradient inequality,

$$f(z) \geq f(x^*) + \nabla f(x^*)^T (z - x^*)$$

which implies  $f(z) \geq f(x^*)$ . So  $x^*$  must be a global minimiser. □

However, the condition above is necessary only when  $C = \mathbb{R}^n$ , the entire space. Because when  $C$  is not the entire space, the minimiser may appear at the boundaries.

Recall that in 1-D case, the gradient inequality (gradient is non-decreasing) is equivalent to  $f''(x) > 0$  for all  $x$ . This is a more handy characterisation for checking whether a function is convex. There is a similar characterisation in higher dimension

**Theorem 4.7** (Second-order characterisation). *If  $f : C \rightarrow \mathbb{R}$  ( $C \subseteq \mathbb{R}^n$  is open convex) is twice continuously differentiable, then  $f$  convex  $\Leftrightarrow \nabla^2 f(x) \succeq 0 \forall x \in C$ .*

*Proof.* Gradient inequality will be the bridge to prove this theorem.

Assume  $\nabla^2 f(x) \succeq 0$ , using linear approximation of  $f$ , given  $x, y \in C$  there is  $z \in [x, y] \cap C$  s.t.

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

the last term is non-negative as  $\nabla^2 f(x) \succeq 0$ , so gradient inequality is proved.

For the opposite direction, assuming  $f$  is convex, aim to show the positive semi-definiteness of Hessian. The openness of  $C$  is required, as we need to do quadratic approximation in a neighbourhood of  $x$  and this is only possible with  $C$  open. So given  $x \in C$  and  $y \in \mathbb{R}^n$ , there is small enough  $\epsilon$  s.t.  $\forall \lambda \in (0, \epsilon)$ ,  $x + \lambda y \in C$ . (Treat  $y$  as a direction here, instead of a point) Then quadratic approximation gives

$$f(x + \lambda y) = f(x) + \nabla f(x)^T \lambda y + \frac{\lambda^2}{2} y^T \nabla^2 f(x) y + o(\lambda^2 \|y\|^2)$$

Using gradient inequality,

$$f(x + \lambda y) \geq f(x) + \nabla f(x)^T \lambda y$$

so combining the above two results,

$$\frac{\lambda^2}{2} y^T \nabla^2 f(x) y + o(\lambda^2 \|y\|^2) \geq 0$$

dividing by  $\lambda^2$  and sending  $\lambda \rightarrow 0^+$  yields

$$y^T \nabla^2 f(x) y \geq 0$$

$y$  was chosen arbitrarily, so  $\nabla^2 f(x) \succeq 0$ . □

You can easily see that a small modification to the above proof gives

**Theorem 4.8.** *If  $f : C \rightarrow \mathbb{R}$  ( $C \subseteq \mathbb{R}^n$  is open convex) is twice continuously differentiable, then  $f$  strictly convex  $\Leftrightarrow \nabla^2 f(x) \succ 0 \forall x \in C$ .*

## 4.5 Properties of Convex Functions

It is straightforward to prove with definitions that linear combinations of convex functions are still convex as long as the coefficients are positive. And the linear transformation of independent variable  $y := Ax + b$  will preserve convexity i.e. if  $f$  is convex over convex set  $C$ , then  $g(x) := f(y) = f(Ax + b)$  is convex over  $D := \{x : Ax + b \in C\}$ . Note  $D$  is convex because it is pre-image of a convex set under linear transformation.

Using the above properties, it can be proved that

$$g(x) := \frac{\|Ax + b\|^2}{c^T x + d}$$

is convex over  $\{x : c^T x + d > 0\}$  because  $h(y, t) := \|y\|^2/t$  is convex over  $\{(y, t) : t > 0\}$ .  $g$  can be obtained from  $h$  by two consecutive linear transformations  $y \mapsto Ax + b$  and  $t \mapsto c^T x + d$ . As for the convexity of  $h(y, t)$ , note  $h = \sum_i y_i^2/t$ . The function  $\phi(x, z) := x^2/z$  defined on  $\{(x, z) : x \in \mathbb{R}, z > 0\}$  is convex. (Left as an exercise)

Composition of convex functions may not be convex

**Example 1.**  $g(x) = x^2$ ,  $f(x) = x^2 - 4$  are both convex but  $g(f(x)) = (x^2 - 4)^2$  is not convex. (check the second derivatives of these functions)

But if  $g : I \rightarrow \mathbb{R}$  (where  $I$  is interval on  $\mathbb{R}$  with  $f(C) \subseteq I$ ) is non-decreasing convex function,  $h(x) := g(f(x))$  is convex for any convex function  $f : C \rightarrow \mathbb{R}$ . (Proof left as an exercise) Note in fact,  $g$  only has to be non-decreasing on  $f(C)$ .

Using this composition rule, one can show that  $e^{\|x\|^2}$  is convex.

Using the famous inequality  $\max_i(a_i + b_i) \leq \max_i a_i + \max_i b_i$ , one can prove the following

**Theorem 4.9** (Pointwise Maximum). *If  $f_i : C \rightarrow \mathbb{R}$  (for  $i \in I$ , arbitrary index set) are convex functions, then  $f(x) := \max_i f_i(x)$  is convex.*

Warning: the function  $h(x) = x_{[i]}$  where  $x_{[i]}$  is the  $i$ -th largest entry in  $x$  is NOT convex. But the following function is convex:

$$h_k(x) := x_{[1]} + x_{[2]} + \dots + x_{[k]} := \max\{x_{i_1} + \dots + x_{i_k} : i_j \text{ are distinct indices}\}$$

Convexity is also preserved when some entries of input  $x$  are minimised over a convex set

**Theorem 4.10** (Partial Minimisation). *If  $f : C \times D \rightarrow \mathbb{R}$  is convex function where  $C, D$  are convex,*

$$g(x) := \min_{y \in D} f(x, y)$$

*is convex. Note min should really be inf here, and we assume infimum exists and is finite.*

*Proof.* Given  $x_1, x_2 \in C, \lambda \in [0, 1]$ . Minimum may not be attained by some  $y$  so we need to give a room of  $\epsilon$ : given  $\epsilon > 0$ , there are  $y_1, y_2 \in D$  s.t.

$$f(x_1, y_1) \leq g(x_1) + \epsilon \quad f(x_2, y_2) \leq g(x_2) + \epsilon$$

By convexity of  $f$ ,

$$\begin{aligned} g(\lambda x_1 + (1 - \lambda)x_2) &\leq f(\lambda x_1 + (1 - \lambda)x_2, \lambda y_1 + (1 - \lambda)y_2) \leq \lambda f(x_1, y_1) + (1 - \lambda)f(x_2, y_2) \\ &\leq \lambda(g(x_1) + \epsilon) + (1 - \lambda)(g(x_2) + \epsilon) = \lambda g(x_1) + (1 - \lambda)g(x_2) + \epsilon \end{aligned}$$

The above inequality holds for all  $\epsilon$ , so sending  $\epsilon$  to 0 yields convexity of  $g$ . □

*Remark.* Note if  $f$  is convex,

$$g(x) := \max_{y \in D} f(x, y)$$

$g$  will also be convex over  $C$ . Because we can treat  $y$  as index set and apply Theorem 4.9. Both partial and full maximisation of convex function is convex, but only partial minimisation (over convex set) of convex function is convex. This may not hold for arbitrary minimisation.

### Continuity and differentiability:

Convex functions are continuous on open convex sets, but may not be continuous otherwise. However, convex functions are locally Lipschitz continuous for points in the interior of  $C$ . i.e. if  $x_0 \in \text{int}(C)$ , then exists  $\epsilon > 0, L > 0$  with  $B_\epsilon(x_0) \subseteq C$  and

$$|f(x) - f(x_0)| \leq L \|x - x_0\| \quad \forall x \in B_\epsilon(x_0)$$

Convex functions may not be differentiable, but all directional derivatives exist at interior points, which is enough for most optimisation problems.

See detailed, but rather long proofs of the above results in (Beck 2014, Section 7.6).

## 4.6 Level sets

**Definition 3.** Level set of  $f$  at level  $\alpha$  is

$$\text{Lev}(f, \alpha) = \{x : f(x) \leq \alpha\}$$

**Theorem 4.11.** If  $f : C \rightarrow \mathbb{R}$  is convex function, for all  $\alpha \in \mathbb{R}$ ,  $\text{Lev}(f, \alpha)$  is convex set.

*Proof.* Given  $x, y \in \text{Lev}(f, \alpha)$ , i.e.  $f(x), f(y) \leq \alpha$ . Using convexity,

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) \leq \alpha$$

that means  $\lambda x + (1 - \lambda)y \in \text{Lev}(f, \alpha)$ , so  $\text{Lev}(f, \alpha)$  is convex. □

There is a special class of functions

**Definition 4** (Quasi-convex functions).  $f : C \rightarrow \mathbb{R}$  is quasi-convex if for all  $\alpha$ ,  $\text{Lev}(f, \alpha)$  is convex set.

Convex functions are all quasi-convex, but there are non-convex functions in this class, e.g.  $f(x) = \sqrt{|x|}$  and

$$f(x) = \frac{a^T x + b}{c^T x + d} \quad c \neq 0, \text{ defined on } \{x : c^T x + d > 0\}$$

## 4.7 Maxima of convex functions

We have seen that finding minima for convex functions is very convenient, as any stationary point is a global minimiser. Maxima are even easier to deal with, given the fact that they cannot exist in the interior:

**Theorem 4.12.** If  $f : C \rightarrow \mathbb{R}$  is non-constant convex function,  $f$  does not attain maximum in  $\text{int}(C)$

*Proof.* If  $x^* \in \text{int}(C)$  is global maximiser, by non-constant property, there is  $y \in C$  s.t.  $f(y) < f(x^*)$ . There is  $\epsilon > 0$  s.t.  $z := x^* + \epsilon(x^* - y) \in C$ . Note  $x^* = \epsilon/(\epsilon + 1)y + 1/(1 + \epsilon)z$ , so by convexity,

$$\begin{aligned} f(x^*) &\leq \frac{\epsilon}{\epsilon + 1} f(y) + \frac{1}{\epsilon + 1} f(z) \\ \Rightarrow f(z) &\geq \epsilon(f(x^*) - f(y)) + f(x^*) > f(x^*) \end{aligned}$$

Contradiction as  $x^*$  is a global maximiser. □

As you may have guessed, maximisers of convex functions will only appear on the boundary of convex sets, and surprisingly, one of them must be an extreme point.

**Theorem 4.13.** If  $f : C \rightarrow \mathbb{R}$  is convex, continuous and  $C$  is convex, compact set, then there is at least one maximiser that is an extreme point.

