

# מבוא לאופטימיזציה

## 236330

הרצאות ווידאו, אביב 2010  
פרופסור חבר מיכאל ציבולבסקי

נכתב ע"י רמי נודלמן

אביב 2015

## תוכן עניינים

3	הרצאה 1a – אלגוריתמים של אופטימיזציה לא-לינארית
8	הרצאה 1b – חזרה על אלגברה לינארית
13	הרצאה 02-03 – אלגוריתמים של אופטימיזציה לא לינארית. נגזרות של פונקציות מרובות משתנים: גרדיאנט והסיאן
25	הרצאה 04-05 – Convex Sets and Functions, קבוצות קמורות ופונקציות קמורות
32	הרצאה 06 – Local and Global Minimum, מינימום מקומי וגלובאלי
	הרצאה 07 – Iterative Methods of One Dimensional Optimization, אלגוריתמים איטרטיביים (נומריים) של אופטימיזציה במשתנה יחיד
35	הרצאה 08 – Multidimensional, Unconstrained Optimization Methods, אלגוריתמים לאופטימיזציה ללא אילוצים של פונקציות בעלות מספר משתנים
41	הרצאה 09 – Another view of Newton's Meth. Via solution of system of nonlinear equations, מבט נוסף על שיטת ניוטון באמצעות פתירת מערכת משוואות לא לינאריות
47	הרצאה 10 – Conjugate Gradient Method, שיטת הווקטורים האורתוגונליים לגרדיאנט
52	הרצאה 11 – Conjugate Gradient Method 2, שיטת הווקטורים האורתוגונליים לגרדיאנט - המשך
56	הרצאה 12 – Sequential Subspace Optimization (SESOP) and Quasi-Newton's Method
64	הרצאה 13 – Summary Of Unconstrained Optimization And Intro To Optimization With Constraints
71	הרצאה – Penalty Function Method and Augmented Lagrangian Method For Constrained Optimization
78	הרצאה 14 – Lagrange Multipliers and Penalty Function Method. Augmented Lagrangian
82	הרצאה 15 – Minmax, Game Theory, Lagrangian Duality
89	הרצאה 16 – Conic Programming, תכנות קוני
96	הרצאה 17 – Conversion of different problems to SDP (Semi Definite Programming)
112	

## הרצאה 1a – אלגוריתמים של אופטימיזציה לא-לינארית

הגדרות:

אופטימיזציה ללא אילוצים:

אופטימיזציה של פונקציה מסוימת  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  ללא אילוצים, זה למעשה מציאת הנקודה  $\bar{x}$  אשר תביא את ערך הפונקציה  $f(\bar{x})$  למינימום שלה, כלומר:

$$\arg \min_{\bar{x} \in \mathbb{R}^n} f(\bar{x})$$

כאשר אין אילוצים על התחום של הנקודות  $\bar{x}$  (בתנאי שהן נמצאות בתחום ההגדרה של הפונקציה). בשלב זה אנחנו מניחים כי  $f(\bar{x})$  היא פונקציה רציפה ו"חלקה" (כלומר גזירה).

אופטימיזציה עם אילוצים:

באופטימיזציה עם אילוצים אנחנו מנסים למצוא את הנקודה  $\bar{x}$  שתביא את הפונקציה  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  לערך המינימום, כלומר:

$$\arg \min_{\bar{x} \in \mathbb{R}^n} f(\bar{x})$$

אך כעת אנחנו נותנים הגבלות על הנקודות האפשריות (מחפשים נקודות רק בתחומים מוגדרים מסוימים), כלומר מגבילים את התחום ע"י אילוצים כלשהם. ואנחנו מניחים כי  $f(x)$  היא פונקציה רציפה ו"חלקה" (גזירה ברציפות) וכן האילוצים הם למשל:

$$g_i(\bar{x}) \geq 0, \quad i = 1, \dots, m$$

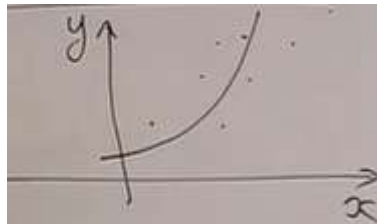
$$h_j(\bar{x}) = 0, \quad j = 1, \dots, l$$

האילוץ הראשון הוא אילוץ שנוצר ע"י אי-שוויון (Inequality constrain) והשני נוצר על ידי שוויון (Equality Constrain).

הקדמה ודוגמאות:

בקורס אנחנו נעסוק בחיפוש התנאים האופטימליים, כלומר חיפוש התנאים על הפונקציה שעבורם נקבל את הפתרון האופטימלי וכן נעסוק באלגוריתמים איטרטיביים (שמתבססים על חישובים נומריים) ששיגו את הפתרון האופטימלי (בין אם האלגוריתם ייתן לנו את הפתרון המדויק או רק נוכל להגיע אליו בצורה "גבולית", בצורה ששואפת לפתרון האופטימלי).

דוגמה: (דוגמה אופטימיזציה ללא אילוצים) (Parametric Regression)



נניח ויש לנו מידע אמפירי, שבו  $x$  הוא הארגומנט ואילו  $y$  אלו המדידות, כלומר יש לנו נקודות על מערכת הצירים. כעת, אנחנו מתבוננים בנקודות הללו ומניחים שהקשר בין  $x$  לבין  $y$  הוא ריבועי, כלומר  $y = f(x) = x^2$ , אך אנחנו רואים כי יש רעש/סטייה בין מה שאנחנו רוצים לראות במדויק לבין תוצאות הניסוי.

ננסה את הבעיה בצורה פורמלית:

בהינתן קבוצה של זוגות  $\{x^{(i)}, y^{(i)}\}_{i=1}^M$  אנחנו רוצים למצוא פונקציה ריבועית  $f(x)$  עם ווקטור מקדמים  $\bar{w}$ :

$$f(x, \bar{w}) = w_0 + w_1 x + w_2 x^2$$

כעת אנחנו רוצים לייצג את הבעיה כבעיה אופטימיזציה והאופטימיזציה היא למעשה מציאת הווקטור  $\bar{w}$  אשר יביא לנו שגיאה מינימלית – מהי השגיאה? הביטוי לשגיאה הוא הביטוי שאנחנו נגדיר עבור כל בעיה לגופה. נניח ואנחנו רוצים להתבונן בשגיאה הריבועית, כלומר סכום ריבועי ההפרשים בין הפונקציה שאנחנו מציעים לבין המדידות בניסוי:

$$\min_{\bar{w}} \sum_{i=1}^M \left( f(x^{(i)}, \bar{w}) - y^{(i)} \right)^2$$

כאשר  $y^{(i)}$  זו המדידה ה- $i$ -ית.

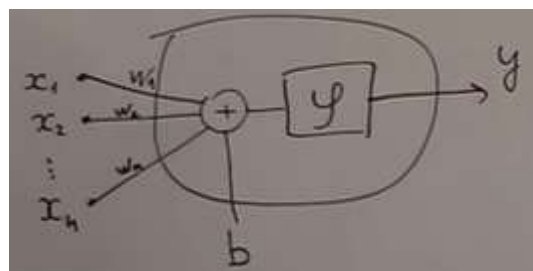
במקרה הכללי יותר, במקום להניח כי לפונקציה  $f(x)$  יש מקדמים קבועים, ניתן להניח כי המקדמים הם בעצמם פונקציות של הארגומנט  $x$ , למשל:

$$f(x, \bar{w}) = \sum_k w_k g_k(x)$$

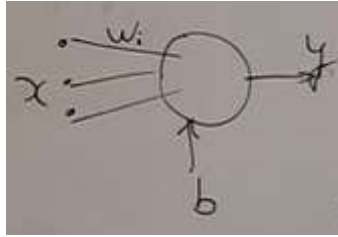
כאשר הפונקציות  $\{g_k\}$  ידועות. במקרה זה, הבעיה נקראת Linear Regression Problem. אמנם הפונקציות  $\{g_k\}$  אינן בהכרח לינאריות אך הפונקציה  $f(x)$  היא לינארית ביחס לווקטור  $\bar{w}$ .

לעיתים נתקלים בבעיות בהן הפונקציה  $f(x)$  היא לא לינארית ביחס לווקטור  $\bar{w}$ , כלומר כל אחד מהרכיבים של הווקטור  $\bar{w}$  הוא לא מקדם קבוע אלא פונקציה בעצמו, כלומר  $\bar{w} = \bar{w}(x)$ .

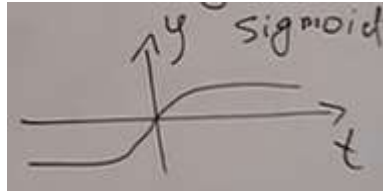
דוגמה: (Non Linear Regression) – למשל Neural Networks (מערכות עצבים)



נניח ואנחנו רוצים למצוא אופטימיזציה של פונקציה בעלת מספר משתנים,  $\bar{x} = \{x_1, x_2, \dots, x_n\}$  כאשר  $b$  הוא קבוע. נתייחס לרכיבים של ווקטור  $\bar{x}$  כאל "כניסות למערכת", כאשר המערכת היא נורון בודד והפלט של נורון בודד הוא פונקציה של הכניסות שלו. נניח כי הפלט הוא סכום של כניסות מסוימות כאשר לכל כניסה יש משקל שונה לאחר מכן, הסכום הזה נכנס לפונקציה  $\phi$  ואנחנו מקבלים פלט  $y$ . ניתן לכתוב את המערכת הנ"ל בצורה הבאה:



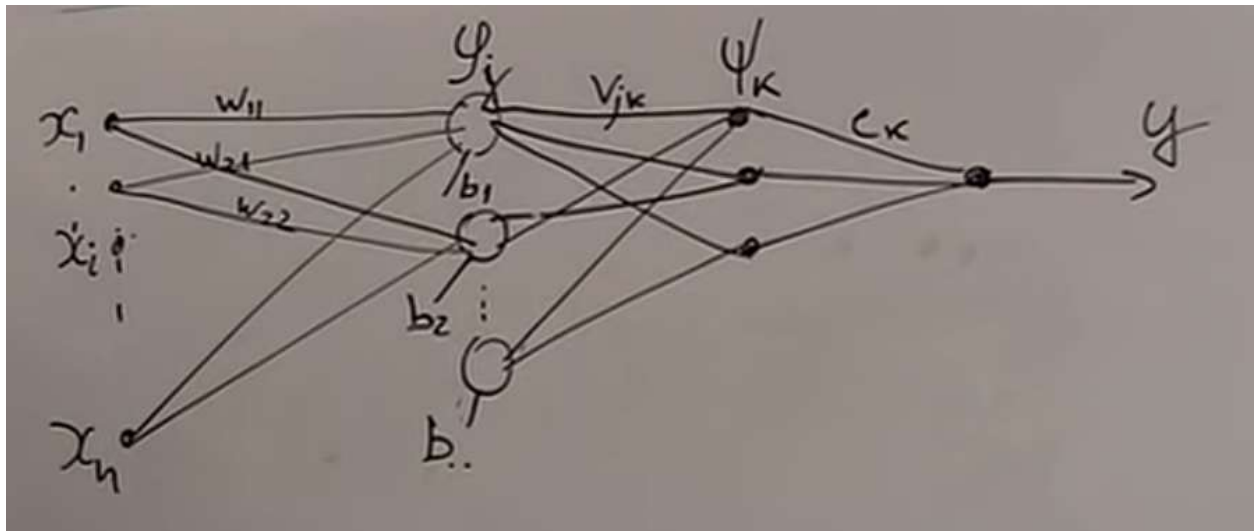
פונקציה זו, שהנורון פולט נקראת sigmoid והיא נראת כמו  $\arctan(x)$  :



ובצורה פורמלית ניתן לתאר את הפלט של הנורון בצורה הבאה:

$$y = \varphi(\bar{w}^T \cdot \bar{x} + b)$$

נתבונן כעת ברשת של נוירונים. נניח כי יש לנו "כניסות"  $\bar{x} = \{x_1, x_2, \dots, x_n\}$  וכן הנורון ה- $i$  מיוצג ע"י הפונקציה  $\varphi_i$ . כלומר נוירון מקבל את כל הכניסות אבל במשקלים שונים. לאחר הרמה הראשונה של הנוירונים, ייתכן ויש עוד נוירונים אך אנחנו נפשט את המודל ונניח כי יש לנו רק רמה אחת של נוירונים ולאחר מכן יש לנו מערכות לינאריות,  $\psi_k$  ונקבל:



כלומר פלט הרשת כולה הוא:

$$y(\bar{x}) = \sum_k \psi_k \left( \sum_j \varphi_j \left( \left( \sum_i x_i w_{ij} \right) + b_j \right) \right)$$

נסמן את הרכיב של ווקטור  $\bar{u}$  להיות:  $u_j = \left( \sum_i x_i w_{ij} \right) + b_j$  ובכתיב מטריצי נקבל:  $u_j = W^T \cdot \bar{x} + b_j$ , וכן נקבל:

$$y(\bar{x}) = \sum_k \psi_k \left( \sum_j \varphi_j(u_j(\bar{x})) \right)$$

נוכל לכתוב זאת גם בצורה הבאה:

$$\varphi(u) \triangleq \begin{pmatrix} \varphi_1(u_1(\bar{x})) \\ \varphi_2(u_2(\bar{x})) \\ \vdots \end{pmatrix}$$

ולכן נקבל:

$$y = f(\bar{x}, W, \bar{b}, V, \bar{c}) = c^T \psi_k \left( V^T \varphi_j(W^T \cdot \bar{x} + b) \right)$$

(כאשר אותיות גדולות הן מטריצות)

כלומר קיבלנו מערכת שאנחנו רוצים לבצע עליה אופטימיזציה ביחס למשתנים:  $W, \bar{b}, V, \bar{c}$ .

כעת נניח ויש לנו שוב זוגות של קואורדינטות ונתונים של ניסוי אך הפעם הארגומנט הוא ווקטור ותוצאת הניסוי היא סקלר  
 כלומר:

$$\{\bar{x}^{(i)}, y^{(i)}\}_{i=1}^M$$

ואנחנו נרצה למצוא פונקציה (רשת נוירונים) כזו שעבור סט נתונים מסויים, הפונקציה תשערך את המדידות  $y^{(i)}$  באופן הקרוב ביותר. ובאופן פורמלי, נרצה למצוא למשל את סכום השגיאות הריבועיות המנימלי:

$$\min_{W, \bar{b}, V, \bar{c}} \sum_{i=1}^M \left[ f(\bar{x}, W, \bar{b}, V, \bar{c}) - y^{(i)} \right]^2$$

זו כמובן בעיה אופטימיזציה לא ליניארית ביחס לפרמטרים  $W, \bar{b}, V, \bar{c}$ , פרט ביחס לווקטור  $\bar{c}$ , שביחס אליו הבעיה היא דווקא כן ליניארית שכן נזכור כי הווקטור  $\bar{c}$  נמצא מחוץ לפונקציה  $\psi$ . בקורס זה נלמד, איך למצוא אופטימיזציה גם לסוג בעיות כאלה.

דוגמה: בעיה אופטימיזציה עם אילוצים (Resource Assignment)

למשל, יש לנו מספר תחנות כוח ומספר לקוחות ויש מגבלות על כמות הכוח שכל תחנה מסוגלת לספק וכן לכל לקוח יש דרישות שיש לקיים וגם יש עלויות שונות עבור ההובלה/שינוע של האנרגיה מתחנה מסויימת ללקוח. אנחנו נרצה לבצע אופטימיזציה על העלות של הבעיה, כלומר למצוא פתרון אופטימלי של חלוקת האנרגיה בין תחנות הכוח השונות לבין הלקוחות השונים במחיר הזול ביותר.

נתונים:

יש לנו  $M$  תחנות כוח.  $s_m$  זה הכוח המקסימלי שתחנת כוח  $m$  מסוגלת לייצר.

יש לנו  $N$  לקוחות.  $c_n$  זו דרישת הכוח של הלקוח  $n$ -י.

$p_{mn}$  זו העלות של שינוע יחידת אנרגיה מתחנת כוח  $m$  ללקוח  $n$ .

המטרה:

למצוא את מטריצת הנעלמים  $X$  שבה  $x_{m,n}$  מייצג את כמות יחידות האנרגיה שתחנה  $m$  מספקת ללקוח  $n$ .

באופן פורמלי הבעיה היא:

$$\min_X \sum p_{mn} \cdot x_{mn}$$

האילוצים הם:

- כל תחנה מספקת לכל היותר את כל הכוח שהיא מסוגלת:  $\forall 1 \leq m \leq M : \sum_{n=1}^N x_{mn} \leq s_m$ .
- כל לקוח מקבל את כל האנרגיה שהוא דרש:  $\forall 1 \leq n \leq N : \sum_{m=1}^M s_{mn} = c_n$ .
- כל תחנה מספקת אנרגיה ללקוח (ולקוח לא מספק אנרגיה לתחנה הכוח) כלומר מתקיים:  $x_{mn} \geq 0$ .

בבעיה זו אנחנו רואים כי יש אילוצים מסוג אי-שוויון וגם שוויון.

כל האילוצים הנ"ל הם לינאריים ביחס ל-  $X$  ולכן זו בעיה אופטימיזציה לינארית, נקראת: Linear Programming Problem.

## הרצאה 1b – חזרה על אלגברה לינארית

$\mathbb{R}$  - מרחב אוקלידי מממד 1 (הציר הממשי).

$\mathbb{R}^n$  - מרחב אוקלידי מממד  $n$ .

$$\bar{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^n \text{ : (ווקטור עמודה): } \bar{x} \in \mathbb{R}^n$$

מרחב ווקטורי (לינארי):

יהי  $S$  מרחב ווקטורי וכן  $\bar{x}, \bar{y} \in S$  אז מתקיים:  $\alpha\bar{x} + \beta\bar{y} \in S$  לכל  $\alpha, \beta \in \mathbb{R}$ .

מרחב ווקטורי אפיני:

נאמר כי  $Q$  מרחב אפיני וכן  $\bar{b} \in S$  אם מתקיים:

$$Q = S + \bar{b} = \{ \bar{y} = \bar{b} + \bar{x} \mid \bar{x} \in S \}$$

אם למשל המרחב  $S$  הוא מממד 2 אז אפשר לומר כי  $Q$  הוא מרחב אפיני אם הוא שווה למרחב  $S$  אך מוזן ע"י הווקטור  $\bar{b}$ .

מרחב אפיני מקיים את תכונת הלינאריות, כלומר אם  $x, y \in Q$  אז:

$$\alpha\bar{x} + \beta\bar{y} \in Q$$

כאשר האילוח על  $\alpha, \beta$  הוא:  $\alpha + \beta = 1$  לכל  $\alpha, \beta$ .

הגדרה: (נורמה ווקטורית)

נניח כי  $\bar{x} \in \mathbb{R}^n$ . הנורמה של הווקטור  $\bar{x}$  היא פונקציה שמעבירה אלמנטים ממרחב  $\mathbb{R}^n$  למרחב  $\mathbb{R}$ . הנורמה של  $\bar{x} \in \mathbb{R}^n$  מסומנת להיות  $\|\bar{x}\|$ .

התכונות של כל נורמה:

$$(1) \text{ לכל } \bar{x} \in \mathbb{R}^n \text{ מתקיים: } \|\bar{x}\| \geq 0$$

$$(2) \text{ לכל } \bar{x} \in \mathbb{R}^n, \alpha \in \mathbb{R} \text{ מתקיים: } \|\alpha\bar{x}\| = |\alpha| \cdot \|\bar{x}\|$$

$$(3) \|\bar{x}\| = 0 \Leftrightarrow \bar{x} = \bar{0}$$

$$(4) \text{ אי-שיוויון המשולש: } \|\bar{x} + \bar{y}\| \leq \|\bar{x}\| + \|\bar{y}\|$$

דוגמאות לנורמות:

• נורמה אוקלידית ( $\mathbb{R}^n$ ):

$$\forall \bar{x} \in \mathbb{R}^n : \|\bar{x}\|_2 = \sqrt{\bar{x}^T \cdot \bar{x}} = \sqrt{\sum_{i=1}^n |x_i|^2}$$



• ההכללה של נורמה אוקלידית למרחב  $\mathbb{R}^p$  ( $l_p$  norm) :

$$\forall \bar{x} \in \mathbb{R}^p : \|\bar{x}\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{1/p}$$

מקרה פרטי חשוב של  $l_p$  norm הוא  $l_1$  וזה למעשה:

$$\forall \bar{x} \in \mathbb{R}^1 = \mathbb{R} : \|\bar{x}\|_1 = \sum_{i=1}^n |x_i|$$

ומקרה פרטי נוסף של  $l_p$  norm הוא  $l_\infty$  וזה למעשה:

$$\forall \bar{x} \in \mathbb{R}^1 = \mathbb{R} : \|\bar{x}\|_\infty = \lim_{x \rightarrow \infty} \|\bar{x}\|_p = \max_i |x_i|$$

נורמה של מטריצות:

יש כמה סוגים של נורמות של מטריצות. אחד מהם הוא דומה לנורמה אוקלידית, שמחשיב למעשה את המטריצה כווקטור אחד ארוך. סוג שני נקרא הנורמה המטריצית לפי פרוביני (Frobenius Norm) :

$$A = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{pmatrix}, \quad \|A\|_F \triangleq \sqrt{\sum_{i,j} (a_{ij})^2}$$

ניתן לראות את ההגדרה האחרונה גם בצורה אחרת וקומפקטית יותר:

$$\text{Trace}(A^T A)$$

$$(\text{Trace}(A) = \sum_{i=1}^n a_{ii} \text{ כאשר})$$

ניתן גם להגדיר סוג אחר של נורמה על מטריצות אם נחשוב על מטריצות כעל אופרטורים לינארים. מטריצות הן למעשה אופרטורים שמעבירים ווקטורים ממרחב אחד לאחר ובין היתר משנים את הנורמה של הווקטור עליו הן פועלות. נגדיר את הנורמה של מטריצה להיות הנורמה המקסימלית שווקטור באורך יחידה שעליו פועלת המטריצה מקבל (יוסבר בהמשך). נורמה כזו נקראת Induced matrix norm או במילים אחרות Operator matrix norm. בצורה פורמלית:

$$\bar{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}, A = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{pmatrix}, \quad \|A\|_{p,q} \triangleq \max_{\|\bar{x}\|_p=1} \|A\bar{x}\|_q$$

כלומר הנורמה  $\|A\|_{p,q}$  היא הנורמה המקסימלית של הווקטור  $A\bar{x}$  תחת נורמה  $q$  כאשר  $\bar{x}$  הוא כל ווקטורי היחידה האפשריים.

מכפלה פנימית של מטריצות:

נניח  $A, B$  מטריצות ריבועיות מאותו סדר, אז המכפלה הפנימית של  $A, B$  :

$$\langle A, B \rangle \triangleq \sum_{i,j} a_{ij} b_{ij} = \text{Trace}(A^T B) = \text{Trace}(B A^T)$$

(הערה: המכפלה הפנימית הנ"ל מוגדרת בדומה לאיך שמוגדרת מכפלה פנימית בין שני ווקטורים ולמעשה מכפלה פנימית זו מתייחסת למטריצות כאל ווקטורים "ארוכים")

הערה: מה שמאפשר את שני המעברים האחרונים היא תכונת ההזזה הציקלית תת אופרטור ה-Trace:

$$\text{Trace}(VW) = \text{Trace}(WV)$$

ערכים עצמיים של מטריצות:

תהי  $A$  מטריצה מסדר  $n \times n$  וכן  $\bar{x} \in \mathbb{R}^n$  וגם  $\lambda \in \mathbb{R}$ . אם מתקיים:

$$A\bar{x} = \lambda\bar{x}$$

נאמר כי  $\bar{x}$  הוא ווקטור עצמי (ו"ע) וכן  $\lambda$  הוא ערך עצמי (ע"ע). המשמעות היא שכאשר מטריצה  $A$  פועלת על ווקטור  $\bar{x}$  היא לא משנה את כיוונו אלא רק משפיעה על מגמתו ו/או ערכו.

למטריצה מסדר  $n \times n$  יש לכל היותר  $n$  ו"ע-ים בלתי תלויים ו- $n$  ע"ע.

נרשום את הו"ע-ים ואת הע"ע-ים בצורה קומפקטית וכן נניח כי כל הו"ע הם בלתי תלויים. נרשום את הו"ע-ים של מטריצה  $S$  בתור עמודות המטריצה:

$$s_i = \begin{pmatrix} s_{1i} \\ \vdots \\ s_{ni} \end{pmatrix} \Rightarrow S = (\bar{s}_1 \quad \cdots \quad \bar{s}_n)$$

וכן נרשום:

$$\Lambda = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}$$

ולכן נוכל לרשום:

$$AS = S\Lambda$$

ואם הו"ע-ים הם בלתי תלויים אז המטריצה  $S$  מקיימת:  $\det(S) \neq 0$  ולכן יש למטריצה  $S$  מטריצה הופכית ולכן נוכל לכפול במטריצה ההופכית ונקבל:

$$AS = S\Lambda \Rightarrow A = SAS^{-1}$$

בנוסף, אם המטריצה  $A$  היא סימטרית אז כל הו"ע-ים של  $A$  הם אורתוגונליים ולכן נקבל:

$$S \cdot S^T = I$$

(כאשר מנרמלים את הו"ע העצמיים להיות בגודל של ווקטורי יחידה)

מהמשוואה האחרונה אנחנו רואים כי:  $S^T = S^{-1}$  ולכן נוכל במקרה שמטריצה  $A$  היא סימטרית לרשום:

$$A = SAS^{-1} = SAS^T$$

דוגמה:

באמצעות ע"ע נוח לבצע חישוב מסויימים על מטריצות – למשל חישוב מטריצה הופכית.

אם נניח כי  $A = SAS^{-1}$  אז זה אומר כי  $A^{-1} = S\Lambda^{-1}S^{-1}$  כי ניתן לראות כי מתקיים:

$$AA^{-1} = (S\Lambda S^{-1})(S\Lambda^{-1}S^{-1}) = S\Lambda(S^{-1}S)\Lambda^{-1}S^{-1} = S(\Lambda\Lambda^{-1})S^{-1} = SS^{-1} = I$$

בנוסף אפשר לחשב את החזקות של מטריצה. למשל:

$$A^2 = A \cdot A = SAS^{-1} \cdot SAS^{-1} = S\Lambda^2S^{-1}$$

וכיוון שמטריצה  $\Lambda$  היא אלכסונית זה מאד קל כיוון שפשוט מעלים בחזקה כל אחד מהרכיבים שלה:

$$\Lambda = \begin{pmatrix} \lambda_1 & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \lambda_n \end{pmatrix} \Rightarrow \Lambda^p = \begin{pmatrix} (\lambda_1)^p & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & (\lambda_n)^p \end{pmatrix}$$

פונקציות מטריציות: (בהקבלה לפונקציות סקלריות)

אנחנו רגילים שפונקציות פועלות על סקלרים או לכל היותר ווקטורים אך אין מניעה שתהיינה פונקציות שפועלות על מטריצות.