*Proof.* By the Weierstrass theorem, there must be a maximiser  $x^*$  of  $f$  over  $C$  as  $C$  is compact. Assume  $x^*$  is not extreme point, by Krein-Milman,  $C = \text{conv}(\text{ext}(C))$ , by definition, there are  $x_1, \dots, x_k \in \text{ext}(C)$ ,  $\lambda \in \Delta_k$  s.t.

$$x^* = \sum_{i=1}^k \lambda_i x_i$$

So by the convexity of  $f$ ,

$$f(x^*) \leq \sum_{i=1}^k \lambda_i f(x_i) \Rightarrow \sum_i \lambda_i (f(x_i) - f(x^*)) \geq 0$$

but  $x^*$  is maximiser,  $f(x_i) \leq f(x^*)$ , so the sum

$$\sum_i \lambda_i (f(x_i) - f(x^*)) \leq 0$$

that means the sum is indeed 0 and  $f(x_i) = f(x^*)$ . So  $x_i$  are all maximisers.  $\square$

The theorem also tells us that if there is only one maximiser, then it must be an extreme point.

**Example 2.** Consider optimisation problem

$$\max\{\|Ax\| : \|x\|_1 \leq 1\}$$

where  $A$  is a matrix. The function  $x \mapsto \|Ax\|_1$  is convex as 1-norm and matrix-multiplication is convex. Further,  $\{x : \|x\|_1 \leq 1\}$  is compact convex set. So there must be a maximiser among extreme points:  $\pm e_i$  for  $i = 1, \dots, n$ . Comparing the values at these extreme points, noting  $\|Ae_j\| = \|A(-e_j)\|$ , we have

$$\max\{\|Ax\| : \|x\|_1 \leq 1\} = \max_j \|Ae_j\|_1$$

which means the matrix norm  $\|A\|_{1,1}$  is the maximum column sum of absolute values of the entries.

## 5 Convex Optimisation

Convex optimisation problems are problems of the form

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

where  $f, g_i$  are convex,  $h_j$  are affine functions.

Note  $h_j$  being convex is not enough for the set  $\{x \in \mathbb{R}^n : h_j(x) = 0\}$  to be convex. e.g. imagine  $h(x) = \|x\|_2^2 - 1$ , then the set above is the unit circle, not convex. But if  $h$  is affine, the set above is convex. (Prove it)

The feasible set of this problem is

$$\left( \bigcap_{i=1}^m \text{Lev}(g_i, 0) \right) \cap \left( \bigcap_{j=1}^p \{x \in \mathbb{R}^n : h_j(x) = 0\} \right)$$

which is a convex set.

The local minimum point of convex functions on a convex set must be a global minimum, graphically this is due to the shape of convex functions: if  $x^*$  is minimum the shape of  $f$  around  $x = x^*$  would look like a bowl. But the convex function has this shape everywhere, so  $f(x^*)$  will be smaller than the value at any other point.

**Theorem 5.1.** For convex function  $f$ , convex set  $C$ , the set of optimal solutions

$$X^* := \min \{f(\mathbf{x}) : \mathbf{x} \in C\}$$

is convex and if  $f$  is strictly convex, there is at most one optimal solution.

*Proof.* The first part is easy. Now assume  $f$  is strictly convex and  $x, y \in X^*$  ( $x \neq y$ ) and let  $f^* := f(x) = f(y)$ , then

$$f\left(\frac{1}{2}x + \frac{1}{2}y\right) < \frac{1}{2}f(x) + \frac{1}{2}f(y) = f^*$$

contradicts the minimality of  $f^*$ .  $\square$

Linear programming is a special case of convex optimisation where objective function  $f(\mathbf{x})$  is linear,  $g_i(\mathbf{x}) = \mathbf{a}_i^T \mathbf{x} - c_i$  and  $h_j(\mathbf{x}) = \mathbf{b}_j^T \mathbf{x} - d_j$ , so the constraints can be summarised to

$$A\mathbf{x} \leq \mathbf{c}, \quad B\mathbf{x} = \mathbf{d}$$

Linear functions are convex and concave, so theorems on convex minimisation problems apply to both minimisation and maximisation linear problems. (Maximisation of concave function  $f$  over convex set  $C$  is equivalent to minimisation of convex function  $-f$  over  $C$ )

When the objective is a quadratic function, the problem is called a quadratic problem. And if the constraint is also quadratic, we call it QCQP (quadratically constrained quadratic problem). It is not necessarily convex unless all matrices are positive and semi-definite and all constraints are inequality.

## 5.1 Classification using linear separator

One should be able to convert a practical problem into a convex optimisation problem, e.g. suppose given two types of data  $\mathbf{x}_1, \dots, \mathbf{x}_m$  (type 1) and  $\mathbf{x}_{m+1}, \mathbf{x}_{m+2}, \dots, \mathbf{x}_{m+p}$  (type 2) and we wish to find a hyperplane  $\{\mathbf{x} : \mathbf{w}^T \mathbf{x} + \beta = 0\}$  that separates the two sets of data. i.e.

$$\mathbf{w}^T \mathbf{x}_i + \beta < 0 \text{ for type 1, } \mathbf{w}^T \mathbf{x}_i + \beta > 0 \text{ for type 2}$$

and that the hyperplane is kept as far from all points as possible, i.e. the problem is

$$\begin{aligned} \max \left\{ \min_{i=1,2,\dots,m+p} \frac{|\mathbf{w}^T \mathbf{x}_i + \beta|}{\|\mathbf{w}\|} \right\} \\ \mathbf{w}^T \mathbf{x}_i + \beta < 0, \quad i = 1, \dots, m \\ \mathbf{w}^T \mathbf{x}_i + \beta > 0, \quad i = m+1, \dots, m+p \end{aligned}$$

the expression

$$\frac{|\mathbf{w}^T \mathbf{x}_i + \beta|}{\|\mathbf{w}\|}$$

represents the minimal distance from the hyperplane to point  $\mathbf{x}_i$ , so the first line means "maximising the minimum distance from  $\mathbf{x}_i$  to the hyperplane". All optimal solutions can be scaled, i.e. if  $\alpha \neq 0$ ,  $(\alpha \mathbf{w}, \alpha \beta)$  is optimal if  $(\mathbf{w}, \beta)$  is optimal. (because  $\{\mathbf{x} : \alpha \mathbf{w}^T \mathbf{x} + \alpha \beta = 0\} = 0$  still represents the same hyperplane) So the constraint

$$\min_i |\mathbf{w}^T \mathbf{x}_i + \beta| = 1$$

can be added to the problem without changing the result. But this new constraint is equivalent to

$$\mathbf{w}^T \mathbf{x}_i + \beta \leq -1, \quad i = 1, \dots, m$$

$$\mathbf{w}^T \mathbf{x}_i + \beta \geq 1, \quad i = m+1, \dots, m+p$$

and objective  $f$  becomes  $1/\|\mathbf{w}\|$ . Maximising this is equivalent to minimising  $\frac{1}{2}\|\mathbf{w}\|^2$ . Now the problem

$$\begin{aligned} \min \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & \mathbf{w}^T \mathbf{x}_i + \beta \leq -1, \quad i = 1, \dots, m \\ & \mathbf{w}^T \mathbf{x}_i + \beta \geq 1, \quad i = m+1, \dots, m+p \end{aligned}$$

is a convex optimisation problem.

## 5.2 Stationary points

For continuously differentiable  $f$  defined over closed convex function  $C$ ,  $\mathbf{x}^* \in C$  is stationary point if for all other  $\mathbf{x} \in C$ , the directional derivative along  $\mathbf{x} - \mathbf{x}^*$

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

that means along any direction from  $\mathbf{x}^*$ ,  $f$  has a positive gradient. Stationary points are important because if given  $f$  is continuously differentiable convex function,  $C$  is non-empty closed and convex set, for convex optimisation problem

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{x} \in C \end{aligned}$$

$\mathbf{x}^*$  is optimal solution iff  $\mathbf{x}^*$  is stationary point.

As shown in lectures, when  $C = \mathbb{R}^n$ ,  $\mathbf{x}$  is stationary point iff  $\nabla f(\mathbf{x}) = 0$ . Let us consider three other cases

**Example 3.** Let  $C := \{\mathbf{x} : x_i \geq 0 \ \forall i\}$ , stationary condition can be written as

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

which is equivalent to the problem  $\mathbf{a}^T \mathbf{x} + b \geq 0$  for all  $\mathbf{x} \geq 0$  where  $\mathbf{a} = \nabla f(\mathbf{x}^*)$ ,  $b = -(\nabla f(\mathbf{x}^*)^T \mathbf{x}^*)$ . Let  $\mathbf{x} = 0$ , we can show  $b \geq 0$ . Let  $\mathbf{x} = (1 - \frac{b}{a_i})\mathbf{e}_i \geq 0$  (it is always positive as  $b, a_i$  has opposing signs), then

$$\mathbf{a}^T \mathbf{x} + b = a_i - b + b = a_i \geq 0$$

so  $a_i \geq 0$  for all  $i$ . Therefore, the stationary condition is equivalent to

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

but  $\mathbf{x}^* \geq 0$ , so the latter condition is in fact equality and cleaning up, we have

$$\frac{\partial f}{\partial x_i}(\mathbf{x}^*) \begin{cases} = 0 & \text{if } x_i^* > 0 \\ \geq 0 & \text{if } x_i^* = 0 \end{cases}$$

**Example 4.** Let  $C := \{\mathbf{x} : \mathbf{e}^T \mathbf{x} = \sum_i x_i = 1\}$  (*unit-sum set*) where  $\mathbf{e}$  is the vector with all entries being 1. Stationary condition on  $\mathbf{x}^* \in C$  is equivalent to

$$\frac{\partial f}{\partial x_i}(\mathbf{x}^*) = \frac{\partial f}{\partial x_j}(\mathbf{x}^*) \quad \forall i, j$$

Proof: If the above hold, for all  $\mathbf{x} \in C$ ,

$$\nabla f(\mathbf{x}^*)^T (\mathbf{x} - \mathbf{x}^*) = \sum_{i=1}^n \frac{\partial f}{\partial x_i}(\mathbf{x}^*) (x_i - x_i^*) = \frac{\partial f}{\partial x_1}(\mathbf{x}^*) \left( \sum_i x_i - \sum_i x_i^* \right) = 0$$