נניח ויש לנו פונקציה סקלרית כלומר  $\varphi: \mathbb{R} \rightarrow \mathbb{R}$  שמוגדרת באופן הבא:

$$\varphi(t) = \sum_i c_i t^i, \quad c_i \in \mathbb{R}, 0 \leq i \leq \infty$$

אז אפשר להגדיר באופן דומה פונקציה  $\varphi_A: n \times n \rightarrow n \times n$

$$\varphi(A) = \sum_i c_i A^i, \quad c_i \in \mathbb{R}, 0 \leq i \leq \infty$$

כיוון שראינו שאפשר לחשב חזקות של מטריצה באמצעות הע"ע של המטריצה ובאמצעות הו"ע-ים של המטריצה אפשר לרשום:

$$\begin{aligned} \varphi(A) &= \sum_i c_i A^i = \sum_i c_i S\Lambda^i S^{-1} = S \left( \sum_i c_i \Lambda^i \right) S^{-1} = S \begin{pmatrix} \sum_i c_i (\lambda_1)^i & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \sum_i c_i (\lambda_n)^i \end{pmatrix} S^{-1} \\ &= S \begin{pmatrix} \varphi(\lambda_1) & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \varphi(\lambda_n) \end{pmatrix} S^{-1} = S\varphi(\Lambda)S^{-1} \end{aligned}$$

הגדרה: חיוביות ואי-שליליות של מטריצות סימטריות

נאמר כי המטריצה  $A$  חיובית,  $A \succ 0$  (Positive definite), או אי-שלילית,  $A \succeq 0$  (Positive semi-definite) אם

$$\forall \bar{x} \neq \bar{0}: \quad \bar{x}^T A \bar{x} > 0 \quad (\text{Positive definite})$$

$$\forall \bar{x} \neq \bar{0}: \quad \bar{x}^T A \bar{x} \geq 0 \quad (\text{Positive (semi)definite})$$

משפט:

מטריצה היא חיובית (אי-שלילית) אם"מ כל הערכים העצמיים של המטריצה  $A$  הם חיוביים (אי-שליליים).

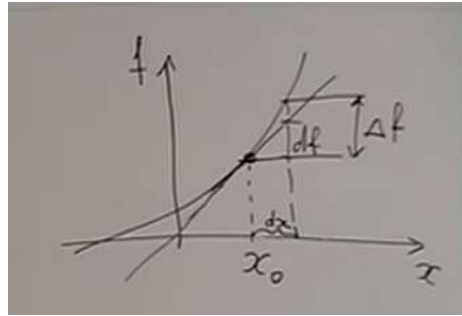
הוכחה:

$$\forall \bar{x} \neq \bar{0}: \quad \underbrace{\bar{x}^T A \bar{x}}_{\equiv \bar{y}^T} = \underbrace{\bar{x}^T S}_{\equiv \bar{y}} \Lambda \underbrace{S^T \bar{x}}_{\equiv \bar{y}} = \bar{y}^T \Lambda \bar{y} = \sum_i \lambda_i (y_i)^2 \geq 0 \quad \Leftrightarrow \quad \forall \lambda_i \geq 0$$

## הרצאה 02-03 – אלגוריתמים של אופטימיזציה לא לינארית. נגזרות של פונקציות מרובות משתנים: גרדיאנט והסיאן

תזכורת: פונקציה של משתנה יחיד

נניח כי  $f: \mathbb{R} \rightarrow \mathbb{R} : (x \in \mathbb{R})$ .



כאשר נשאוף את  $dx$  לאפס, כלומר  $dx \rightarrow 0$  אז נוכל לומר כי מתקיים:

$$\Delta f = df + o(dx)$$

ואנחנו קוראים לביטוי  $df$  דיפרנציאל ולכן אפשר לומר כי הדיפרנציאל של פונקציה  $f$  הוא החלק הלינארי של השינוי בפונקציה  $f$ .

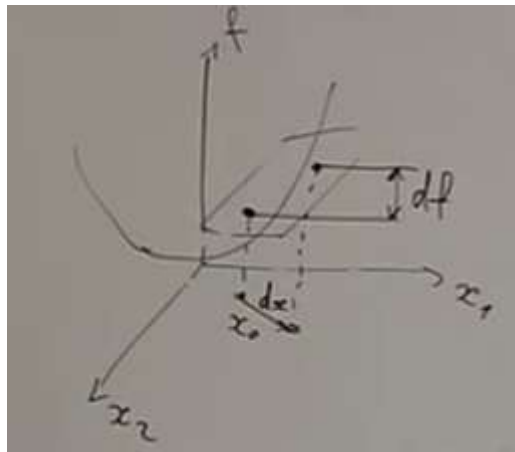
כמו כן, אנחנו יודעים מחדו"א שהמשיק מקיים:

$$df = f'(x_0)dx$$

פונקציה של מספר משתנים:

נניח כי  $f: \mathbb{R}^2 \rightarrow \mathbb{R} : (x_1, x_2)$  ונניח שאנחנו מסתכלים על המישור המשיק לנקודה  $(x_1, x_2)$  ואנחנו מתקדמים מנקודה

$$(x_1, x_2) \text{ ע"י וקטור } d\bar{x} = (dx_1, dx_2):$$



מחדו"א 2 ידוע כי הדיפרנציאל של פונקציה דיפרנציאבילית מרובת משתנים,  $f(x_1, x_2, \dots, x_n)$ , הוא הסכום:

$$df(\bar{x}) = \frac{\partial f}{\partial x_1} dx_1 + \frac{\partial f}{\partial x_2} dx_2 + \dots + \frac{\partial f}{\partial x_n} dx_n = \sum_{i=1}^n \frac{\partial f}{\partial x_i} dx_i$$

וניתן לרישום בצורה קומפקטית יותר כמכפלה פנימית:

$$df(\bar{x}) = \left\langle \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}, \begin{pmatrix} dx_1 \\ \vdots \\ dx_n \end{pmatrix} \right\rangle = \left\langle \overline{\left( \frac{\partial f}{\partial x} \right)}, d\bar{x} \right\rangle$$

ואפשר להבחין למעשה מתקיים כי הווקטור  $\overline{\left( \frac{\partial f}{\partial x} \right)}$  הוא למעשה  $\nabla_{\bar{x}} f(\bar{x})$  שנסמנו  $\bar{g}(\bar{x})$  וזה למעשה הגרדיאנט של

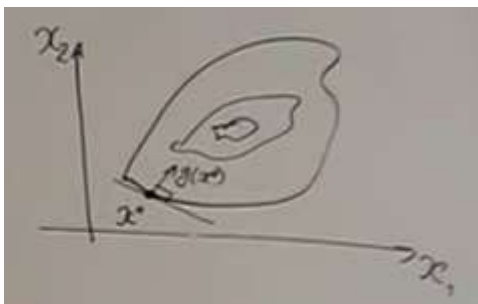
הפונקציה  $f(\bar{x})$ . לכן ניתן לרשום:

$$df(\bar{x}) = \langle \bar{g}(\bar{x}), d\bar{x} \rangle$$

הגדרה זו מאד שימושית כי היא תעזור לנו לבנות גרדיאנטים של פונקציות מסובכות שניתקל בהן במהלך הקורס, כי לבנות את הגרדיאנט לפי הגזרות החלקיות יהיה מאד מסובך.

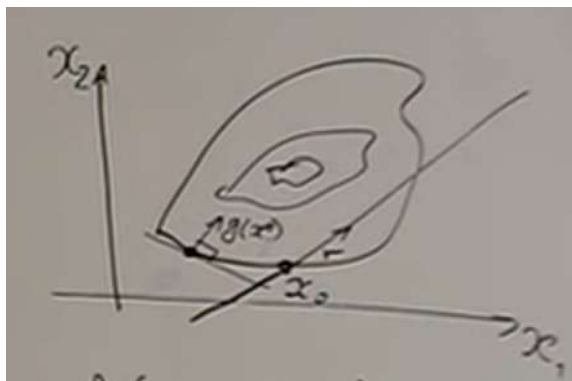
נגזרת כיוונית בכיוון כלשהו:

נתבונן שוב בפונקציה של שני משתנים:  $f(x_1, x_2): \mathbb{R}^2 \rightarrow \mathbb{R}$  ונצייר אותה בצורה של מפת טופוגרפית (הקווים בשרטוט הם קווי שווי ערך/גובה):



אנחנו נראה כי הגרדיאנט (שהוא בעצמו וקטור) מצביע על כיוון העליה החדה ביותר בערכי הפונקציה, וכן הוא מאונך למישור המשיק לנקודה.

אנחנו בעיקר נתעניין בהתנהגות הפונקציה כאשר אנחנו מתחילים בנקודה מסויימת  $\bar{x}_0$  בתחום ומתקדמים לאורך וקטור  $\bar{f}$ .



כלומר אנחנו רוצים להעריך את הפונקציה  $f$  בנקודות  $\bar{x}_0 + \alpha \bar{r}$  (כאשר  $\alpha \in \mathbb{R}$ ) ונוכל למעשה להגדיר פונקציה  $\varphi(\alpha): \mathbb{R} \rightarrow \mathbb{R}$  בצורה הבאה:

$$\varphi(\alpha) \triangleq f(\bar{x}_0 + \alpha \bar{r})$$

ואנחנו נרצה למצוא את הנגזרת של הפונקציה  $\varphi(\alpha)$  שהגדרנו בנקודה  $\bar{x}_0$  (כלומר בנקודה ההתחלתית שבה אנחנו מתבוננים בפונקציה  $f(\bar{x})$ ), כלומר למצוא את:

$$\varphi'(\alpha)|_{\alpha=0}$$

הביטוי  $\varphi'(\alpha)$  נקרא הנגזרת הכיוונית של  $f(\bar{x}_0)$  בכיוון  $\bar{r}$  ונסמנה להיות:

$$\boxed{\varphi'(\alpha)|_{\alpha=0} = \frac{\partial}{\partial \alpha} f(\bar{x}_0 + \alpha \bar{r}) \triangleq f_{\bar{r}}(\bar{x}_0)}$$

זוהי תוצאה מאד חשובה בקורס שנשתמש בה רבות.

דיפרנציאל וגרדיאנט של פונקציה מרובת משתנים:

השאלה שתעסיק אותנו כעת זה האם בהינתן הגרדיאנט בנקודה מסוימת  $\bar{g}(\bar{x}_0) \equiv \nabla f(\bar{x}_0)$  האם אנחנו יכולים לדעת מהי הנגזרת הכיוונית בכיוון כלשהו  $\bar{r}$ , כלומר האם ע"י הגרדיאנט ניתן לדעת מה היא  $f_{\bar{r}}(\bar{x}_0)$ ?

ניתן למצוא את התשובה באמצעות הצבה ישירה למשוואות שמצאנו עד כה אך גם ניתן להשתמש בהגדרה של הדיפרנציאל של הפונקציה  $f(\bar{x})$ .

לפי מה שראינו עבור משתנה יחיד:

$$df(\bar{x}) = f'(\bar{x}_0) d\bar{x}$$

מתקיים באופן דומה עבור פונקציה של מספר משתנים:

$$\boxed{df(\bar{x}) = \underbrace{\bar{g}^T(\bar{x}_0) \cdot d\bar{x}}_{\text{scalar if } c \in \mathbb{R} \Rightarrow c=c^T} = (d\bar{x})^T \cdot \bar{g}(\bar{x}_0)}$$

כעת, כיוון שהגדרנו  $\varphi(\alpha) \triangleq f(\bar{x}_0 + \alpha \bar{r})$ , אזי  $\varphi$  היא פונקציה של משתנה יחיד וע"י החלפת משתנים,

כאשר נחליף  $\bar{x} \triangleq \bar{x}_0 + \alpha \bar{r}$  ולכן נקבל  $d\bar{x} = d\alpha \bar{r}$  (כי  $\bar{x}_0$  קבוע), כלומר:

$$d\varphi(\alpha) = g^T(\bar{x}_0) \cdot d\bar{x} = g^T(\bar{x}_0) \cdot d\alpha \bar{r} = d\alpha \cdot g^T(\bar{x}_0) \cdot \bar{r}$$

$$\Rightarrow \frac{d\varphi}{d\alpha} = g^T(\bar{x}_0) \cdot \bar{r}$$

ולכן לפי הגדרת הנגזרת לפי הדיפרנציאל נקבל:

$$\boxed{f'_{\bar{r}}(\bar{x}_0) \triangleq \varphi'(\alpha)|_{\alpha=0} = \frac{d\varphi(\alpha)}{d\alpha} \Big|_{\alpha=0} = g^T(\bar{x}_0) \cdot \bar{r}}$$

(למעשה ניזכר כי ראינו תוצאה זו גם בחדו"א 2, כאשר אמרנו כי הנגזרת הכיוונית היא ההיטל של הגרדיאנט בכיוון שאנחנו מחפשים, בחדו"א 2 ביצענו גם נרמול לאורך הגרדיאנט ופה לא)

מטריצת ההסיאן (Hessian):

תזכורת: פונקציה של משתנה יחיד

הדיפרנציאל של פונקציה בעלת משתנה יחיד מקיים:

$$df = f'(x)dx$$

והנגזרת של דיפרנציאל של פונקציה של משתנה יחיד היא:

$$df' = f''(x)dx$$

כעת ננסה למצוא את הנגזרת של הדיפרנציאל של פונקציה עם מספר משתנים. ראינו כבר כי מתקיים עבר פונקציה עם מספר משתנים:

$$df(\bar{x}) = \bar{g}^T(\bar{x}) \cdot d\bar{x}$$

כעת נמצא את ההקבלה למשוואה של הנגזרת השניה של פונקציה עם משתנה יחיד שכאמור היא  $df' = f''(x)dx$ . אנחנו למעשה מחפשים את הדיפרנציאל של הגרדיאנט (כפי שעברנו מהביטוי  $df$  לביטוי  $df'$ ), אך הגרדיאנט הוא ווקטור ולכן גם הדיפרנציאל שלו הוא ווקטור, ולכן אנחנו מצפים שבמקום ווקטור שמוכפל ב-  $d\bar{x}$  אנחנו צריכים מטריצה (משיקולי מימדים), וזהו בדיוק ההסיאן:

$$\boxed{d\bar{g}(\bar{x}) = H(\bar{x}) \cdot d\bar{x}}$$

$$\text{ונהוג לסמן: } H(\bar{x}) \triangleq \nabla^2 f(\bar{x}) \equiv \nabla_{\bar{x}}^2 f(\bar{x}) \text{ . כאשר } \nabla^2 = \left( \frac{\partial^2}{\partial x_1^2}, \dots, \frac{\partial^2}{\partial x_n^2} \right)$$

הוכחה:

לפי ההגדרה הפורמלית של מטריצת ההסיאן שלומדים בחדו"א 2, מטריצת ההסיאן היא:



$$H(\bar{x}) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{pmatrix}$$

ונבחין במקרה בו הפונקציה  $f(x)$  דיפרנציאבלית, אז לפי משפט שוורץ אפשר להחליף את סדר הגזירה ולכן מטריצת ההסיאן היא סימטרית.