Conversely, assume stationary condition holds but for some  $i, j$

$$\frac{\partial f}{\partial x_i}(\mathbf{x}^*) > \frac{\partial f}{\partial x_j}(\mathbf{x}^*)$$

then define a new vector  $\mathbf{x}$  as follows

$$x_i = x_i^* + 1, x_j = x_j^* - 1, x_k = x_k^* \text{ if } k \notin \{i, j\}$$

if  $\mathbf{x}^* \in C$  then  $\mathbf{x} \in C$ . And

$$\begin{aligned} \nabla f(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) &= \frac{\partial f}{\partial x_i}(\mathbf{x}^*)(x_i - x_i^*) + \frac{\partial f}{\partial x_j}(\mathbf{x}^*)(x_j - x_j^*) \\ &= -\frac{\partial f}{\partial x_i}(\mathbf{x}^*) + \frac{\partial f}{\partial x_j}(\mathbf{x}^*) < 0 \end{aligned}$$

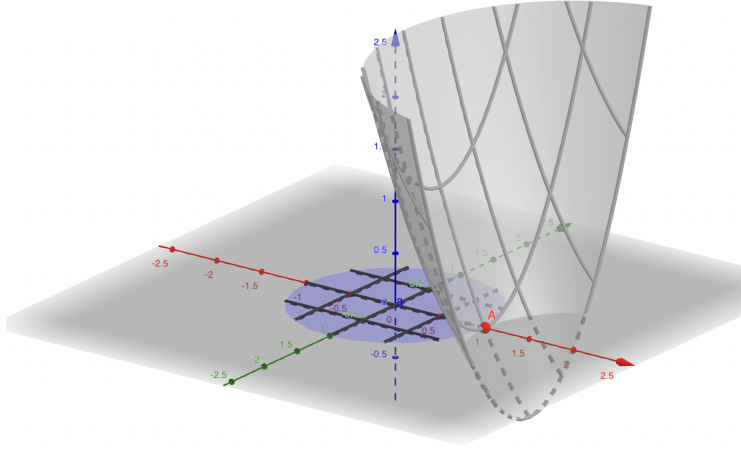
this is a contradiction.

So stationary condition is that all partial derivatives are the same.

**Example 5.**  $C = B(0, 1)$  the unit ball. Stationary condition is equivalent to

$$\nabla f(\mathbf{x}^*) = 0 \text{ or } \|\mathbf{x}^*\| = 1 \text{ and } \exists \lambda \leq 0 \text{ s.t. } \nabla f(\mathbf{x}^*) = \lambda \mathbf{x}^*$$

this means if the stationary point is in the interior of the unit ball, the gradient must be 0. But if it is on the boundary, the gradient should be 0 or point towards the centre of the unit ball. For example, consider optimisation of  $f(\mathbf{x}) = (x_1 - 2)^2 + x_2^2 - 1$  over the unit ball, the minimum is at  $(1, 0)$  (see Figure 3) with gradient  $(-2, 0) = -2*(1, 0)$



**Figure 3: Optimisation of  $(x_1 - 2)^2 + x_2^2 - 1$  over unit ball**

### 5.3 Orthogonal Projection

Projection of  $\mathbf{x}$  to convex set, i.e. the point in convex set  $C$  with closest distance to  $\mathbf{x}$ , is

$$P_C(\mathbf{x}) := \text{Argmin}\{\|\mathbf{y} - \mathbf{x}\|^2 : \mathbf{y} \in C\} = \text{Argmin}\{\|\mathbf{y} - \mathbf{x}\| : \mathbf{y} \in C\}$$

and note this is a minimisation problem of convex quadratic function  $\|\mathbf{y} - \mathbf{x}\|^2$  over convex set  $C$ . This is well-defined (first projection theorem), i.e. optimal point exists and is unique.

*Proof.*

$$f(\mathbf{y}) = \|\mathbf{y} - \mathbf{x}\|^2 = \mathbf{y}^T I \mathbf{y} - 2\mathbf{x}^T \mathbf{y} + \|\mathbf{x}\|^2$$

it is quadratic which is coercive (because  $I$  is positive definite matrix), so at least one optimal solution exists. Further,  $f$  is strictly convex, so there is at most one optimal solution by Theorem 5.1. That means  $P_c(\mathbf{x})$  is uniquely defined.  $\square$



Distance from  $\mathbf{x}$  to convex set  $C$  is

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

Following are examples of orthogonal projection

**Example 6.** Let  $C = \mathbb{R}_+^n := \{\mathbf{x} \in \mathbb{R}^n : x_i \geq 0\}$ .  $P_C(\mathbf{x})$  is an optimisation problem where the object function  $\sum_{i=1}^n (y_i - x_i)^2$  and constraints  $\forall i, y_i \geq 0$  are separable. So consider the individual optimisation problems

$$\min\{(y_i - x_i)^2 : y_i \geq 0\}$$

The solution is

$$y_i^* = [x_i]_+ := \begin{cases} x_i & x_i \geq 0 \\ 0 & x_i < 0 \end{cases}$$

So  $P_{\mathbb{R}_+^n}(\mathbf{x}) = [\mathbf{x}]_+$  where  $[\mathbf{x}]_+ := ([x_1]_+, [x_2]_+, \dots, [x_n]_+)^T$ .

**Example 7** (Box projection). Let  $C = [l_1, u_1] \times [l_2, u_2] \times \dots \times [l_n, u_n]$  where  $l_i \leq u_i$ . Note the optimisation problem of  $\text{Argmin}_{\mathbf{y}}\{\|\mathbf{y} - \mathbf{x}\| : \mathbf{y} \in C\}$  is actually separable as long as constraint  $\mathbf{y} \in C$ . So in this case, like the previous example, consider the individual optimisation problems

$$\min\{(y_i - x_i)^2 : l_i \leq y_i \leq u_i\}$$

the optimal solution is

$$y_i = \begin{cases} u_i & x_i \geq u_i \\ x_i & l_i < x_i < u_i \\ l_i & x_i \leq l_i \end{cases}$$

**Example 8** (Projection to balls). Let  $C = B(0, r)$ . This time the optimisation problem  $\text{Argmin}_{\mathbf{y}}\{\|\mathbf{y} - \mathbf{x}\| : \mathbf{y} \in C\}$  is no longer separable. When  $\|\mathbf{x}\| \leq r$ ,  $\mathbf{y} = \mathbf{x}$  is optimal solution. Otherwise,  $\mathbf{y}$  must be on the boundary of the ball. Assume the contrary, it is a stationary point in the interior of the unit ball, i.e.  $2(\mathbf{y} - \mathbf{x}) = 0 \Rightarrow \mathbf{y} = \mathbf{x}$ . But  $\mathbf{x} \notin C$ , contradiction.

So the minimisation problem is equivalent to

$$\text{Argmin}_{\mathbf{y}}\{-2\mathbf{x}^T \mathbf{y} + r^2 + \|\mathbf{x}\|^2 : \|\mathbf{y}\| = r\} = \text{Argmin}_{\mathbf{y}}\{-2\mathbf{x}^T \mathbf{y} : \|\mathbf{y}\| = r\}$$

where the norm is squared, expanded using the inner product, and constant terms are removed. Using Cauchy-Schwarz,

$$-2\mathbf{x}^T \mathbf{y} \geq -2\|\mathbf{x}\|\|\mathbf{y}\| = -2r\|\mathbf{x}\|$$

and this bound is attained at  $\mathbf{y} = r\mathbf{x}/\|\mathbf{x}\|$ . So

$$P_{B(0,r)} = \begin{cases} \mathbf{x} & \|\mathbf{x}\| \leq r, \\ r \frac{\mathbf{x}}{\|\mathbf{x}\|} & \|\mathbf{x}\| > r \end{cases}$$

**Theorem 5.2** (Second Projection Theorem). *Given convex set  $C$ , any point  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{y} \in C$ , the angle between  $\mathbf{x} - P_C(\mathbf{x})$ ,  $\mathbf{y} - P_C(\mathbf{x})$  is greater than or equal to 90 degrees. (Try to draw it to see why this is true!) i.e.*

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

*Proof.* Use the fact that  $P_C(\mathbf{x})$  is solution to optimisation

$$\begin{aligned} \min \quad & g(\mathbf{y}) := \|\mathbf{y} - \mathbf{x}\|^2 \\ \text{s.t.} \quad & \mathbf{y} \in C \end{aligned}$$

which means  $P_C(\mathbf{x})$  is a stationary point, i.e.

$$\nabla g(P_C(\mathbf{x}))(\mathbf{y} - P_C(\mathbf{x})) \geq 0 \quad \forall \mathbf{y} \in C$$

which is equivalent to the required inequality. □

Orthogonal projection gives an equivalent condition to stationarity

**Theorem 5.3** (Representation of Stationarity). *For any continuously differentiable  $f$ , closed convex set  $C$ ,  $s > 0$ ,  $\mathbf{x}^*$  is stationary iff*

$$\mathbf{x}^* = P_C(\mathbf{x}^* - s\nabla f(\mathbf{x}^*))$$

*note  $\mathbf{x}^* - s\nabla f(\mathbf{x}^*)$  is a step in gradient descent with step size  $s$ . So the equation means the projection of the point after gradient descent back to the convex set is  $\mathbf{x}^*$  itself.*

*Proof.* Use second projection theorem on the pair  $\mathbf{x}^* - s\nabla f(\mathbf{x}^*)$ ,  $\mathbf{x}$  where  $\mathbf{x} \in C$ . The second theorem is useful for proving properties with inequality.  $\square$

## 5.4 Gradient Projection

Theorem 5.3 gives an iterative formulae that numerically finds minimal value of  $f$  over  $C$ , i.e.

$$x_{k+1} = P_C(x_k - t_k \nabla f(x_k))$$

Similar to gradient descent, there are many strategies to choose the step sizes  $t_k$ . But before that, we need to prove gradient projection method has sufficient decrease (find a lower bound for  $f(\mathbf{x}) - f(P_C(\mathbf{x} - t\nabla f(\mathbf{x})))$ )

**Lemma 5.4** (Sufficient Decrease of gradient projection). *Given  $f \in C^{1,1}(C)$  with Lipschitz constant  $L$  where  $C$  is closed, convex. For any  $\mathbf{x} \in C$ ,  $t \in (0, 2/L)$ ,*

$$f(\mathbf{x}) - f(P_C(\mathbf{x} - t\nabla f(\mathbf{x}))) \geq t \left(1 - \frac{Lt}{2}\right) \left\| \frac{1}{t}(\mathbf{x} - P_C(\mathbf{x} - t\nabla f(\mathbf{x}))) \right\|^2$$

*This is very similar to the sufficient decrease lemma for gradient descent (lemma 3.10)*

*Proof.* Let  $\mathbf{x}^+ := P_C(\mathbf{x} - t\nabla f(\mathbf{x}))$ , using descent lemma (lemma 3.9),

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

using the second projection theorem,

$$(\mathbf{x} - t\nabla f(\mathbf{x}) - \mathbf{x}^+)^T(\mathbf{x} - \mathbf{x}^+) \leq 0$$

so

$$\nabla f(\mathbf{x})^T(\mathbf{x}^+ - \mathbf{x}) \leq -\frac{1}{t} \|\mathbf{x}^+ - \mathbf{x}\|^2$$

plugging this to the first inequality,

$$\begin{aligned} f(\mathbf{x}^+) &\leq f(\mathbf{x}) + \left(\frac{L}{2} - \frac{1}{t}\right) \|\mathbf{x}^+ - \mathbf{x}\|^2 \\ \Rightarrow f(\mathbf{x}) - f(\mathbf{x}^+) &\geq \frac{1}{t} \left(1 - \frac{Lt}{2}\right) \|\mathbf{x}^+ - \mathbf{x}\|^2 \end{aligned}$$

which is the same as required inequality.  $\square$

So an ideal choice for constant step size is  $1/L$  where  $L$  is the Lipschitz constant of objective  $f$ . To simplify notations, let

$$G_L(\mathbf{x}) := L \left[ \mathbf{x} - P_C \left( \mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x}) \right) \right] \quad L > 0$$

note if  $C = \mathbb{R}^n$ ,  $G_L(\mathbf{x}) = \nabla f(\mathbf{x})$ . And by representation theorem of stationarity,  $G_M(\mathbf{x}) = 0 \Leftrightarrow \mathbf{x}$  is a stationary point. The sufficient decrease now can be restated as

$$f(\mathbf{x}) - f(P_C(\mathbf{x} - t\nabla f(\mathbf{x}))) \geq t \left(1 - \frac{Lt}{2}\right) \|G_{1/t}(\mathbf{x})\|^2$$

This new general version of the gradient function has the following property

**Lemma 5.5** (Monotonicity). *If  $f$  is continuously differentiable function defined on closed and convex set  $C$ , and  $L_1 \geq L_2$ , then*

$$\begin{aligned} \|G_{L_1}(\mathbf{x})\| &\geq \|G_{L_2}(\mathbf{x})\| \\ \frac{\|G_{L_1}(\mathbf{x})\|}{L_1} &\leq \frac{\|G_{L_2}(\mathbf{x})\|}{L_2} \end{aligned}$$

### 5.4.1 Backtracking

Similar to gradient descent, we need to choose initial step size  $s > 0$ , tracking rates  $\alpha, \beta \in (0, 1)$ . At iteration  $k$  of gradient projection, let  $t_k := s$ . Continue update by  $t_k = \beta t_k$  if

$$f(\mathbf{x}_k) - f(P_C(\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k))) < \alpha t_k \|G_{1/t_k}(\mathbf{x}_k)\|^2$$

process stops when

$$f(\mathbf{x}_k) - f(P_C(\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k))) \geq \alpha t_k \|G_{1/t_k}(\mathbf{x}_k)\|^2 \quad (*)$$

again this is very similar to the formulae for gradient descent.

If  $f \in C^{1,1}(C)$ , backtracking eventually stops. Using sufficient decrease lemma for gradient projection,

$$f(\mathbf{x}_k) - f(P_C(\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k))) \geq \left(1 - \frac{L t_k}{2}\right) t_k \|G_{1/t_k}(\mathbf{x}_k)\|^2$$

when  $1 - L t_k / 2 \geq \alpha$ , the stopping inequality (\*) is satisfied. This is equivalent to  $t_k \leq 2(1 - \alpha)/L$ . So after  $t_k$  is shrunk below this boundary, backtracking stops.

**Lower bound on  $t$ :**  $t_k$  is either  $s$ , or backtracking invoked. If backtracking is invoked, the step size in the last iteration of backtracking  $t_k/\beta > 2(1 - \alpha)/L$ , i.e.  $t_k > 2(1 - \alpha)\beta/L$ , so

$$t_k \geq \min \left\{ s, \frac{2(1 - \alpha)\beta}{L} \right\}$$

**Lemma 5.6** (Sufficient decrease of backtracking). *If  $(\mathbf{x}_k)_{k \geq 0}$  is the sequence of  $\mathbf{x}$  generated by gradient projection,*

$$f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \geq M \|G_d(\mathbf{x}_k)\|^2$$

where

$$M = \alpha \min \left\{ s, \frac{2(1 - \alpha)\beta}{L} \right\}, \quad d = \frac{1}{s}$$

*Proof.* By definition of backtracking,

$$f(\mathbf{x}_k) - f(P_C(\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k))) \geq \alpha t_k \|G_{1/t_k}(\mathbf{x}_k)\|^2$$

so the value of  $M$  follows from lower bound of  $t_k$  derived above. Note  $t_k \leq s$ , so by lemma 5.5,  $\|G_{1/t_k}(\mathbf{x}_k)\| \geq \|G_{1/s}(\mathbf{x}_k)\|$  which yields the result above.  $\square$

Finally, for convergence of the gradient projection method, we need to verify  $f(\mathbf{x}_k)$  converges.

**Theorem 5.7** (Convergence of gradient projection). *Assume  $f \in C^{1,1}(C)$  is bounded below,  $\mathbf{x}_k$  is sequence generated by gradient projection, then*

- (1)  $f(\mathbf{x}_k)$  converges
- (2)  $G_d(\mathbf{x}_k) \rightarrow 0$ , i.e. the lower bound in sufficient decrease lemma does not block the convergence of  $\mathbf{x}_k$ .

*Proof.*

$$f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \geq M \|G_d(\mathbf{x}_k)\|^2 \quad (*)$$

for some  $M > 0$  by the previous lemma and sufficient decrease lemma. So directly we have  $f(\mathbf{x}_k) \geq f(\mathbf{x}_{k+1})$  with equality holds iff  $\mathbf{x}_k$  is stationary point.

$f(\mathbf{x}_k)$  must converge as it is bounded below and non-increasing. So  $f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \rightarrow 0$ , combined with the above inequality (\*) gives  $\|G_d(\mathbf{x}_k)\| \rightarrow 0$ .  $\square$

## 5.5 Sparsity constraint

Sparsity in data means that most entries are 0 while only few important entries are kept. Given continuously differentiable  $f$  which is bounded below, sparsity of optimal solution  $\mathbf{x}$  is realised by the following NON-convex optimisation problem (called *Sparsity constrained problem*)

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & \|\mathbf{x}\|_0 \leq s \quad s \in \mathbb{N}(s \neq 0) \end{aligned}$$

where the 0-norm (not actually a norm) counts number of non-zero components

$$\|\mathbf{x}\|_0 := |\{i : x_i \neq 0\}|$$

Since this problem is non-convex, theorems derived in previous sections do not apply.

Define the set  $C_s := \{\mathbf{x} : \|\mathbf{x}\|_0 \leq s\}$ .

Motivated by Theorem 5.3, we say  $\mathbf{x}^*$  is  $L$ -stationary point if

$$\mathbf{x}^* \in P_{C_s} \left( \mathbf{x}^* - \frac{1}{L} \nabla f(\mathbf{x}^*) \right)$$

Equality is not used, because for non-convex set  $C_s$ , projection to  $C_s$  is not well-defined, i.e. it could give several points. For example,

$$P_{C_2}((3, 1, 1)) = \{(3, 1, 0), (3, 0, 1)\}$$

this projection sends the smallest  $n - s$  entries to 0.

This stationarity condition depends on  $L$ , smaller  $L$  makes it stricter. A natural question that arises is for what  $L$ ,  $L$ -stationarity is related to optimality.

Given  $f \in C_{L_f}^{1,1}(\mathbb{R}^n)$ , assume  $L > L_f$ , the  $L$ -stationarity condition

$$\mathbf{x}^* \in P_{C_s} \left( \mathbf{x}^* - \frac{1}{L} \nabla f(\mathbf{x}^*) \right)$$

is equivalent to

$$\mathbf{x}^* \in \text{Argmin}_{\mathbf{y} \in C_s} \left\| \mathbf{y} - \left( \mathbf{x}^* - \frac{1}{L} \nabla f(\mathbf{x}^*) \right) \right\|^2$$

Recall in descent lemma for Lipschitz continuous function,

$$f(\mathbf{y}) \leq f(\mathbf{x}^*) + \nabla f(\mathbf{x}^*)^T (\mathbf{y} - \mathbf{x}^*) + \frac{L}{2} \|\mathbf{x}^* - \mathbf{y}\|^2 =: h_L(\mathbf{y}, \mathbf{x}^*)$$

rearranging  $h_L(\mathbf{y}, \mathbf{x}^*)$  yields

$$= \frac{L}{2} \left\| \mathbf{y} - \left( \mathbf{x}^* - \frac{1}{L} \nabla f(\mathbf{x}^*) \right) \right\|^2 + f(\mathbf{x}^*) - \frac{1}{2L} \|\nabla f(\mathbf{x}^*)\|^2$$

where the last two terms are constant w.r.t.  $\mathbf{y}$ . So the optimisation problem is equivalent to

$$\mathbf{x}^* \in \text{Argmin}_{\mathbf{y} \in C_s} h_L(\mathbf{y}, \mathbf{x}^*)$$

It can be proved that if  $L > L_f$  (the Lipschitz constant of  $f$ ), then  $L$ -stationarity is the necessary condition for  $\mathbf{x}^*$  to be optimal solution of the sparsity constrained problem.

A generalisation of gradient projection algorithm known as *iterative hard-threshold* (IHT) method, begins with constant  $L > L_f$  and  $\mathbf{x}_0 \in C_s$ . Then the iterative step is given by

$$\mathbf{x}^{k+1} \in P_{C_s} \left( \mathbf{x}^k - \frac{1}{L} \nabla f(\mathbf{x}^k) \right)$$

and this is equivalent to

$$\mathbf{x}^{k+1} \in \text{Argmin}_{\mathbf{x} \in C_s} h_L(\mathbf{x}, \mathbf{x}^k)$$

as derived above. It has been proved in (Sun and Cheng, 2016) that IHT converges.

## 6 KKT conditions

Optimality condition we studied:

$$\nabla f(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) \geq 0 \text{ for all } \mathbf{x} \in C$$

may not be easy to verify. But alternative formulation KKT can be used, which is easier to solve. We will study a specific family of problems: linearly constrained problems (LCP) first.

LCP takes the form

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

with  $\mathbf{a}_i \in \mathbb{R}^n$ ,  $b_i \in \mathbb{R}$ . We always assume  $f$  is continuously differentiable.