כמו כן, ניתן לקבל מתוך ההגדרה שאנחנו פיתחנו כי מתקיים:

$$H(\bar{x}) = \begin{pmatrix} \frac{\partial g_1(\bar{x})}{\partial x_1} & \dots & \frac{\partial g_1(\bar{x})}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_n(\bar{x})}{\partial x_1} & \dots & \frac{\partial g_n(\bar{x})}{\partial x_n} \end{pmatrix} = \begin{pmatrix} \text{Row Vector} & \text{Row Vector} \\ \nabla \cdot \overbrace{g_1^T(\bar{x})} & \\ \vdots & \\ \nabla \cdot \overbrace{g_n^T(\bar{x})} & \end{pmatrix}$$

(כאשר  $\nabla = \left( \frac{\partial}{\partial x_1}, \dots, \frac{\partial}{\partial x_n} \right)$  ולכן  $\nabla \cdot \overbrace{g_1^T(\bar{x})}^{\text{Row Vector}} = \text{Row Vector}$  (והמכפלה היא איבר-איבר)).

והביטוי האחרון למטריצת ההסיאן נכון כי ראינו שמתקיים:  $g(\bar{x}) = (g_1(\bar{x}) \dots g_n(\bar{x}))^T = \left( \frac{\partial f}{\partial x_1} \dots \frac{\partial f}{\partial x_n} \right)^T$

כלומר  $g(\bar{x})$  הוא פונקציה וקטורית, וכל רכיב של הווקטור הוא למעשה פונקציה סקלרית של מספר משתנים (כמספר הרכיבים של ווקטור  $\bar{x}$ ). בנוסף, ואפשר לראות בקלות כי מתקיים:

$$\frac{\partial^2 f(\bar{x})}{\partial x_i \partial x_j} = \frac{\partial}{\partial x_i} \left( \frac{\partial f(\bar{x})}{\partial x_j} \right) = \frac{\partial}{\partial x_i} (g_j(\bar{x})) = \frac{\partial g_j(\bar{x})}{\partial x_i}$$

ולכן:

$$H(\bar{x}) = \begin{pmatrix} \frac{\partial^2 f(\bar{x})}{\partial x_1 \partial x_1} & \dots & \frac{\partial^2 f(\bar{x})}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\bar{x})}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f(\bar{x})}{\partial x_n \partial x_n} \end{pmatrix} = \begin{pmatrix} \frac{\partial g_1(\bar{x})}{\partial x_1} & \dots & \frac{\partial g_1(\bar{x})}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_n(\bar{x})}{\partial x_1} & \dots & \frac{\partial g_n(\bar{x})}{\partial x_n} \end{pmatrix} = \begin{pmatrix} \text{Row Vector} & \text{Row Vector} \\ \nabla \cdot \overbrace{g_1^T(\bar{x})} & \\ \vdots & \\ \nabla \cdot \overbrace{g_n^T(\bar{x})} & \end{pmatrix}$$

כמו כן, נפשט את הביטוי  $H(\bar{x}) \cdot d\bar{x}$  ונקבל:

$$H(\bar{x}) \cdot d\bar{x} = H(\bar{x}) \begin{pmatrix} dx_1 \\ \vdots \\ dx_n \end{pmatrix} = \begin{pmatrix} \overbrace{\nabla g_1^T(\bar{x}) \cdot d\bar{x}}^{\text{Row Vector}} & \overbrace{d\bar{x}}^{\text{Column Vector}} \\ \vdots \\ \nabla g_n^T(\bar{x}) \cdot d\bar{x} \end{pmatrix} = \begin{pmatrix} \overbrace{dg_1(\bar{x})}^{\text{Scalar}} \\ \vdots \\ dg_n(\bar{x}) \end{pmatrix} = d\bar{g}(\bar{x})$$

מש"ל.

הנגזרת של הנגזרת הכיוונית בכיוון כלשהו: (Second Directional Derivative)

ראינו כי אם נגדיר  $\varphi(\alpha) = f(\bar{x}_0 + \alpha \cdot \bar{r})$  אז הנגזרת הכיוונית של  $f$  בנקודה  $\bar{x}_0$  בכיוון  $\bar{r}$  מסומנת להיות  $f_{\bar{r}}(\bar{x}_0)$  והיא שווה לביטוי:

$$f_{\bar{r}}(\bar{x}_0) = \varphi'(\alpha)|_{\alpha=0} = g^T(\bar{x}_0) \cdot \bar{r}$$

(כאשר  $g^T(\bar{x}_0)$  זה הגרדיאנט (שהוא ווקטור) של הפונקציה  $f(\bar{x})$  בנקודה  $\bar{x}_0$ ).

הנגזרת הכיוונית השניה מוגדרת באופן דומה ונקבל:

$$\varphi''(\alpha)|_{\alpha=0} = f_{\bar{r}\bar{r}}^*(\bar{x}_0) = \left( \frac{\partial^2}{\partial \alpha^2} f(\bar{x}_0 + \alpha \bar{r}) \right) \Big|_{\alpha=0} = \left( \frac{\partial}{\partial \alpha} f_r' \left( \underbrace{\bar{x}_0 + \alpha \bar{r}}_{\triangleq \bar{x}} \right) \right) \Big|_{\alpha=0}$$

אנחנו יודעים כי מתקיים:

$$f_{\bar{r}}(\bar{x}_0) = g^T(\bar{x}_0) \cdot \bar{r} = \bar{r}^T \cdot g(\bar{x}_0)$$

ולכן הדיפרנציאל של הנגזרת הכיוונית הוא (נזכור כי למעשה  $g(\bar{x}_0) \equiv \bar{g}(\bar{x}_0)$  הוא גרדיאנט ולכן הוא ווקטור!):

$$df_{\bar{r}}'(\bar{x}_0) = \bar{r}^T \cdot d\bar{g}(\bar{x}_0) = \bar{r}^T \cdot H(\bar{x}_0) \cdot d\bar{x} = \bar{r}^T \cdot H(\bar{x}_0) \cdot d\alpha \bar{r} = \bar{r}^T \cdot H(\bar{x}_0) \cdot \bar{r} \cdot d\alpha$$

ולכן לפי הגדרת הנגזרת באמצעות הדיפרנציאל נקבל כי מתקיים:

$$\boxed{f_{\bar{r}\bar{r}}''(\bar{x}_0) = \frac{df_{\bar{r}}'(\bar{x}_0)}{d\alpha} = \bar{r}^T H(\bar{x}_0) \bar{r}}$$

וזו תבנית ריבועית שכבר ראינו (בהקשרים של ערכים עצמיים ווקטורים עצמיים)

דיפרנציאל של אופרטור לינארי:

נניח כי  $\bar{y} = A\bar{x}$ , כאשר  $A$  זו מטריצה. הדיפרנציאל של אופרטור לינארי הוא:

$$d\bar{y} \triangleq A(\bar{x} + d\bar{x}) - A\bar{x} = A d\bar{x}$$

גרדיאנט של פונקציה לינארית:

נניח כי  $\bar{b}$  הוא ווקטור עמודה. אם נניח כי הפונקציה הסקלרית היא  $f(\bar{x}) = \bar{b}^T \cdot \bar{x}$  אז הדיפרנציאל שלה הוא:

$$\forall d\bar{x}: df(\bar{x}) = \bar{b}^T d\bar{x}$$

ואם נשווה את הביטוי האחרון להגדרה שראינו  $df(\bar{x}) = \bar{g}^T(\bar{x}) \cdot d\bar{x}$  נקבל:

$$\bar{g}(\bar{x}) = \nabla f(\bar{x}) = \bar{b}$$

דוגמה: (מציאת גרדיאנט של פונקציה בעלת תבנית ריבועית)

נניח כי נתונה לנו הפונקציה הסקלרית (התוצאה שלה היא סקלר) של מספר משתנים:  $f(\bar{x}) = \bar{x}^T A \bar{x}$ , כאשר  $\bar{x}$  הוא ווקטור עמודה מימד  $n$  והמטריצה  $A$  היא ריבועית מימד  $n \times n$ . אפשר לרשום בצורה מפורשת את הפונקציה בצורה הבאה:

$$f(\bar{x}) = \sum_{i,j} x_i a_{i,j} x_j$$

אנחנו רוצים למצוא את הגרדיאנט של הפונקציה הנ"ל בנקודה כלשהי. ניתן למצוא את הגרדיאנט לפי ההגדרה, לגזור את הפונקציה כל פעם לפי אחד המשתנים ובסופו של דבר לקבל את הגרדיאנט אך אנחנו נמצא אותו באמצעות הגדרת הדיפרנציאל שלמדנו. לפי שיטות של מציאת דיפרנציאל, אנחנו צריכים לסכום מספר איברים כמספר ההופעות של  $\bar{x}$  בפונקציה  $f(\bar{x})$  ובכל פעם להתייחס למופע אחד של  $\bar{x}$  בתור משתנה, ואל השאר בתור קבועים (דומה לנגזרת של מכפלה). הדיפרנציאל של פונקציה  $f(\bar{x})$  הוא:

$$df(\bar{x}) = \underbrace{(d\bar{x}^T)}_{\text{Constant}} \cdot \underbrace{A\bar{x}}_{\text{Constant}} + \underbrace{\bar{x}^T A}_{\text{Scalar}} \cdot (d\bar{x})$$

נבחין כי הביטוי  $\bar{x}^T A \cdot (d\bar{x})$  הוא סקלר ולכן ניתן לעשות עליו את פעולת ה-Transpose מבלי שהוא ישנה משהו (כיוון שמתקיים  $(\forall c \in \mathbb{R}: c^T = c)$ ) ונקבל:

$$df(\bar{x}) = (d\bar{x}^T) \cdot A\bar{x} + \underbrace{\bar{x}^T A \cdot (d\bar{x})}_{\text{Scalar}} = (d\bar{x}^T) \cdot A\bar{x} + (d\bar{x}^T) A^T \cdot \bar{x} = d\bar{x}^T (A + A^T) \bar{x}$$

כעת אפשר לזהות מתוך הביטוי האחרון את הגרדיאנט, שכן ראינו לפי ההגדרה כי הגרדיאנט הוא פונקציה שתלויה ב- $\bar{x}$  ומוכפלת בדיפרנציאל כפי שראינו בהגדרה:  $df(\bar{x}) = \bar{g}^T(\bar{x}) \cdot d\bar{x} = d\bar{x}^T \cdot \bar{g}(\bar{x})$  ולכן קיבלנו כי מתקיים עבור הפונקציה  $f(\bar{x}) = \bar{x}^T A \bar{x}$  כי הגרדיאנט שלה הוא:

$$\bar{g}(\bar{x}) = (A + A^T) \bar{x}$$

וכן במקרה בו  $A$  היא סימטרית, אז מתקיים  $A = A^T$  ונקבל:

$$\bar{g}(\bar{x}) = 2A\bar{x}$$

(נבחין כי אם הפונקציה  $f(\bar{x}) = \bar{x}^T A \bar{x}$  הייתה של משתנה יחיד אז היה מתקיים  $f(x) = Ax^2$  וברור כי  $f'(x) = 2Ax$ )

כעת נמצא את מטריצת ההסיאן של הפונקציה  $f(\bar{x}) = \bar{x}^T A \bar{x}$ . לפי מה שלמדנו מתקיים כי  $d\bar{g}(\bar{x}) = H(\bar{x}) \cdot d\bar{x}$  ולכן אנחנו רוצים למצוא את הדיפרנציאל של  $\bar{g}(\bar{x}) = (A + A^T) \bar{x}$ :

$$d\bar{g}(\bar{x}) = (A + A^T) d\bar{x}$$

ולכן אנחנו מקבלים (ע"י השוואה לביטוי  $(d\bar{g}(\bar{x}) = H(\bar{x}) \cdot d\bar{x})$  :

$$H(\bar{x}) = A + A^T$$

וכן במקרה בו  $A = A^T$  היא סימטרית, אז מתקיים  $A = A^T$  ונקבל:

$$H(\bar{x}) = 2A$$

מדוגמה זו אנחנו למדים כי השיטה של מציאת הגרדיאנט ומטריצת ההסיאן של פונקציות מרובות משתנים באמצעות הגדרת הגרדיאנט והדיפרנציאל כפי שראינו אותה קודם נוחה ופשוטה יחסית. בעתיד נוכל להשתמש בשיטה זו גם עבור פונקציות מסובכות הרבה יותר, למשל פונקציות המתארות רשתות עצביות (Neural Networks).