The KKT condition is like a generalisation of Lagrange multiplier,

**Theorem 6.1** (KKT condition for LCP(necessary condition)). *If  $\mathbf{x}^*$  is local minimum point of (LCP), exists  $\lambda_1, \dots, \lambda_m \geq 0$  s.t.*

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \mathbf{a}_i = 0 \quad (KKT1)$$

and

$$\lambda_i (\mathbf{a}_i^T \mathbf{x}^* - b_i) = 0 \quad (KKT2)$$

*Remark.* (KKT1) is a modification to  $\nabla f(\mathbf{x}^*) = 0$  by adding terms  $\lambda_i \nabla c_i(\mathbf{x}^*)$  where  $c_i(\mathbf{x}) = \mathbf{a}_i^T \mathbf{x} - b_i$  are constraints,  $\lambda_i \geq 0$  ensures that the direction of inequality  $c_i(\mathbf{x}) \leq 0$  is not changed. Later you will see that multipliers of constraints like  $h_i(\mathbf{x}) = 0$  do not have to be non-negative. Set of  $m$  equations in (KKT2), also called *complementary slackness conditions*, are satisfied iff  $\lambda_i = 0$  or  $\mathbf{a}_i^T \mathbf{x}^* = b_i$ . Recall the examples given in the last chapter for optimality conditions on unit-sum set or unit ball, either some(or all) partial derivative(s) to be 0, or  $\mathbf{x}^*$  sits at the boundary. In this case, the boundary is  $\{\mathbf{x} : \mathbf{a}_i^T \mathbf{x} = b_i\}$ , and we will call the set  $\{i : \mathbf{a}_i^T \mathbf{x} = b_i\} =: I(\mathbf{x})$  the set of active constraints.

If given  $f$  is convex, then the KKT condition becomes sufficient. In practice, even if  $f$  is not convex, you can solve the equations given by KKT conditions, find all possible candidates of  $\mathbf{x}^*$ , and pick the one giving smallest  $f(\mathbf{x})$  value.

KKT condition for LCP is built on Farkas' lemma, which requires the separation theorem to prove. Separation theorem simply means there is always a hyperplane separating a point from a closed, convex set.

**Theorem 6.2.** *Given non-empty, closed, convex  $C$  and  $\mathbf{y} \notin C$ , exists hyperplane*

$$\{\mathbf{x} : \mathbf{p}^T \mathbf{x} = \alpha\} \quad \mathbf{p} \in \mathbb{R}^n \setminus \{\mathbf{0}\}, \text{ and } \alpha \in \mathbb{R}$$

*that separates  $\mathbf{y}$  and  $C$ , i.e.*

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

*Proof.* And let  $\mathbf{p}$  be  $\mathbf{y} - P_C(\mathbf{y})$ , the vector connecting  $\mathbf{y}$  to  $C$  which is orthogonal to  $C$ , and  $\alpha := (\mathbf{y} - P_C(\mathbf{y}))^T P_C(\mathbf{y})$ . The choice of  $\alpha$  ensures  $\mathbf{p}^T P_C(\mathbf{y}) = \alpha$ . Then second orthogonal projection theorem ensures  $\mathbf{p}^T \mathbf{x} \leq \alpha$  for all other  $\mathbf{x} \in C$ . Closeness of  $C$  is to ensure  $P_C(\mathbf{y}) \in C$ . See detailed proof in lecture notes(page 94).  $\square$

Farkas' lemma is stated as below

**Lemma 6.3.** *Given  $A \in \mathbb{R}^{m \times n}$ ,  $\mathbf{c} \in \mathbb{R}^n$ , exactly one of the following hold*

- I. *exists  $\mathbf{y} \geq 0$  s.t.  $A^T \mathbf{y} = \mathbf{c}$*
- II. *there exists  $\mathbf{x}$  s.t.  $A\mathbf{x} \leq \mathbf{0}$ ,  $\mathbf{c}^T \mathbf{x} > 0$ .*

It means that given convex cone  $C_0 := \{A^T \mathbf{y} : \mathbf{y} \in \mathbb{R}_+^n\} = \{\sum y_i \mathbf{r}_i^T\}$  where  $\mathbf{r}_i$  are rows of  $A$  (i.e. set of linear combinations of rows of  $A$  with non-negative coefficients). For any vector  $\mathbf{c}$ , if it is not in  $C_0$  (case I not satisfied), then you can find hyperplane  $\{\mathbf{y} : \mathbf{x}^T \mathbf{y} = 0\}$  s.t. separates cone  $C_0$  and  $\mathbf{c}$  (condition II). To see this,

$$A\mathbf{x} \leq \mathbf{0} \Leftrightarrow \begin{pmatrix} \mathbf{r}_1^T \mathbf{x} \\ \vdots \\ \mathbf{r}_m^T \mathbf{x} \end{pmatrix} \leq \mathbf{0} \Leftrightarrow \mathbf{r}_i^T \mathbf{x} < 0 \Rightarrow \mathbf{y}^T \mathbf{x} < 0 \forall \mathbf{y} \in C_0$$

Farkas' lemma is formulated in the way "either A or B", which can be treated as B iff not A.

**Lemma 6.4** (Farkas' Lemma). *Given  $A \in \mathbb{R}^{m \times n}$ ,  $\mathbf{c} \in \mathbb{R}^n$ , the following are equivalent:*

- (A)  $A\mathbf{x} \leq \mathbf{0} \Rightarrow \mathbf{c}^T \mathbf{x} \leq 0$
- (B) *exists  $\mathbf{y} \in \mathbb{R}_+^m$  s.t.  $A^T \mathbf{y} = \mathbf{c}$*

*Proof.* (B)  $\Rightarrow$  (A) is easy because you are using an inequality  $A\mathbf{x} \leq \mathbf{0}$  and an equality  $A^T \mathbf{y} = \mathbf{c}$  to derive an inequality, a simple substitution will deal with this case. (see detailed proof on page 95 of lecture notes)

The other direction is easier with contradiction. Assume (A) is true, (B) is false, aim to find  $\mathbf{x}$  s.t.  $(A\mathbf{x})^T = \mathbf{x}^T A^T \leq 0$  but  $\mathbf{c}^T \mathbf{x} = \mathbf{x}^T \mathbf{c} > 0$ . This is exactly a separation problem using hyperplane  $\{\mathbf{y} : \mathbf{x}^T \mathbf{y} = 0\}$  where  $\mathbf{x}$  is the vector to be found. So we need to justify the cone

$$C_0 := \{A^T \mathbf{y} : \mathbf{y} \in \mathbb{R}_+^n\}$$

is convex, closed, non-empty. It is closed because it is a conic hull of a finite set of vectors  $\{\mathbf{r}_i\}_{i=1, \dots, m}$  (search for proof this result if you are interested).  $C_0$  is clearly convex, non-empty, so separation theorem applies.  $\square$

A different theorem is required for KKT condition of general optimisation problem,

**Theorem 6.5** (Gordan's alternative theorem). *Given  $A \in \mathbb{R}^{m \times n}$ , exactly one of the following hold:*

- (A)  $\exists \mathbf{x}$  s.t.  $A\mathbf{x} < \mathbf{0}$
- (B)  $\exists \mathbf{p} \geq \mathbf{0} (\mathbf{p} \neq \mathbf{0})$  s.t.  $A^T \mathbf{p} = \mathbf{0}$

*Proof.* If (A) holds, and for contradiction assume (B) also holds, then

$$\mathbf{x}^T A^T \mathbf{p} = 0 \Rightarrow (A\mathbf{x})^T \mathbf{p} = 0$$

but  $A\mathbf{x} < \mathbf{0}, \mathbf{p} \geq \mathbf{0}$  and  $\mathbf{p} \neq \mathbf{0}$ , so this is impossible. (if you are not convinced, write it element-wise)

If (A) fails, note (A) is equivalent to

$$A\mathbf{x} + s\mathbf{e} \leq \mathbf{0}, \quad s > 0$$

for small enough  $s > 0$ , and  $\mathbf{e} = (1, 1, \dots, 1)^T$ . Defining  $\mathbf{c} := e_{n+1}$  Writing these in matrix form:

$$\underbrace{\begin{pmatrix} A & \mathbf{e} \end{pmatrix}}_{=: \tilde{A}} \underbrace{\begin{pmatrix} \mathbf{x} \\ s \end{pmatrix}}_{=: \mathbf{w}} \leq \mathbf{0}, \quad \mathbf{c}^T \begin{pmatrix} \mathbf{x} \\ s \end{pmatrix} > 0$$

So if (A) does not have solution,

$$\tilde{A}\mathbf{w} \leq \mathbf{0}, \quad \mathbf{c}^T \mathbf{w} > 0$$

has no solution for  $\mathbf{w} \in \mathbb{R}^{n+1}$ , so by Farkas' lemma, exists  $\mathbf{z} \in \mathbb{R}_+^{n+1}$  s.t.

$$\begin{aligned} \tilde{A}^T \mathbf{z} &= \begin{pmatrix} A^T \\ \mathbf{e}^T \end{pmatrix} \mathbf{z} = \mathbf{c} \\ \Rightarrow A^T \mathbf{z} &= \mathbf{0}, \quad \mathbf{e}^T \mathbf{z} = 1 \end{aligned}$$

clearly  $\mathbf{z} \neq \mathbf{0}$ , so  $\mathbf{z}$  satisfies condition (B).  $\square$

Now we are ready to prove the KKT condition, Theorem 6.1.

*Proof.* If  $\mathbf{x}^*$  is local minimum, then  $\mathbf{x}^*$  is stationary point. Let  $\mathbf{y} := \mathbf{x} - \mathbf{x}^*$ , stationarity condition becomes  $\nabla f(\mathbf{x}^*)^T \mathbf{y} \geq 0$  for any  $\mathbf{y}$  satisfying the constraint. (i.e.  $\mathbf{a}_i^T(\mathbf{y} + \mathbf{x}^*) \leq b_i$ ) Recall we defined the set of active constraints  $I(\mathbf{x})$  (i.e. indices of constraints  $i$  s.t.  $\mathbf{x}$  falls on the boundary) by  $I(\mathbf{x}) := \{i : \mathbf{a}_i^T \mathbf{x} = b_i\}$ , using these we can partition the inequalities into two categories. Note the constraints are equivalent to  $\mathbf{a}_i^T \mathbf{y} \leq b_i - \mathbf{a}_i^T \mathbf{x}^*$ , so

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

constraints not in  $I(\mathbf{x}^*)$  can be removed: assume  $\mathbf{a}_i^T \mathbf{y} \leq 0$  for  $i \in I(\mathbf{x}^*)$ . For  $i \notin I(\mathbf{x}^*)$ ,  $b_i - \mathbf{a}_i^T \mathbf{x}^* > 0$ , so exists small enough  $\alpha > 0$  s.t.  $\mathbf{a}_i^T(\alpha \mathbf{y}) \leq b_i - \mathbf{a}_i^T \mathbf{x}^*$ . Together with the fact that  $\mathbf{a}_i^T(\alpha \mathbf{y}) \leq 0$  for  $i \in I(\mathbf{x}^*)$  (because  $\mathbf{a}_i^T \mathbf{y} \leq 0$ ), stationary condition for  $\nabla f(\mathbf{x}^*)^T(\alpha \mathbf{y}) \geq 0$  is satisfied. That yields  $\nabla f(\mathbf{x}^*)^T \mathbf{y} \geq 0$ . So indeed the condition on  $i \notin I(\mathbf{x}^*)$  is not required, i.e.