פיתוח טור טיילור של פונקציה מרובת משתנים והקשר לגרדיאנט ולמטריצת ההסיאן:

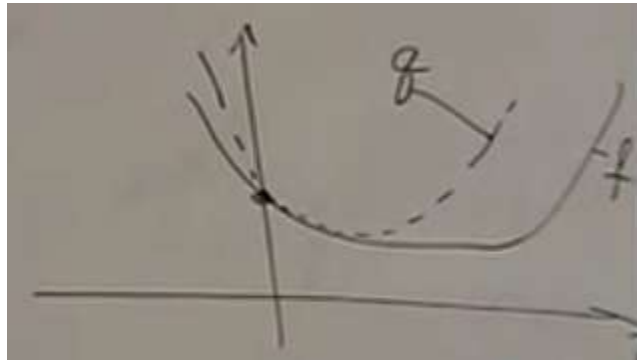
האיברים הראשונים (עד סדר שני) של פיתוח טיילור של פונקציה של מספר משתנים ניתן לרישום בצורה הבאה:

$$f(\bar{x}_0 + \bar{r}) = \underbrace{f(\bar{x}_0) + \bar{g}^T(\bar{x}_0) \cdot \bar{r} + \frac{1}{2} \bar{r}^T H(\bar{x}_0) \bar{r}}_{q(\bar{r})} + \dots$$

נסמן את האיברים הראשונים הללו להיות:

$$q(\bar{r}) = f(\bar{x}_0) + \bar{g}^T(\bar{x}_0) \cdot \bar{r} + \frac{1}{2} \bar{r}^T H(\bar{x}_0) \bar{r}$$

נתבונן על הפונקציה  $f(\bar{x}_0 + \bar{r})$  (שהיא פונקציה של  $\bar{r}$ , כי  $\bar{x}_0$  הוא קבוע), אזי זה אומר כי ציר  $x$  של הגרף הוא  $\bar{r}$ :



נבחין כי עבור  $\bar{r} = 0$  נקבל בדיוק את הנקודה  $f(\bar{x}_0)$ . כמו כן, נבחין כי  $q(\bar{r})$  הוא פיתוח טיילור עד סדר שני, ולכן הוא מנסה לקרב את פונקציה  $f(\bar{x})$  בסביבת הנקודה  $\bar{x}_0$  ולפונקציה  $q(\bar{r})$  יש את אותה הנגזרת הראשונה וגם השנייה בנקודה  $\bar{x}_0$  בדיוק כמו לפונקציה  $f(\bar{x})$ , אך כיוון שאנחנו עוסקים בפונקציות של מספר משתנים הנגזרת הראשונה היא הגרדיאנט(!), נוכל לסכם זאת כך:

$$\begin{aligned}\bar{r} = \bar{0} &\Rightarrow f(\bar{x}_0) = q(\bar{r} = \bar{0}) \\ \nabla f(\bar{x}_0) &= \bar{g}(\bar{x}_0) = \nabla q(\bar{r}) \\ \nabla^2 f(\bar{x}_0) &= \nabla^2 q(\bar{r})\end{aligned}$$

ניתן בקלות להוכיח את שלושת המשוואות לעיל ע"י הצבה ישירה של  $\bar{r} = \bar{0}$  בהגדרה של  $q(\bar{r})$  ולגזור פעם ופעמיים עבור שתי המשוואות האחרונות ולהיווכח שזה נכון, נבצע זאת:

המשוואה הראשונה היא הצבה ישירה.

המשוואה השנייה (נזכור כי מטריצת ההסיאן היא סימטרית, כיוון שאנחנו מניחים שהפונקציה  $f(\bar{x})$  היא דיפרנציאבילית ולכן מתקיים משפט שוורץ שמאפשר החלפת סדר גזירה):

$$\nabla q(\bar{r}) \Big|_{\substack{\bar{r}=0 \\ f(\bar{x}_0) \text{ is a constant}}} = \left( \bar{g}(\bar{x}_0) + \frac{1}{2} \cdot 2H(\bar{x}_0) \bar{r} \right) \Big|_{\bar{r}=0} = \bar{g}(\bar{x}_0) = \nabla f(\bar{x}_0)$$

$$\left( \frac{\partial}{\partial \bar{x}} \bar{a}^T \cdot \bar{x} = \bar{a} \text{ ובאופן כללי: } \frac{\partial}{\partial \bar{r}} \bar{g}^T(\bar{x}_0) \cdot \bar{r} = \bar{g}(\bar{x}_0) \right) \text{ נבחין כי מתקיים:}$$

המשוואה השלישית:

$$\nabla^2 q(\bar{r}) \Big|_{\substack{\bar{r}=0 \\ \bar{g}^T(x) \text{ is a constant}}} = \frac{1}{2} \cdot 2H(\bar{x}_0) = H(\bar{x}_0)$$

והוכחנו את הנדרש.

בקורס זה לא נשתמש בסדרים גבוהים יותר של טור טיילור אך בפיתוח עד סדר שני נשתמש רבות באלגוריתמים לאופטימיזציה.

פונקציות של מטריצות:

נגדיר פונקציה סקלרית  $f: n \times n \rightarrow \mathbb{R}$ , כלומר אם  $X$  היא מטריצה מסדר  $n \times n$  אז  $f(X)$  הוא סקלר.

באופן דומה לאיך שהגדרנו את הגרדיאנט של פונקציה של מספר משתנים (וקטור של משתנים), בתור וקטור של הנגזרות הראשונות של הפונקציה ביחס לכל אחד מהמשתנים, ניתן להגדיר את הגרדיאנט של פונקציה של מטריצה באמצעות מטריצת הנגזרת הראשונה של הפונקציה ביחס לכל אחד מהמשתנים באופן הבא:

$$\text{Grad } f = G(X) = \begin{pmatrix} \frac{\partial f}{\partial x_{11}} & \dots & \frac{\partial f}{\partial x_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial x_{n1}} & \dots & \frac{\partial f}{\partial x_{nn}} \end{pmatrix}$$

מכפלה פנימית/סקלרית של מטריצות:

עד היום הכרנו את ההגדרה של מכפלה סקלרית בין ווקטורים. מכפלה סקלרית/פנימית בין שתי מטריצות מוגדרת באופן הבא:

$$\langle A, B \rangle = \sum_{i,j} a_{ij} b_{ij} = \text{Trace}(A^T B)$$

הדיפרנציאל של פונקציה סקלרית של מטריצות:

נניח כי  $f: \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$  אז מתקיים כי דיפרנציאל של מטריצה  $f$  הוא:

$$df = \sum_{i,j} \frac{\partial f}{\partial x_{ij}} \cdot dx_{ij} = \text{Trace}(G^T(X) \cdot dX)$$

כאשר נגדיר:

$$dX = \begin{pmatrix} dx_{11} & \cdots & dx_{1n} \\ \vdots & \ddots & \vdots \\ dx_{n1} & \cdots & dx_{nn} \end{pmatrix}$$

אך נבחין כי לפי ההגדרה של מכפלה פנימית בין מטריצות נקבל כי מתקיים:

$$df = \sum_{i,j} \frac{\partial f}{\partial x_{ij}} \cdot dx_{ij} = \text{Trace}(G^T(X) \cdot dX) = \langle G(X), dX \rangle$$

אך זה לא דבר חדש לנו, שכן ראינו כי בדיוק אותו הדבר מתקיים עבור דיפרנציאל של פונקציה של מספר משתנים.

דוגמה: (גרדיאנט של רשת עצבים)

נזכיר כי רשת עצבים (בעלת שכבה אחת) פשוטה ניתנת לתיאור ע"י הפונקציה הסקלרית (מרבית המשתנים) הבאה:

$$f(\bar{x}, W, \bar{b}, \bar{v}) = v^T \varphi \left( \underbrace{W^T \bar{x} + \bar{b}}_{\hat{u}} \right)$$

כאשר  $\bar{\varphi}: \mathbb{R}^n \rightarrow \mathbb{R}^n$  היא פונקציה ווקטורית, כלומר:

$$\bar{\varphi}(\bar{u}) = \begin{pmatrix} \varphi(u_1) \\ \vdots \\ \varphi(u_n) \end{pmatrix} \Rightarrow d\bar{\varphi} = \begin{pmatrix} \varphi'(u_1) du_1 \\ \vdots \\ \varphi'(u_n) du_n \end{pmatrix} = \begin{pmatrix} \varphi'(u_1) & & 0 \\ & \ddots & \\ 0 & & \varphi'(u_n) \end{pmatrix} \begin{pmatrix} du_1 \\ \vdots \\ du_n \end{pmatrix}$$

האופטימיזציה שאנחנו מחפשים היא למעשה לחפש את המקדמים של מטריצה  $W$ , וכן את המקדמים של הווקטורים  $\bar{b}, \bar{v}$ . בכדי למצוא את המקדמים, אנחנו צריכים את הנגזרות ביחס לכל אחד מהמשתנים  $\bar{b}, \bar{v}, W$  אך בדוגמה זו נמצא רק את הגרדיאנט של  $f$  ביחס למטריצה  $W$ .

בשביל למצוא את הגרדיאנט של  $f$  ביחס למטריצה  $W$  נרצה למצוא את הדיפרנציאל של פונקציה  $f$  ביחס לדיפרנציאל (שינוי) באחד המשתנים, ובמקרה שלנו, אנחנו מחפשים שינוי במטריצה  $W$  (כאשר כל שאר המשתנים הם קבועים) ולכן נסמן את הדיפרנציאל של  $f$  ביחס למטריצה  $W$  להיות  $d_W f$ .

תחילה נסמן:

$$\mathcal{G}' \triangleq \begin{pmatrix} \varphi'(u_1) & & 0 \\ & \ddots & \\ 0 & & \varphi'(u_n) \end{pmatrix}$$

ולכן, לפי הסימון נקבל:

$$d\bar{\varphi} = \mathcal{G}' \cdot d\bar{u}$$

כעת אנחנו נמצא את הדיפרנציאל של  $f$  ביחס למטריצה  $W$ . אנחנו רואים לפי הסימון של  $\bar{u}$ , שהוא פונקציה של  $W$  והדיפרנציאל של  $\bar{u}$  ביחס למטריצה  $W$  הוא:

$$d_w \bar{u} = (dW^T) \bar{x}$$

ולכן:

$$d_w f = \bar{v}^T d\bar{\varphi} = \bar{v}^T \mathcal{G}' d\bar{u} = \bar{v}^T \mathcal{G}' (dW^T) \bar{x}$$

אנחנו יודעים כי לפי הגדרת הגרדיאנט, אנחנו נרצה כעת כי הביטוי  $dW^T$  יהיה מוכפל בביטוי אחר משמאל (שיהווה את הגרדיאנט של פונקציה  $f$ ), אך כעת אנחנו רואים כי הוא במרכז ולכן נרצה לבודד אותו. נבחין כי הביטוי  $\bar{v}^T \mathcal{G}' (dW^T) \bar{x}$  הוא סקלר וכן אנחנו יודעים כי לכל סקלר מתקיים  $c = \text{Trace}(c)$  ולכן נקבל:

$$d_w f = \bar{v}^T \mathcal{G}' (dW^T) \bar{x} = \text{Trace}(\bar{v}^T \mathcal{G}' (dW^T) \bar{x})$$

אך תחת  $\text{Trace}$  יש לנו את תכונת הקומוטטיביות הציקלית (כלומר  $\text{Trace}(A \cdot B) = \text{Trace}(B \cdot A)$ ) גם עבור מטריצות וגם עבור ווקטורים וגם עבור סקלרים)) ולכן נוכל לרשום:

$$d_w f = \bar{v}^T \mathcal{G}' (dW^T) \bar{x} = \text{Trace} \left( \underbrace{\bar{v}^T}_{A} \underbrace{\mathcal{G}' (dW^T) \bar{x}}_B \right) = \text{Trace} \left( \underbrace{\mathcal{G}' (dW^T) \bar{x}}_C \underbrace{\bar{v}^T}_D \right) = \text{Trace}(\bar{x} \cdot \bar{v}^T \mathcal{G}' (dW^T))$$

נשתמש בתכונת הקומוטטיביות פעם נוספת ונקבל:

$$d_w f = \bar{v}^T \mathcal{G}' (dW^T) \bar{x} = \text{Trace} \left( \underbrace{\bar{x} \cdot \bar{v}^T \mathcal{G}' (dW^T)}_A \right) = \text{Trace} \left( \underbrace{(dW^T)}_C \underbrace{\bar{x} \cdot \bar{v}^T \mathcal{G}'}_D \right) = \langle dW, \bar{x} \cdot \bar{v}^T \mathcal{G}' \rangle$$

כמו כן, ראינו כי מתקיים:

$$df = \sum_{i,j} \frac{\partial f}{\partial x_{ij}} \cdot dx_{ij} = \text{Trace}(G^T(X) \cdot dX) = \langle G(X), dX \rangle \quad \Rightarrow \quad df = \langle G(X), dX \rangle$$

ולכן נוכל להשוות בין הביטויים האחרונים ונקבל:

$$G_w(X) = \bar{x} \cdot \bar{v}^T \cdot \mathcal{G}'$$

מבוא לאופטימיזציה, 236330, הרצאות ווידאו מאביב 2010 של מיכאל ציבולבסקי, נכתב ע"י רמי נודלמן, אביב 2015  
עמוד 24

כמו כן ניתן באופן דומה לחשב את הגרדיאנטים ביחס למשתנים האחרים  $\bar{b}, \bar{v}$ , כלומר לחשב את הביטויים  $\nabla_{\bar{b}} f, \nabla_{\bar{v}} f$ .



## הרצאה 04-05, Convex Sets and Functions, קבוצות קמורות ופונקציות

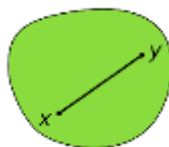
### קמורות

הגדרה: (קבוצה קמורה)

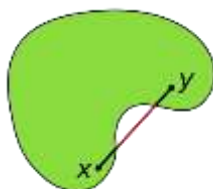
קבוצה אשר הקו הישר בין כל שתי נקודות בקבוצה נמצא כולו בקבוצה. ובאופן פורמלי, קבוצה  $C$  תקרא קבוצה קמורה, אם לכל שתי נקודות  $x, y \in C$  מתקיים:

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

דוגמה לקבוצה קמורה:



דוגמה לקבוצה לא קמורה:



הגדרה: (פונקציה קמורה)

פונקציה  $f(x)$  (של משתנה יחיד) תקרא פונקציה קמורה, אם לכל שתי נקודות  $x_1, x_2$ , הקו הלינארי שמחבר בין  $f(x_1), f(x_2)$  יהיה מעל (גדול או שווה) הפונקציה  $f(x)$  בתחום  $[x_1, x_2]$  ובצורה פורמלית:

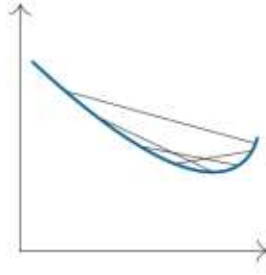
הפונקציה  $f(x)$  היא פונקציה קמורה היא לכל  $x_1, x_2$  מתקיים:

$$f(\alpha x_1 + (1 - \alpha)x_2) \leq \alpha f(x_1) + (1 - \alpha)f(x_2), \quad \forall \alpha \in [0, 1]$$

(הערה: פונקציה תקרא קמורה ממש, תדרוש אי שוויון חזק בהגדרה, במקום אי-שוויון חלש)

(הערה: פונקציה קמורה יכולה להיות גם פונקציה של מספר משתנים, וההגדרה המילולית משתנה להיות במקום קו ישר, אלא הישר שעובר על המרחק בין ערכי הפונקציה וכן צריך לשפר את ההגדרה הפורמלית)

דוגמה לפונקציה קמורה:



הגדרה: (פונקציה קעורה)

פונקציה  $f(x)$  (של משתנה יחיד) תקרא פונקציה קעורה אם  $(-f(x))$  היא פונקציה קמורה.

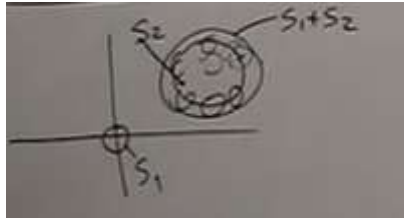
(הערה: פונקציה תקרא קעורה ממש, אם  $(-f(x))$  היא קמורה ממש)

תכונות של קבוצות קמורות: (להוכיח כשיעורי בית את התכונות)

נניח כי הקבוצות  $C_i$  הן קבוצות קמורות. מתקיימות התכונות הבאות:

(1) הקבוצה  $\bigcap_i C_i$  היא קבוצה קמורה.

(2) הקבוצה  $C_i + C_j = \{x_p + x_k \mid x_p \in C_i, x_k \in C_j\}$  היא קבוצה קמורה.



(אפשר להבין זאת אם נחשוב על נקודות כעל ווקטורים, ולכן הקבוצה  $C_i + C_j$  היא הקבוצה של כל הנקודות מתוך  $C_i$  כאשר אנחנו מוסיפים לכל נקודה כזו את כל הנקודות מתוך  $C_j$ , ואפשר להתייחס לכל הנקודות בתור ווקטורים, ולכן למעשה קיבלנו כי סביב כל נקודה בקבוצה  $C_i$  קיבלנו את הקבוצה  $C_j$ ).

(3) הקבוצה  $\{A\bar{x} \mid \bar{x} \in C_i, A: n \times n\}$  הקבוצה שנוצרת ע"י העתקה (מטריצה מייצגת) היא קבוצה קמורה.

(4) קבוצות שמייצגות את כל התחום הפנימי של ערך זהה של פונקציה קמורה (פונקצית של מספר משתנים), ובצורה פורמלית:  $\{\bar{x} \in C_i \mid f(\bar{x}): C \rightarrow \mathbb{R}, f(\bar{x}) \leq \alpha\}$ .

הוכחה לתכונה (1):

נתבונן בשתי נקודות כלשהן  $x, y \in \bigcap_i C_i$ . אם  $x, y \in \bigcap_i C_i$  אז לפי ההגדרה של חיתוך, שתי הנקודות שייכות לכל קבוצה קמורה  $C_i$  ולכן הקו הישר העובר בין שתי הנקודות שייך לכל קבוצה  $C_i$ , אך זה אומר, לפי ההגדרה, שכל קו כזה שייך גם

לחיתוך  $\bigcap_i C_i$  ולכן הוכחנו כי לכל שתי נקודות  $x, y \in \bigcap_i C_i$  הקו הישר המחבר בין הנקודות נמצא בקבוצה  $\bigcap_i C_i$  ולכן

$\bigcap_i C_i$  הוא קבוצה קמורה לפי ההגדרה. מש"ל.

תכונות של פונקציות קמורות: (להוכיח כשיעורי בית)

נניח כי  $f(\bar{x}), g(\bar{x})$  פונקציות קמורות.

(1) נגדיר את הפונקציה  $\varphi: \mathbb{R} \rightarrow \mathbb{R}$  בצורה הבאה:  $\varphi(\alpha) = f(\bar{x} + \alpha \bar{r})$  לכל  $\alpha \in \mathbb{R}$  ולכל  $\bar{x}, \bar{r} \in \mathbb{R}^n$ , אזי

הפונקציה  $\varphi(\alpha)$  היא פונקציה קמורה.

(2) (הרחבה לתכונה 1) אם לכל  $\alpha \in \mathbb{R}$  ולכל  $\bar{x}, \bar{r} \in \mathbb{R}^n$  ולכל פונקציה  $\varphi: \mathbb{R} \rightarrow \mathbb{R}$  שמוגדרת בצורה הבאה:

$\varphi(\alpha) = f(\bar{x} + \alpha \bar{r})$ , הפונקציה  $\varphi(\alpha)$  היא פונקציה קמורה אז גם הפונקציה  $f(\bar{x})$  היא פונקציה קמורה.

(3) ("לינאריות") הפונקציה  $v(\bar{x}) = \alpha f(\bar{x}) + \beta g(\bar{x})$  היא פונקציה קמורה לכל  $\alpha, \beta > 0$ .

(4) הפונקציה  $m(\bar{x}) = \max \{f(\bar{x}), g(\bar{x})\}$  היא פונקציה קמורה.

(5) נניח כי  $h: \mathbb{R} \rightarrow \mathbb{R}$  פונקציה קמורה ומונוטונית עולה אז הפונקציה  $h(f(\bar{x}))$  היא פונקציה קמורה.

פונקציה קמורה מוכללת (Extended value convex function):

נניח כי  $C \in \mathbb{R}^n$  קבוצה קמורה וכן  $f(\bar{x}): \mathbb{R}^n \rightarrow \mathbb{R}$  היא פונקציה קמורה. הפונקציה הקמורה המוכללת מוגדרת להיות

פונקציה  $\bar{f}: \mathbb{R}^n \rightarrow \mathbb{R}$  ומוגדרת באופן הבא:

$$\bar{f}(\bar{x}) = \begin{cases} f(\bar{x}) & , \quad x \in C \\ \infty & , \quad x \notin C \end{cases}$$

וההגדרה היא:

$$\bar{f}(\alpha \bar{x}_1 + (1-\alpha) \bar{x}_2) \leq \alpha f(\bar{x}_1) + (1-\alpha) f(\bar{x}_2)$$

או במובן שקול (ע"י אלגברה פשוטה):

$$\bar{f}(\bar{x}_2 + \alpha(\bar{x}_1 - \bar{x}_2)) \leq f(\bar{x}_2) + \alpha(f(\bar{x}_1) - f(\bar{x}_2))$$

הערה: נבחין כי גם אם הפונקציה המקורית  $f(\bar{x})$  היא לא פונקציה קמורה (כי למשל היא מורכבת רק משני תחומים שבהם היא

קמורה), אז לפי ההגדרה לעיל, עדיין אנחנו נקבל כי הפונקציה המוכללת היא קמורה למרות שניתן לחשוב בטעות כי הקו הנראה בשרטוט (למטה) הוא הקו הלינארי, אך זה לא נכון כי הקו הלינארי הוא בעל ערך אינסופי לפי ההגדרה לעיל:



הערה: הגדרה זו נוחה כאשר אנחנו רוצים להתייחס לפונקציה שקמורה רק בתחום מסויים שלה (אך לא בכל בתחום) ולכן נוהג להתייחס לפונקציה הקמורה המוכללת המוגדרת על ידה, כי התחום של הפונקציה המוכללת הקמורה שלה היא כל  $\mathbb{R}^n$ .

הגדרה: (אפיגרף (epigraph))

אפיגרף של פונקציה  $f(\bar{x})$ , הוא קבוצת הנקודות שנמצאות על או מעל הגרף  $f(\bar{x})$ . ובאופן פורמלי: נניח כי  $C \in \mathbb{R}^n$  היא קבוצה קמורה, וכן  $f(\bar{x}): C \rightarrow \mathbb{R}$  אז מתקיים:

$$epi(f) = \{(\bar{x}, y) | x \in C, y \in \mathbb{R}, y \geq f(\bar{x})\}$$

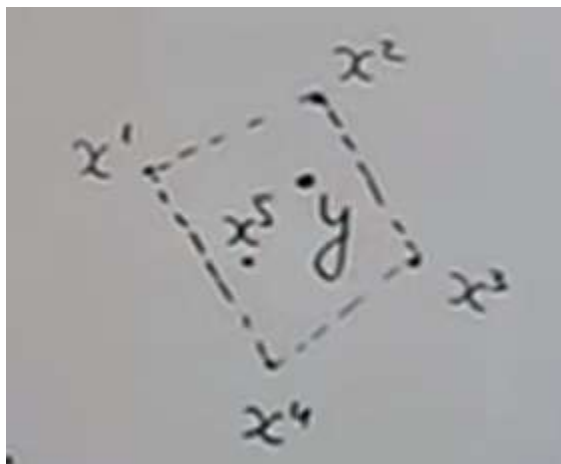
משפט:

$f(\bar{x})$  היא פונקציה קמורה אם"מ  $epi(f)$  היא קבוצה קמורה. (הוכחה: שיעורי בית)

הגדרה: (קומבינציה קמורה (convex combination))

תהי קבוצה  $C \in \mathbb{R}^n$  ויהיו  $\bar{x}_i \in \mathbb{R}^n$  נקודות וכן  $\alpha_i \in \mathbb{R}$  לכל  $1 \leq i \leq \infty$  כאשר  $\forall i: \alpha_i \geq 0$  וכן  $\sum_i \alpha_i = 1$  אז אם

הנקודה  $y = \sum_{i=1}^m \alpha_i \bar{x}_i$  (כאשר  $1 \leq m \leq \infty$ ) נמצאת בתוך הקבוצה  $C$ , לכל הנקודות  $\bar{x}_i \in \mathbb{R}^n$  ולכל  $\alpha_i \in \mathbb{R}$  אז הקבוצה  $C$  היא קומבינציה קמורה.



(הנקודה  $y = \sum_{i=1}^m \alpha_i \bar{x}_i$  נמצאת בתוך התחום שנוצר ע"י הנקודות  $\bar{x}_i \in \mathbb{R}^n$ ).

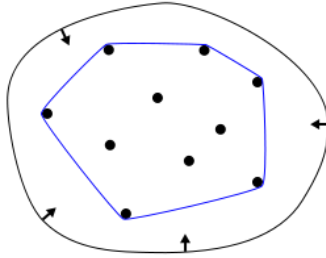
הגדרה: (קמור, convex hull)

הקמור של גוף או אוסף של גופים הוא הגוף הקמור המינימלי המכיל אותם. הקמור של נקודות במישור, הוא סוג של גומיה שנמתחת מעל כל הנקודות המישור ונמתחה כך שתכיל את כל הנקודות ולאחר מכן שוחררה. סימון:

$$conv\{\bar{x}_i\}_{i=1}^m$$

(הסבר לסימון: הקבוצה הקמורה המינימלית שמכילה את כל הנקודות  $\{\bar{x}_i\}_{i=1}^m$ )

דוגמה לקמור:



הגדרה שקולה לקמור:

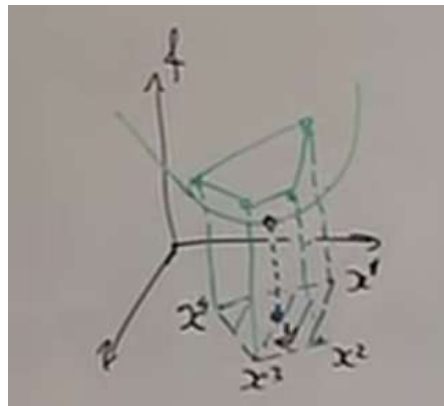
הקבוצה הקמורה המינימלית (בשטח/נפח וכו') שמכילה את כל הנקודות  $\{\bar{x}_i\}_{i=1}^m$ .