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

You may recognise this as condition (A) of Farkas' lemma (Lemma 6.4, correspondence:  $A = (\mathbf{a}_i^T)_{i \in I(\mathbf{x}^*)}$ ,  $\mathbf{c} = -\nabla f(\mathbf{x}^*)$  where negative sign ensures inequality is in the direction  $\leq 0$ ). So we can find  $\lambda_i \geq 0$  (corresponds to  $\mathbf{y}_i$  in condition (B) of Farkas' lemma) s.t.

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

For  $i \notin I(\mathbf{x}^*)$ , let  $\lambda_i = 0$ . Now we have  $\lambda_i(\mathbf{a}_i^T \mathbf{x}^* - b_i) = 0$  for all  $i$  and

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \mathbf{a}_i = 0$$

**Further Assumption** If  $f$  is convex, we want to prove  $\mathbf{x}^*$  is optimal solution iff  $\mathbf{x}^*$  satisfies KKT condition. Only have to prove sufficiency (necessity is proved above, because  $\mathbf{x}^*$  is optimal solution means it is local minimum point) Suppose  $\mathbf{x}^*$  satisfies KKT conditions, define

$$L(\mathbf{x}) := f(\mathbf{x}) + \sum_{i=1}^m \lambda_i(\mathbf{a}_i^T \mathbf{x} - b_i)$$

(KKT1) means  $\nabla L(\mathbf{x}) = \mathbf{0}$ .  $L$  is convex, so gradient being 0 is sufficient for  $\mathbf{x}^*$  to be global minimiser of  $L$ . Combine with (KKT2), for any  $\mathbf{x}$  satisfying the constraints  $\mathbf{a}_i^T \mathbf{x} \leq b_i$ ,

$$f(\mathbf{x}^*) + 0 = f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i(\mathbf{a}_i^T \mathbf{x}^* - b_i) = L(\mathbf{x}^*) \leq L(\mathbf{x}) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i(\mathbf{a}_i^T \mathbf{x} - b_i) \leq f(\mathbf{x})$$

where last inequality follows by constraints. So  $\mathbf{x}^*$  is also global optimal point for minimising  $f$  with constraints  $\mathbf{a}_i^T \mathbf{x} \leq b_i$ .  $\square$

Further constraints  $\mathbf{a}^T \mathbf{x} = b$  can be added to LCP, and observe that it is equivalent to  $\mathbf{a}^T \mathbf{x} \leq b$ ,  $-\mathbf{a}^T \mathbf{x} \leq -b$ . So KKT conditions can also be imposed.

**Theorem 6.6** (KKT condition for LCP with equalities). *Given the following optimisation problem*

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

where  $f$  is continuously differentiable,  $\mathbf{a}_i, \mathbf{c}_j \in \mathbb{R}^n$ ,  $b_i, d_j \in \mathbb{R}$ .

(a) If  $\mathbf{x}^*$  is local minimum point, exists  $\lambda_i \geq 0$ ,  $\mu_j \in \mathbb{R}$  s.t.

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \mathbf{a}_i + \sum_{j=1}^p \mu_j \mathbf{c}_j = 0 \quad (KKT1)$$

and

$$\lambda_i(\mathbf{a}_i^T \mathbf{x}^* - b_i) = 0 \quad (\text{KKT2})$$

note actually there are other  $p$  hidden equations: the constraints!

(b) If  $f$  is convex and  $\mathbf{x}^*$  is feasible solution (i.e. satisfies the constraints), if there are  $\lambda_i \geq 0$ ,  $\mu_j \in \mathbb{R}$  satisfying (KKT1), (KKT2), then  $\mathbf{x}^*$  is optimal solution.

*Proof.* The proof basically uses previous KKT theorem, in the spirit of

$$\mathbf{a}^T \mathbf{x} \leq b, -\mathbf{a}^T \mathbf{x} \leq -b \Leftrightarrow \mathbf{a}^T \mathbf{x} = b$$

though a small trick is required: let  $\mu^+ := \max\{\mu, 0\} \geq 0$ ,  $\mu^- := \max\{-\mu, 0\} \geq 0$ , then  $\mu = \mu^+ - \mu^-$ .  
Proof left as exercise.  $\square$

By treating orthogonal projections as optimisation problems, one can find the formulae for projections to affine space  $C = \{\mathbf{x} : A\mathbf{x} = \mathbf{b}\}$  (assume rows of  $A$  are independent), hyperplane  $H := \{\mathbf{x} : \mathbf{a}^T \mathbf{x} = b\}$  and half-space  $H^- := \{\mathbf{x} : \mathbf{a}^T \mathbf{x} \leq b\}$  are

$$P_C(\mathbf{y}) = \mathbf{y} - A^T(AA^T)^{-1}(A\mathbf{y} - \mathbf{b})$$

$$P_H(\mathbf{y}) = \mathbf{y} - \frac{\mathbf{a}^T \mathbf{y} - b}{\|\mathbf{a}\|^2} \mathbf{a}$$

$$P_{H^-}(\mathbf{y}) = \mathbf{y} - \frac{[\mathbf{a}^T \mathbf{y} - b]_+}{\|\mathbf{a}\|^2} \mathbf{a}$$

note  $AA^T$  is guaranteed to be invertible because  $A$  has full row rank by our assumption.

## 6.1 General KKT Condition

In general, optimisation problem with constraints takes the form

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

where  $f, g_i, h_j$  are continuously differentiable.

In this general case, (KKT1) is often written as  $\nabla_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = 0$  where

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) := f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) + \sum_{j=1}^p \mu_j h_j(\mathbf{x})$$

For the LCP case, with the following definitions

$$A := \begin{pmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_m^T \end{pmatrix}, \quad C := \begin{pmatrix} \mathbf{c}_1^T \\ \vdots \\ \mathbf{c}_p^T \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix}, \quad \mathbf{d} = \begin{pmatrix} d_1 \\ \vdots \\ d_p \end{pmatrix}$$

constraints become  $A\mathbf{x} \leq \mathbf{b}$ ,  $C\mathbf{x} = \mathbf{d}$  and Lagrangian becomes

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

with KKT condition

$$\nabla f(\mathbf{x}) + A^T \boldsymbol{\lambda} + C^T \boldsymbol{\mu} = 0$$

A definition makes life easier when dealing with general cases:



**Definition 5.** When considering problem

$$\begin{aligned} \min \quad & h(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{x} \in C \end{aligned}$$

where  $h$  is continuously differentiable defined over closed and convex  $C$ .  $\mathbf{d} \neq \mathbf{0}$  is called feasible descent direction at  $\mathbf{x} \in C$  if  $\nabla h(\mathbf{x})^T \mathbf{d} < 0$  and exists  $\epsilon > 0$  s.t.  $\mathbf{x} + t\mathbf{d} \in C$  for all  $t \in [0, \epsilon]$ .

It can be shown that  $\mathbf{x}^*$  is local optimal solution iff no feasible descent direction at  $\mathbf{x}^*$  exists. (either you cannot descend anymore, or descend direction is outside the constraints)

Now we begin to state theorems for KKT conditions in non-linear case, but will not prove them.

The original KKT condition for problem

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

(where  $f, g_i$  are continuously differentiable) is given by Fritz-John, proved using Gordan's theorem. Fritz-John condition: exists  $\lambda_0, \dots, \lambda_m \geq 0$  which are not all zero s.t.

$$\lambda_0 \nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \nabla g_i(\mathbf{x}^*) = \mathbf{0} \quad (\text{FJ1})$$

$$\lambda_i g_i(\mathbf{x}^*) = 0 \quad \text{for } i = 1, \dots, m \quad (\text{FJ2})$$

The problem is if  $\lambda_0 = 0$ , (FJ1) simply becomes

$$\sum_{i=1}^m \lambda_i \nabla g_i(\mathbf{x}^*) = \mathbf{0}$$

which means  $\{\nabla g_i(\mathbf{x}^*)\}_{i \in I(\mathbf{x}^*)}$  (similar to before, we define  $I(\mathbf{x}) := \{i : g_i(\mathbf{x}) = 0\}$ ) is linearly dependent set. Adding the constraint that this set is linearly independent, and setting  $\lambda_0 = 1$  gives us a useful kKT condition

**Theorem 6.7** (KKT for inequality constraints). *If  $\mathbf{x}^*$  is local minimum of*

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

where  $f, g_i$  are continuously differentiable. *If  $\{\nabla g_i(\mathbf{x}^*)\}_{i \in I(\mathbf{x}^*)}$  is linearly independent, then exists  $\lambda_1, \dots, \lambda_m \geq 0$  s.t.*

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \nabla g_i(\mathbf{x}^*) = \mathbf{0}$$

$$\lambda_i g_i(\mathbf{x}^*) = 0 \quad \text{for } i = 1, \dots, m$$

Now back to the full problem (P), with equality constraints,

**Theorem 6.8** (KKT full theorem). *If  $\mathbf{x}^*$  is local minimum of (P), where  $f, g_i, h_j$  are continuously differentiable. If*

$$\{\nabla g_i(\mathbf{x}^*)\}_{i \in I(\mathbf{x}^*)} \cup \{\nabla h_j(\mathbf{x}^*)\}_{j=1, \dots, p}$$

*is linearly independent, then exists  $\lambda_1, \dots, \lambda_m \geq 0, \mu_1, \dots, \mu_p \in \mathbb{R}$  s.t.*

$$\nabla L(\mathbf{x}^*, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{0} \quad (\text{KKT1})$$

$$\lambda_i g_i(\mathbf{x}^*) = 0 \quad \text{for } i = 1, \dots, m \quad (\text{KKT2})$$

*Remark.* Convexity of  $f, g_i, h_j$  are not required here.

The additional assumption is defined as regularity

**Definition 6** (Regularity). For optimisation problem (P), feasible point  $\mathbf{x}^*$  is regular if set of gradients of active constraints of inequality constraints and equality constraints,

$$\{\nabla g_i(\mathbf{x}^*)\}_{i \in I(\mathbf{x}^*)} \cup \{\nabla h_j(\mathbf{x}^*)\}_{j=1, \dots, p}$$

is linearly independent.

With this definition, Theorem 6.8 means necessary optimality condition for a regular point to be local minima is KKT condition. If there is no irregular point, every local minima satisfies KKT condition. The regularity condition is not required for LCP. (If gradients of two linear constraints are linearly dependent, then essentially they are the same constraint!)

Regularity can also be used to derive necessary second-order condition, when  $f, g_i, h_j$  are all twice continuously differentiable, if  $\mathbf{x}^*$  is local minimum and regular, then KKT conditions are satisfied, and

$$\mathbf{d}^T \left[ \nabla^2 f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \nabla^2 g_i(\mathbf{x}^*) + \sum_{j=1}^p \mu_j \nabla^2 h_j(\mathbf{x}^*) \right] \mathbf{d} \geq 0 \quad \forall \mathbf{d} \in \Lambda(\mathbf{x}^*)$$

where  $\Lambda(\mathbf{x}^*)$  is the set of vectors perpendicular to gradient of active inequality constraints and all equality constraints

$$\Lambda(\mathbf{x}^*) := \{\mathbf{d} \in \mathbb{R}^n : \nabla g_i(\mathbf{x}^*)^T \mathbf{d} = 0 \quad \forall i \in I(\mathbf{x}^*), \nabla h_j(\mathbf{x}^*)^T \mathbf{d} = 0 \text{ for } j = 1, 2, \dots, p\}$$

## 6.2 KKT for Convex Problems

Things become easier when convexity assumption is added. KKT conditions are always sufficient for  $\mathbf{x}^*$  to be optimal condition as long as  $f, g_i$  are continuously differentiable convex functions, and  $h_j$  are affine functions. Recall these conditions mean feasible set  $C$  is convex. The proof is surprisingly simple:

*Proof.* Note

$$s(\mathbf{x}) =: s(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) + \sum_{j=1}^p \mu_j h_j(\mathbf{x})$$

is convex. By the (KKT1),

$$\nabla s(\mathbf{x}^*) = \nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{0}$$

so  $\mathbf{x}^*$  must be minimiser of  $s$  over  $\mathbb{R}^n$ . So for any  $\mathbf{x} \in \mathbb{R}^n$ ,

$$\begin{aligned} f(\mathbf{x}^*) &= f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}^*) + \sum_{j=1}^p \mu_j h_j(\mathbf{x}^*) \quad \text{derived by (KKT2) and constraints } h_j(\mathbf{x}^*) = 0 \\ &= s(\mathbf{x}^*) \leq s(\mathbf{x}) \quad \text{because } \mathbf{x}^* \text{ is minimiser of } s \\ &= f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) + \sum_{j=1}^p \mu_j h_j(\mathbf{x}) \\ &\leq f(\mathbf{x}) \quad \text{because } \lambda_i \geq 0, g_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) = 0 \end{aligned}$$

so  $\mathbf{x}^*$  is optimal solution. □

Further, with convexity an very convenient *Slater's condition* to replace regularity condition, it only requires existence of one single point in feasible set to be away from the boundary of constraints  $g_i(\mathbf{x}) \leq 0$ , i.e.  $\exists \hat{\mathbf{x}} \in C(\text{feasible set})$  s.t.  $g_i(\hat{\mathbf{x}}) < 0$  for all  $i$ . Note this condition is automatically satisfied if  $g_i$  are affine(linear) maps. So consider the following problem:

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

where  $h_j, s_k$  are affine,  $f, g_i$  are continuously differentiable convex functions, (generalised) Slater's condition is exists  $\hat{\mathbf{x}} \in \mathbb{R}^n$  s.t.

$$\begin{aligned} g_i(\hat{\mathbf{x}}) &< 0, \quad i = 1, 2, \dots, m \\ h_j(\hat{\mathbf{x}}) &\leq 0, \quad i = 1, 2, \dots, p \\ s_k(\hat{\mathbf{x}}) &= 0, \quad i = 1, 2, \dots, q \end{aligned}$$

**Theorem 6.9** (Necessity of KKT for convex problem with constraints). *If problem (CP) satisfies generalised Slater's condition, then exists  $\lambda_1, \dots, \lambda_m, \eta_1, \eta_2, \dots, \eta_p \geq 0, \mu_1, \mu_2, \dots, \mu_q \in \mathbb{R}$  s.t.*

$$\nabla L(\mathbf{x}^*, \boldsymbol{\lambda}, \boldsymbol{\eta}, \boldsymbol{\mu}) = 0$$

and

$$\lambda_i g_i(\mathbf{x}^*) = 0, \quad \eta_j h_j(\mathbf{x}^*) = 0 \quad \forall i, j$$

Below is a summary about relation of KKT conditions with optimality in various conditions on various optimisation problems

$f(\text{objective})$	$g_i(\text{inequality})$	$h_j(\text{equality})$	condition	result
cont. diff.	affine	affine		optimal $\Rightarrow$ KKT
convex	affine	affine		optimal $\Leftrightarrow$ KKT
convex	convex	affine		optimal $\Leftarrow$ KKT
convex	convex	affine	Slater	optimal $\Leftrightarrow$ KKT
cont. diff.	cont. diff.	cont. diff.	regularity	optimal $\Rightarrow$ KKT

Note: cont. diff. means continuously differentiable. When the "condition" is empty, it means there is no extra condition required.

## 7 Duality

Given the following optimisation problem

$$\begin{aligned} (P) \min \quad & x_1^2 + x_2^2 + 2x_1 \\ \text{s.t.} \quad & x_1 + x_2 = 0 \end{aligned}$$

using basic techniques it is easy to prove that  $(-1/2, 1/2)$  is a global minimiser and minimal value  $f^* = -1/2$ . Suppose the constraint is not this easy, we aim to find a lower bound for the minimal value  $f^*$ . One can turn (P) into unconstrained problem by

$$(UP) \min \quad x_1^2 + x_2^2 + 2x_1 + \mu(x_1 + x_2)$$

where  $\mu \in \mathbb{R}$ . Since on the feasible set of (P),  $x_1 + x_2 = 0$ , the optimal value of (UP) is lower bound of  $f^*$ . (UP) is another convex optimisation problem, so finding stationary point yields  $(-1 - \mu/2, -\mu/2)$ . The corresponding optimal value, denote by  $q(\mu)$ , is  $q(\mu) = -\mu^2/2 - \mu - 1$ . Now, we aim to find the best lower bound, i.e. find solution to the following problem

$$(D) \quad \max_{\mu \in \mathbb{R}} \{q(\mu)\}$$

and the value of  $q$  at optimal point of (D) is denoted as  $q^*$ . It is clear that  $q^* \leq f^*$  (called weak duality).

In this particular case, optimal value of (D) is given by  $\mu = -1$ , and  $q^* = -1/2 = f^*$ , so solving the *dual problem* (D) gives the solution to the primal(original) problem (P). This very useful property is called *Strong duality*. (Conditions for strong duality to hold will be studied later).

You may worry that similar techniques cannot be used on inequality constraints, because after adding the term with multiplier, you are not able to remove the constraint without changing the solutions. New primal problem:

$$(P) \min \quad x_1^2 + x_2^2 + 2x_1 \\ \text{s.t. } x_1 + x_2 \leq 0$$

after adding Lagrange multiplier:

$$\max \quad x_1^2 + x_2^2 + 2x_1 + \mu(x_1 + x_2) \\ x_1 + x_2 \leq 0$$

But weak duality still holds! Suppose feasible set of (P) is  $S$ .

$$\begin{aligned} q(\mu) &= \min_{x \in \mathbb{R}^2} x_1^2 + x_2^2 + 2x_1 + \mu(x_1 + x_2) \\ &\leq \min_{x \in S} x_1^2 + x_2^2 + 2x_1 + \mu(x_1 + x_2) \\ &\leq \min_{x \in S} x_1^2 + x_2^2 + 2x_1 \quad \text{because } x_1 + x_2 \leq 0 \text{ for } x \in S \\ &= f^* \end{aligned}$$

Now we rigorously state the definition of dual problem: Given general model (called primal problem in this section)

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

where  $f, g_i, h_j$  are functions with NO assumption (for now). Recall the Lagrangian  $L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$  defined in the chapter 6, define dual objective function

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

note all constraints are dropped for this minimisation. You may notice that there may not be a value  $\mathbf{x}$  that attains this infimum, and value of  $q$  can be  $-\infty$ . So an additional requirement that  $q(\boldsymbol{\lambda}, \boldsymbol{\mu}) > -\infty$  is required.

$$\begin{aligned} (D) \min \quad & q(\boldsymbol{\lambda}, \boldsymbol{\mu}) \\ \text{s.t. } & (\boldsymbol{\lambda}, \boldsymbol{\mu}) \in \text{dom}(q) \end{aligned}$$

where

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

Assume  $f, g_i, h_j$  in (P) are finite-valued. It can be shown that  $\text{dom}(q)$  is convex

*Proof.* Given  $(\boldsymbol{\lambda}_1, \boldsymbol{\mu}_1), (\boldsymbol{\lambda}_2, \boldsymbol{\mu}_2) \in \text{dom}(q)$ , since Lagrangian  $L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$  is affine w.r.t.  $\boldsymbol{\lambda}, \boldsymbol{\mu}$ ,

$$\begin{aligned} q(\alpha \boldsymbol{\lambda}_1 + (1 - \alpha) \boldsymbol{\lambda}_2, \alpha \boldsymbol{\mu}_1 + (1 - \alpha) \boldsymbol{\mu}_2) &= \min_{\mathbf{x} \in X} L(\mathbf{x}, \alpha \boldsymbol{\lambda}_1 + (1 - \alpha) \boldsymbol{\lambda}_2, \alpha \boldsymbol{\mu}_1 + (1 - \alpha) \boldsymbol{\mu}_2) \\ &= \min_{\mathbf{x} \in X} [\alpha L(\mathbf{x}, \boldsymbol{\lambda}_1, \boldsymbol{\mu}_1) + (1 - \alpha) L(\mathbf{x}, \boldsymbol{\lambda}_2, \boldsymbol{\mu}_2)] \\ &\geq \alpha \min_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\lambda}_1, \boldsymbol{\mu}_1) + (1 - \alpha) \min_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\lambda}_2, \boldsymbol{\mu}_2) \\ &= \alpha q(\boldsymbol{\lambda}_1, \boldsymbol{\mu}_1) + (1 - \alpha) q(\boldsymbol{\lambda}_2, \boldsymbol{\mu}_2) \\ &> -\infty \quad \text{because } (\boldsymbol{\lambda}_1, \boldsymbol{\mu}_1), (\boldsymbol{\lambda}_2, \boldsymbol{\mu}_2) \in \text{dom}(q) \end{aligned}$$

□

From the above proof, you may have also noticed that since  $L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$  is affine,  $q(\boldsymbol{\lambda}, \boldsymbol{\mu})$  is minimum of concave functions, so  $q$  is concave. (follows from Theorem 4.9) Then (D) is maximisation problem of concave  $q$  over convex set, which is essentially minimisation of  $-q$ . So (D) is a convex problem.