אי-שיוויון ינסן: (נלמד בקורס הסתברות ח')

נניח כי  $f(\bar{x})$  פונקציה קמורה (של  $n$  משתנים), וכן  $\bar{x}_i \in \mathbb{R}^n$  נקודות בתחום, אז לכל  $1 \leq i \leq n$ ,  $\alpha_i \geq 0$  וכן  $\sum_{i=1}^n \alpha_i = 1$  מתקיים:

$$f\left(\sum_{i=1}^n \alpha_i \bar{x}_i\right) \leq \sum_{i=1}^n \alpha_i f(\bar{x}_i)$$

ולשם המחשה:



ערך הפונקציה בנקודה  $\sum_{i=1}^n \alpha_i \bar{x}_i$  הוא בהכרח קטן או שווה לקומבינציה הלינארית של הערכים של הפונקציה בנקודות  $\bar{x}_i$ .

דוגמה: (דוגמה לשימוש באי-שיוויון ינסן)

נוכיח כי הממוצע החשבוני גדול או שווה לממוצע הגיאומטרי של קבוצת נקודות (חלק ממשפט אי-שיוויון הממוצעים), כלומר נוכיח כי:

$$\frac{1}{m} \sum_{i=1}^m x_i \geq \left( \prod_{i=1}^m x_i \right)^{1/m}$$

הוכחה:

נזכור כי הפונקציה  $\log(\cdot)$  היא פונקציה קעורה, ולכן ע"י הפעלת פונקצית  $(-\log(\cdot))$  (מינוס לוג שהיא פונקציה קמורה לפי ההגדרה) על שני האגפים ופישוט אי-השוויון נוכל להשתמש באי-שוויון ינסן ולקבל את מה שנדרשנו להוכיח.

נגזרות מסדר ראשון ושני (גרדיאנטים) של פונקציות קמורות:

אי-שוויון הגרדיאנט:

הפונקציה  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  היא פונקציה קמורה ו"חלקה" אם לכל  $\bar{d} \in \mathbb{R}^n$  מתקיים

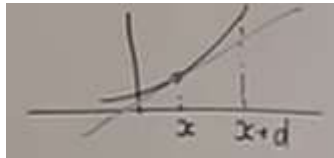
$$f(\bar{x} + \bar{d}) \geq f(\bar{x}) + \bar{d}^T \nabla f(\bar{x})$$

או באופן שקול:

$$f(\bar{x} + \bar{d}) - f(\bar{x}) \geq \bar{d}^T \nabla f(\bar{x})$$

(הערה:  $f(\bar{x}) + \bar{d}^T \nabla f(\bar{x})$  הוא פיתוח טיילור מסדר ראשון של  $f(\bar{x})$ )

המחשה של אי-שוויון הגרדיאנט עבור פונקציה של משתנה יחיד:



אנחנו רואים כי המשיק (קו משיק/משטח משיק) נמצא על או מתחת לגרף הפונקציה בנקודה אליה הוא משיק.

נגזרת שנייה של פונקציה קמורה (במקרה של פונקציה של משתנה יחיד):

אם  $f(x)$  היא פונקציה קמורה אז לכל  $x_1 < x_2$  מתקיים  $f'(x_1) \leq f'(x_2)$  וכן לכל  $x$  מתקיים  $f''(x) > 0$ .

נגזרת שנייה של פונקציה קמורה (במקרה של פונקציה של כמה משתנים):

אם  $f(\bar{x})$  היא פונקציה קמורה אז מטריצת ההסיאן של פונקציה  $f(\bar{x})$  מקיימת  $H(\bar{x}) \succeq 0$  (positive semi-definite).

הוכחה:

נתון כי  $f(\bar{x})$  היא פונקציה קמורה, ולכן לכל פונקציה  $\varphi: \mathbb{R}^n \rightarrow \mathbb{R}$  שנגדיר באופן הבא:  $\varphi(\alpha) = f(\bar{x} + \alpha \bar{r})$  (לכל

$\bar{x}, \bar{r} \in \mathbb{R}^n$  וכן לכל  $\alpha \in \mathbb{R}$ ) הפונקציה  $\varphi$  תהיה קמורה ולכן מתקיים לכל  $\alpha$  מתקיים:  $\varphi''(\alpha) \geq 0$ . כמו כן, ראינו כי

הנגזרת הכיוונית, מסדר שני, של הפונקציה  $f(\bar{x})$ , אשר מוגדרת להיות  $f''_{\bar{r}}(\bar{x})$  (והיא גם שווה לפי ההגדרה ל-  $\varphi''(\alpha)|_{\alpha=0}$ )

(ולכן קיבלנו:

$$\forall \bar{r} \in \mathbb{R}^n : \quad \varphi''(\alpha)|_{\alpha=0} = f''_{\bar{r}}(\bar{x}) = \bar{r}^T \cdot \nabla^2 f(\bar{x}) \cdot \bar{r} = \bar{r}^T H(\bar{x}) \bar{r} \geq 0$$

מבוא לאופטימיזציה, 236330, הרצאות ווידאו מאביב 2010 של מיכאל ציבולבסקי, נכתב ע"י רמי נודלמן, אביב 2015  
עמוד 31

ולכן, לפי ההגדרה של מטריצה אי-שלילית (ההגדרה אומרת כי המטריצה  $A$  תהיה אי-שלילית אם לכל ווקטור  $\forall \bar{r} \in \mathbb{R}^n$  מתקיים  $\bar{r}^T \cdot A \cdot \bar{r} \geq 0$ ) קיבלנו כי מתקיים:  $H(\bar{x}) \succeq 0$ . מש"ל.

הערה: אם  $f(\bar{x})$  פונקציה קמורה ממש אז מתקיים  $H(\bar{x}) \succ 0$ .

## הרצאה 06 – Local and Global Minimum , מינימום מקומי וגלובלי

הגדרה: (מינימום גלובלי)

נניח כי  $C \in \mathbb{R}^n$  תחום (קבוצה) קמור וכן פונקציה  $f: C \rightarrow \mathbb{R}$ . נאמר כי  $\bar{x} \in C$  הוא מינימום גלובלי של הפונקציה  $f$  מעל הקבוצה  $C$  אם לכל  $\bar{y} \in C$  מתקיים:

$$f(\bar{x}) \leq f(\bar{y})$$

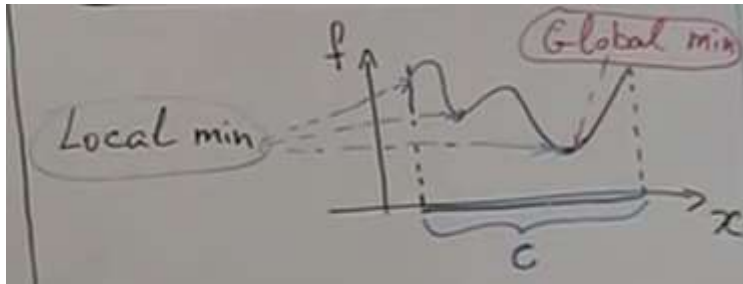
הגדרה: (מינימום מקומי, מינימום לוקאלי)

נניח כי  $C \in \mathbb{R}^n$  תחום (קבוצה) קמור וכן פונקציה  $f: C \rightarrow \mathbb{R}$ . נאמר כי  $\bar{x}$  הוא מינימום מקומי (לוקאלי) אם קיים  $\varepsilon > 0$  כך שמתקיים לכל  $\bar{y} \in B_\varepsilon(\bar{x}) \cap C$ :

$$f(\bar{x}) \leq f(\bar{y})$$

(כאשר  $B_\varepsilon(\bar{x})$  הוא סביבה ברדיוס  $\varepsilon$  של הנקודה  $\bar{x}$ , כלומר  $B_\varepsilon(\bar{x}) = \{\bar{z} \mid \|\bar{z} - \bar{x}\| \leq \varepsilon\}$ )

דוגמה:



נבחין כי גם הנקודה השמאלית ביותר של תחום  $C$  היא מינימום מקומי.

משפט:

נניח  $C \in \mathbb{R}^n$  תחום קמור. לכל פונקציה קמורה  $f: C \rightarrow \mathbb{R}$  מתקיים כי כל מינימום מקומי הוא מינימום גלובלי.

הוכחה: (על דרך השלילה)

נתון כי  $C \in \mathbb{R}^n$  תחום קמור וכן נתונה פונקציה קמורה  $f: C \rightarrow \mathbb{R}$ . נניח בשלילה כי  $\bar{x}_0$  הוא מינימום מקומי אך לא מינימום גלובלי, אם כך, אז קיים מינימום גלובלי אחר, כלומר קיים  $\bar{y}_0 \in C$  כך שמתקיים:  $f(\bar{y}_0) < f(\bar{x}_0)$ . לפי ההגדרה של פונקציה קמורה, מתקיים (ע"י פיתוח פשוט של ההגדרה של פונקציה קמורה שראינו קודם – למעשה מחליפים תפקידים בין שתי הנקודות  $\bar{x}_0$  ו- $\bar{y}_0$ ):

$$\forall \alpha \in [0,1]: f(\bar{x}_0 + \alpha(\bar{y}_0 - \bar{x}_0)) \leq f(\bar{x}_0) + \alpha(f(\bar{y}_0) - f(\bar{x}_0))$$

כך מצאנו כי  $f(\bar{y}_0) < f(\bar{x}_0)$  ולכן מתקיים  $f(\bar{y}_0) - f(\bar{x}_0) < 0$  ולכן:



$$f(\bar{x}_0 + \alpha(\bar{y}_0 - \bar{x}_0)) \leq f(\bar{x}_0) + \alpha(f(\bar{y}_0) - f(\bar{x}_0)) < f(\bar{x}_0) \Rightarrow f(\bar{x}_0 + \alpha(\bar{y}_0 - \bar{x}_0)) < f(\bar{x}_0)$$

כלומר מצאנו סביבה של נקודות מסביב לנקודה  $\bar{x}_0$  עבורן  $f(\bar{x}_0 + \alpha(\bar{y}_0 - \bar{x}_0)) < f(\bar{x}_0)$  ולכן  $\bar{x}_0$  היא כלל לא נקודת מינימום מקומי, בסתירה לנתון. לכן ההנחה כי  $\bar{x}_0$  היא לא מינימום גלובלי שגויה ולכן הוכחנו את הנדרש. מש"ל.

משפט:

נניח  $C \in \mathbb{R}^n$  תחום קמור. לכל פונקציה קמורה ממש  $f: C \rightarrow \mathbb{R}$  יש לכל היותר מינימום גלובלי אחד בלבד.

הוכחה: (על דרך השלילה)

נניח  $C \in \mathbb{R}^n$  תחום קמור. נתונה פונקציה קמורה ממש  $f: C \rightarrow \mathbb{R}$  ונניח כי לפונקציה יש לפחות שתי נקודות שנשמנו להיות  $\bar{x}_0$  ו- $\bar{y}_0$  מינימום גלובליות. אם שתי הנקודות הן מינימום גלובלי אז ערך הפונקציה בנקודה בהכרח זהה ונסמנו להיות:  $f(\bar{x}_0) = f(\bar{y}_0) = \mu$ . נתון כי הפונקציה קמורה ממש ולכן לפי ההגדרה של פונקציה קמורה ממש ועבור  $\alpha = \frac{1}{2}$  נקבל:

$$f\left(\frac{\bar{x}_0 + \bar{y}_0}{2}\right) < \frac{f(\bar{x}_0) + f(\bar{y}_0)}{2} = \frac{\mu + \mu}{2} = \mu$$

ולכן מצאנו נקודה  $\frac{\bar{x}_0 + \bar{y}_0}{2} \in C$  אשר עברה מקבלים ערך נמוך יותר מהערך של המינימום הגלובלי בסתירה להגדרת מינימום גלובלי ולכן לכל פונקציה קמורה ממש קיים לכל היותר מינימום גלובלי אחד.

דוגמה: (פונקציה קמורה ממש ללא מינימום גלובלי)

בתחום  $C = \mathbb{R}^+ \setminus \{0\}$  לפונקציה  $f(x) = \frac{1}{x}$  אין מינימום גלובלי.

תנאים לאופטימליות (של פתרון, כלומר מציאת נקודה בתחום של פונקציה עברה נקבל את ערכה המינימלי של הפונקציה):

משפט:

נניח כי  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  פונקציה קמורה. מתקיים:

$$\nabla f(\bar{x}^*) = 0 \Leftrightarrow \bar{x}^* \in \mathbb{R}^n \text{ היא מינימום גלובלי ב- } \mathbb{R}^n$$

הוכחה:

כיוון  $\Rightarrow$ :

לפי אי-שוויון הגרדיאנט, מתקיים:

$$f(\bar{x}^* + \bar{r}) \geq f(\bar{x}^*) + \bar{r}^T \nabla f(\bar{x}^*)$$

אבל נתון לנו כי  $\nabla f(\bar{x}^*) = 0$  ולכן מתקיים:

$$\forall \bar{r} \in \mathbb{R}^n: f(\bar{x}^* + \bar{r}) \geq f(\bar{x}^*) + \bar{r}^T \nabla f(\bar{x}^*) = f(\bar{x}^*)$$

ולכן  $\bar{x}^*$  היא מינימום גלובלי לפי ההגדרה של מינימום גלובלי.

כיוון  $\Leftarrow$ :

נתון לנו כי  $\bar{x}^*$  הוא מינימום גלובלי ומחדו"א אנחנו יודעים כי הגרדיאנט של כל פונקציה בנקודת מינימום (גם גלובלי וגם מקומי) מתאפס, כלומר  $\nabla f(\bar{x}^*) = 0$ .

מש"ל.

הערה: תנאים הכרחיים עבור פונקציות שאינן קמורות.