As mentioned before, dual problem is used to find lower bound for  $f^*$ , the optimal value of  $f$ ,

**Theorem 7.1** (Weak Duality Theorem). *If  $q^*$  is optimal solution of (D),  $f^*$  is optimal solution of (P), then  $q^* \leq f^*$ .*

*Proof.* Similar to the proof used for minimisation of  $x_1^2 + x_2^2 + 2x_1$  discussed before. See page 102 of Lecture notes for detailed proof.  $\square$

Duality problem may give useless lower bound,

**Example 9.** Consider

$$(P) \min \quad x_1^2 - 3x_2^2 \\ \text{s.t. } x_1 = x_2^3$$

This function is not convex, and the equality constraint is not affine.

Using substitution, one can find that  $(1, 1), (-1, -1)$  are optimal solutions giving  $f^* = -2$ . But the Lagrangian is

$$L(x_1, x_2, \mu) = x_1^2 + \mu x_1 - 3x_2^2 - \mu x_2^3$$

if  $\mu \geq 0$ , plug in  $(0, x_2)$  and send  $x_2 \rightarrow \infty$  gives  $-\infty$ , or if  $\mu < 0$ , plug in  $(x_1, 0)$  and send  $x_2 \rightarrow \infty$  gives  $-\infty$ . So for any  $\mu \in \mathbb{R}$ , minimisation of  $L$  is  $-\infty$ , and that means dual optimal value is  $q^* = -\infty$ .

## 7.1 Strong Duality

Convexity leads us to strong duality property, where optimal value of primal and dual problem coincides. Proof of strong duality theorem requires a variation of separation theorem, where assumption "closed" of the convex set is removed. Now instead of requiring "separation", we allow the point to sit on boundary of the convex set

**Theorem 7.2** (Supporting Hyperplane Theorem). *Given convex set  $C$  and point  $\mathbf{y} \notin C$ . Exists  $\mathbf{p} \neq \mathbf{0}$  s.t.*

$$\mathbf{p}^T \mathbf{x} \leq \mathbf{p}^T \mathbf{y} \quad \forall \mathbf{x} \in C$$

*Proof.* We play some topological tricks to allow us to apply separation theorem on  $\overline{C}$ , the closure of  $C$ , which is closed by definition and convex. Unfortunately, we cannot guarantee  $\mathbf{y} \notin \overline{C}$ , but since  $\mathbf{y} \notin \text{int}(C) = \text{int}(\overline{C})$ , we can find sequence of  $\mathbf{y}_k \notin \overline{C}$  s.t.  $\mathbf{y}_k \rightarrow \mathbf{y}$ . Then applying separation theorem, there exists  $\mathbf{p}_k \neq \mathbf{0}$  s.t.

$$\mathbf{p}_k^T \mathbf{x} < \mathbf{p}_k^T \mathbf{y}_k \quad \forall \mathbf{x} \in \overline{C}$$

already close to our objective. But  $\mathbf{p}_k$  may not converge, and further, this sequence may not be bounded! But, we can derive the following inequality:

$$\frac{\mathbf{p}_k^T}{\|\mathbf{p}_k\|} (\mathbf{x} - \mathbf{y}_k) < 0 \quad \forall \mathbf{x} \in \overline{C}$$

the sequence  $(\frac{\mathbf{p}_k^T}{\|\mathbf{p}_k\|})_{k \in \mathbb{N}}$  is bounded (norm of the whole sequence is just 1), so by Weierstrass theorem, there is a convergent subsequence, say the limit is  $\mathbf{p}$ . We have  $\|\mathbf{p}\| = 1$  so  $\mathbf{p} \neq \mathbf{0}$ . Using this convergent subsequence on the inequality yields

$$\mathbf{p}^T (\mathbf{x} - \mathbf{y}) \leq 0 \quad \forall \mathbf{x} \in \overline{C}$$

the result follows as  $C \subseteq \overline{C}$ .  $\square$

**Theorem 7.3** (Separation of convex sets). *Given non-empty convex sets  $C_1, C_2$  with  $C_1 \cap C_2 = \emptyset$ . Then exists  $\mathbf{0} \neq \mathbf{p}$  s.t.*

$$\mathbf{p}^T \mathbf{x} \leq \mathbf{p}^T \mathbf{y} \quad \forall \mathbf{x} \in C_1, \mathbf{y} \in C_2$$

*Proof.* Use Theorem 7.2 with  $C_1 - C_2$  (prove it is convex) and  $\mathbf{0}$ . Vector  $\mathbf{0} \notin C_1 - C_2$  as  $C_1, C_2$  are disjoint.  $\square$

Farkas lemma can be generalised to non-linear cases, with Slater-type conditions.

**Theorem 7.4** (Nonlinear Farkas Lemma). *Given convex set  $X$  and convex functions  $f, g_1, \dots, g_m$  (not necessarily continuous), assume exists  $\hat{\mathbf{x}} \in X$  s.t.  $g_i(\hat{\mathbf{x}}) < 0 \forall i$ . Given  $c \in \mathbb{R}$ , the following two are equivalent:*

- (a)  $\mathbf{x} \in X, g_i(\mathbf{x}) \leq 0 \Rightarrow f(\mathbf{x}) \geq c$
- (b) exists  $\lambda_1, \dots, \lambda_m \geq 0$  s.t.

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

note the function inside is  $L(\mathbf{x}, \boldsymbol{\lambda})$  when  $g_i$  are used as inequality constraints.

The proof is similar to that of Farkas lemma but rather long. See (page 242-243, Beck, 2014)

**Theorem 7.5** (Strong Duality Property of convex problems). *Given optimisation problem*

$$(P) \min \quad f(\mathbf{x}) \\ \text{s.t. } g_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, m$$

with  $X$  convex,  $f, g_i$  convex functions. Assume exists  $\hat{\mathbf{x}} \in X$  s.t.  $g_i(\hat{\mathbf{x}}) < 0 \forall i$ . (this is exactly Slater's condition). If the optimal value  $f^*$  of problem (P) is finite, then optimal value of dual problem is attained, and

$$q^* = f^*$$

where  $q^*$  is optimal value of the dual problem, i.e.

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

*Proof.* By optimality of  $f^*$ ,

$$\mathbf{x} \in X, g_i(\mathbf{x}) \leq 0 \Rightarrow f(\mathbf{x}) \geq f^* \forall \mathbf{x} \in X$$

So by nonlinear Farkas's lemma, there are  $\lambda_i \geq 0$  s.t.

$$q(\boldsymbol{\lambda}) = \min_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\lambda}) \geq f^*$$

So combining with weak duality theorem,

$$q^* \geq q(\boldsymbol{\lambda}) \geq f^* \geq q^*$$

which means  $f^* = q^*$ . □

When the Slater's condition is not satisfied, dual problem may still give the strong dual property, but optimal value of (D) is not attained.

**Example 10.**

$$\min \quad x_1^2 - x_2 \\ \text{s.t. } x_2^2 \leq 0$$

The problem is convex but not satisfy Slater's condition. Optimal solution is simply (0,0) giving  $f^* = 0$ . The Lagrangian is  $L(x_1, x_2, \lambda) = x_1^2 - x_2 + \lambda x_2^2$  giving dual objective

$$q(\lambda) = \min_{x_1, x_2 \in \mathbb{R}} L(x_1, x_2, \lambda) = \begin{cases} -\infty, & \lambda = 0 \\ -\frac{1}{4\lambda}, & \lambda > 0 \end{cases}$$

So dual problem has optimal value  $q^* = 0$ , but not attained by any  $\lambda \geq 0$ .

Equality conditions can be added to Theorem 7.5

**Theorem 7.6.** *Consider*

$$(P) \min \quad f(\mathbf{x}) \\ \text{s.t. } g_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, m \\ h_j(\mathbf{x}) \leq 0, \quad j = 1, 2, \dots, p \\ s_k(\mathbf{x}) = 0, \quad k = 1, 2, \dots, q \\ \mathbf{x} \in X$$

where  $X$  is convex,  $f, g_i$  are convex,  $h_j, s_k$  are affine. Suppose generalised Slater's condition is satisfied, then if optimal value  $f^*$  of (P) is finite, optimal of dual problem is attained and

$$f^* = q^*$$

where

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

One can prove that complementary slackness condition (similar to that of KKT2) holds as long as  $q^* = f^*$  (even if  $f$  is not convex)

**Theorem 7.7** (Complementary Slackness). *For problem*

$$\begin{aligned} \min \quad & x_1^2 - x_2 \\ \text{s.t.} \quad & x_2^2 \leq 0 \end{aligned}$$

*if  $f^* = q^*$ , then if  $\mathbf{x}^*, \boldsymbol{\lambda}^*$  are optimal solutions of primal and dual problem respectively, then*

$$\mathbf{x}^* \in \text{Argmin} L(\mathbf{x}, \boldsymbol{\lambda}^*)$$

$$\lambda_i^* g_i(\mathbf{x}^*) = 0 \quad \forall i$$

## References

- [1] Beck, Amir. Introduction to Nonlinear Optimization. SIAM, 2014.
- [2] Chen, Xuemei. Kaczmarz Algorithm, Row Action Methods, and Statistical Learning Algorithms. [faculty.sites.iastate.edu/esweber/files/inline-files/KaczmarzSGDrevised.pdf](http://faculty.sites.iastate.edu/esweber/files/inline-files/KaczmarzSGDrevised.pdf).
- [3] S. Kaczmarz. Angenäherte auflösung von systemen linearer gleichungen. Bulletin International de l'Académie Polonaise des Sciences et des Lettres A, 35:355–357, 1937.
- [4] Weisstein, Eric W. “Ellipsoid.” Mathworld.wolfram.com, [mathworld.wolfram.com/Ellipsoid.html](http://mathworld.wolfram.com/Ellipsoid.html).
- [5] Armijo, Larry. “Minimization of Functions Having Lipschitz Continuous First Partial Derivatives.” Pacific Journal of Mathematics, vol. 16, no. 1, Jan. 1966, pp. 1–3, <https://doi.org/10.2140/pjm.1966.16.1>.
- [6] Sun, T., Cheng, L. Convergence of iterative hard-thresholding algorithm with continuation. Optim Lett 11, 801–815 (2017). <https://doi.org/10.1007/s11590-016-1062-0>