(1) עבור פונקציות שאינן קמורות, התנאי  $\nabla f(\bar{x}^*) = 0$  הוא הכרחי, אך לא מספיק! כלומר:

הנקודה  $\bar{x}^* \in \mathbb{R}^n$  היא מינימום גלובלי ב- $\mathbb{R}^n$   $\nabla f(\bar{x}^*) = 0$   $\not\Rightarrow$

למשל עבור הפונקציה  $f(x) = x^3$ , בנקודה אפס הגרדיאנט מתאפס אבל ברור כי הנקודה היא לא מינימום גלובלי.

(2) עבור פונקציות שאינן קמורות, תנאי הכרחי נוסף שחייב להתקיים בכדי שנוכל לומר שנקודה מסוימת היא מינימום מקומי, חייבים לדרוש שעבור כל  $\bar{r} \in \mathbb{R}^n$  נקבל כי הנגזרת הכיוונית מסדר שני בנקודה ה"חשודה כמינימום", תהיה גדולה או שווה לאפס, כלומר  $f''_{\bar{r}}(\bar{x}^*) = \bar{r}^T H(\bar{x}^*) \bar{r} \geq 0$  וזו דרישה שקולה למעשה לדרישה על מטריצת ההסיאן להיות אי-שלילית:  $H(\bar{x}^*) \succeq 0$  (נבחין כי היותה של מטריצת ההסיאן אי-שלילית זה לא תנאי מספיק להיותה של הפונקציה  $f(\bar{x})$  להיות קמורה).

משפט:

נניח כי  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  פונקציה (לא בהכרח קמורה). מתקיים:

$$f \text{ הנקודה } \bar{x}^* \in \mathbb{R}^n \text{ היא מינימום מקומי ב- } \mathbb{R}^n \text{ של הפונקציה } f \Leftarrow \begin{cases} \nabla f(\bar{x}^*) = 0 \\ H(\bar{x}^*) = \nabla^2 f(\bar{x}^*) \succ 0 \end{cases}$$

(הערה: במקום התנאי  $H(\bar{x}^*) = \nabla^2 f(\bar{x}^*) \succ 0$  אפשר לדרוש במקום זו שהפונקציה  $f$  תהיה קמורה בסביבה  $B_\varepsilon(\bar{x}^*)$ )

רעיון ההוכחה: (שיעורי בית)

אם נתון לנו כי מטריצת ההסיאן חיובית בנקודה  $\bar{x}^*$  אז גם קיימת סביבה של הנקודה  $\bar{x}^*$  בה מטריצת ההסיאן היא חיובית, ולכן בסביבה זו אפשר להגדיר פונקציה קמורה ועבור פונקציה קמורה התנאי  $\nabla f(\bar{x}^*) = 0$  הוא תנאי מספיק להיותה של הנקודה  $\bar{x}^*$  מינימום גלובלי (בסביבה של הנקודה  $\bar{x}^*$ ) שלה, אבל הסביבה כולה של הפונקציה הקמורה היא סביבה מקומית של הפונקציה  $f$  ולכן הוכחנו כי  $\bar{x}^*$  הוא מינימום מקומי של הפונקציה  $f$ .

בכך סיימנו את ההכנות המתמטיות הנדרשות לקורס ואפשר להתחיל ללמוד אלגוריתמים של אופטימיזציה.

## **הרצאה 07 – Iterative Methods of One Dimensional Optimization** **אלגוריתמים איטרטיביים (נומריים) של אופטימיזציה במשתנה יחיד**

(החל מדקה 13:15)

לעיתים, ניתן לפתור בעיות אופטימיזציה באופן אנליטי ולהגיע לפתרון אנליטי. פתרון אנליטי טוב כאשר יש לנו מודל מדויק שמתאר את הבעיה, או משוואות מדויקות שאפשר לנתח. ברוב המקרים, המצב הוא לא כזה, ולכן אנחנו נאלצים לבצע חישובים נומריים שלעיתים מקרבים אותנו לפתרון באופן מיטבי ולפעמים פחות. אנחנו נתחיל בבעיה אופטימיזציה של מציאת מינימום של פונקציה של משתנה יחיד.

אלגוריתם Bisection:

נניח כי יש לנו פונקציה  $f(x)$  (משתנה יחיד) אשר יש לה מינימום מקומי (או גלובלי) בנקודה  $x^*$  (ונניח כי הפונקציה קמורה בסביבת הנקודה  $x^*$ ). האלגוריתם מתבסס על חיפוש נקודה בה הנגזרת של הפונקציה שאנחנו מחפשים את המינימום שלה, מתאפסת. אם אנחנו מניחים כי  $x^*$  יש מינימום אז משמאל, הנגזרת  $f'(x)$ , שלילית ממש ומימין לנקודה  $x^*$  הנגזרת  $f'(x)$  חיובית ממש ואילו בנקודה  $x^*$ , הנגזרת מתאפסת ממש. נסתכל על הקטע  $[a, b]$  אשר נקודה  $x^*$  נמצאת בתוכו, כלומר  $a < x^* < b$  וכעת אנחנו מחפשים את השורש של פונקציית הנגזרת, כלומר המרנו את בעית האופטימיזציה שלנו מבעית חיפוש מינימום של פונקציה לחיפוש שורש של פונקציה (אחרת מהפונקציה המקורית). האלגוריתם למציאת השורש, הוא פשוט:

- 1) בחר מספר איטרציות התחלתי  $N$  (או בחר אורך קטע  $[a, b]$  שבו תעצור) ובכל איטרציה בצע:
  - a. אם  $a = b$  החזר את הנקודה  $a$  וסיים את האלגוריתם (או אם אורך הקטע הוא קטן או שווה למה שבחרת).
  - b. בחר את הנקודה האמצעית מבין נקודות הקצה  $[a, b]$ , כלומר בחר את נקודה  $\tau = \frac{a+b}{2}$ .
  - c. אם ערך הנגזרת בנקודה  $\tau$  היא חיובית, החלף את הקטע  $[a, b]$  בקטע  $[a, \tau]$  וחזור לשלב (a).
  - d. אחרת, אם הנגזרת בנקודה  $\tau$  היא שלילית, החלף את הקטע  $[a, b]$  בקטע  $[\tau, b]$  וחזור לשלב (a).

סיבוכיות אלגוריתם Bisection:

אם נגדיר את אורך הקטע ההתחלתי שבו אנחנו מחפשים את המינימום המקומי להיות  $|b - a| = \Delta$  אז גודל הקטע  $\Delta$  קטן בצורה אקספוננציאלית ככל שנבצע יותר איטרציות כיוון שבכל איטרציה אנחנו מקטינים את הקטע פי 2, כלומר גודל הקטע לאחר  $i$  איטרציות יהיה:  $\frac{\Delta}{2^i}$ .

קשיים במימוש האלגוריתם Bisection:

בכדי להשתמש באלגוריתם Bisection אנחנו צריכים לדעת את הקטע ההתחלתי וכן את הנגזרת של הפונקציה לה אנחנו מחפשים פתרון אופטימלי, אך לא תמיד יש בידינו את הנגזרת או את היכולת לחשב אותה ולכן אלגוריתם זה מוגבל יחסית.

אלגוריתם חיתוך הזהב, Golden section:

נניח כי יש בידינו קטע  $[a, b]$  אשר אנחנו יודעים כי נקודת המינימום מוכלת בקטע הזה וכן אנחנו יודעים כי הפונקציה שלה אנחנו מחפשים את המינימום היא קמורה בקטע הזה. בנוסף, נניח כי יש לנו נקודה  $a < c < b$  עבורה גילינו כי

האם יש בידינו מספיק מידע על מנת לקבוע שנקודת המינימום נמצאת בקטע  $[a, c]$  או בקטע  $[c, a]$  בוודאות? לא! למשל:



כלומר, בחירת נקודה אחת בתוך הקטע  $[a, b]$  לא מספיקה, אך מתברר כי בחירת שתי נקודות באופן מסוים, מאפשרת לנו כן לדעת באיזה קטע נמצאת נקודת המינימום, הבחירה נעשית למשל באמצעות אלגוריתם Golden section.

אופן פעולת האלגוריתם: (נניח כי יש לנו קטע התחלתי  $[a, d]$  והמינימום מוכל בקטע זה)

(1) אם  $a = d$  עצור והחזר את נקודה  $a$ .

(2) אחרת, מצא בקטע  $[a, d]$  עוד שתי נקודות  $a < b < c < d$  וחשב את ערך הפונקציה בנקודות הנוספות.

(3) אם מתקיים  $f(b) < f(c)$  החלף את הקטע  $[a, d]$  בקטע  $[a, c]$  וחזור לשלב (1).

(4) אחרת, כלומר אם  $f(b) \geq f(c)$ , החלף את הקטע  $[a, d]$  בקטע  $[b, d]$  וחזור לשלב (1).

נבחין כי בכל פעם שאנחנו מבצעים את שלב (2) אנחנו צריכים לחשב את ערך הפונקציה בעוד שתי נקודות. מתברר שאין צורך בכך, כלומר באיטרציה הראשונה אכן נחשב את הערך של שתי הנקודות  $f(b), f(c)$  אבל בכל האיטרציות הבאות מספיק לחשב את ערך הפונקציה בנקודה נוספת אחת (במקום בשתיים) כיוון שבשלב הקודם של האלגוריתם כבר חישבנו את אחת הנקודות  $f(b), f(c)$  ואנחנו יכולים להשתמש בערכים אלו.

בנוסף, נחליט שאנחנו לא בוחרים באופן רנדומלי שתי נקודות  $b, c$  בכל קטע בשלב (2) אלא נבחר את הנקודות באופן הבא:

(מופיע שרטוט המחשה בסוף ההסבר)

נחליט על פרמטר חלוקה  $\tau \in \mathbb{R}$  קבוע אשר נמצא עליו תנאים נוספים בהמשך.

נסמן את אורך הקטע  $[a, d]$  באיטרציה האפס להיות  $\Delta_0$ . נחלק את הקטע  $[a, d]$  באופן כזה כך שנקבל:

$$[a, d] = [a, b] \cup [b, c] \cup [c, d]$$

וכן אורך כל אחד מהקטעים  $[a, b], [c, d]$  הוא בדיוק  $\tau \cdot \Delta_0$ .

לאחר מכן, באיטרציה הבאה, אם למשל בחרנו להמשיך בקטע  $[a, c]$  נחלק את הקטע  $[a, c]$  באופן כזה, כך שאורך הקטע  $[b, c]$  (כאשר  $b$  היא הנקודה שבחרנו באיטרציה הקודמת) הוא  $\tau \cdot \Delta_1$  כאשר  $\Delta_1 = |c - a|$ , כלומר  $\Delta_1$  הוא אורך הקטע של האיטרציה הראשונה (לאחר האיטרציה האפס).

מכאן, אנחנו מסיקים כי התנא לבחירת  $\tau \in \mathbb{R}$  מתאים הוא:

$$\tau \cdot \Delta_0 = (1 - \tau) \Delta_1$$

אך אפשר לבטא את  $\Delta_1$  באמצעות  $\Delta_0$  באופן הבא:  $\Delta_1 = (1-\tau)\Delta_0$

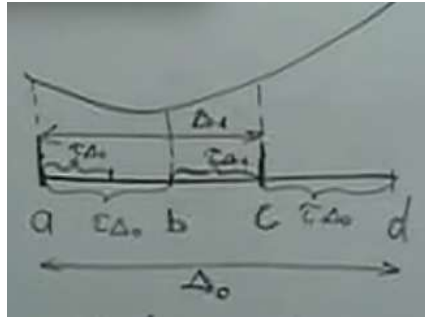
ולכן קיבלנו:

$$\tau \cdot \Delta_0 = (1-\tau)\Delta_1 = (1-\tau)^2 \Delta_0 \Rightarrow \tau = (1-\tau)^2$$

וע"י פתירת משוואה ריבועית קלה אנחנו מקבלים:

$$\tau = \frac{3-\sqrt{5}}{2}$$

וזה למעשה פרמטר החלוקה של האלגוריתם. שרטוט להמחשה:



אלגוריתם אינטרפולציה ריבועית, Quadratic Interpolation:

כאשר השתמשנו באלגוריתם Golden section לא ייחסנו חשיבות לסוג הפונקציה שאנחנו מנתחים, לא ייחסנו חשיבות לעד כמה טוב היא מקורבת לפונקציה האמיתית (במידה וזה אכן קירוב) וכן לא התייחסנו עד כמה היא "חלקה"/"גזירה". אם הפונקציה היא כן "חלקה" ו"טובה" אז אפשר להשתמש בידע שצברנו על נקודות מסויימות שכבר חישובנו עבור חישוב נקודות שעוד לא חישובנו, כלומר לבצע סוג של אינטרפולציה, כלומר אנחנו מקרבים את הפונקציה שיש ברשותנו לפונקציה אחרת אך שניתנת לביטוי אנליטי ואז אנחנו מחפשים את המינימום בצורה אנליטית – כלומר גוזרים ומשווים לאפס וכו'.

נניח ויש לנו פונקציה של משתנה יחיד  $f(x)$  והיא "חלקה" ו"טובה" ונניח שאנחנו יודעים שיש לה מינימום בקטע  $[a, c]$  ואנחנו מחפשים את נקודת המינימום (וכמובן גם את ערך הפונקציה בנקודה זו). נניח שאנחנו יודעים את ערכה של הפונקציה בשלוש נקודות,  $a < b < c$  ומתקיים:  $f(b) < f(a)$  וגם  $f(b) < f(c)$  - כלומר יש לנו שלושה ערכים של הפונקציה. נוכל לקרב את הפונקציה  $f(x)$  בקטע  $[a, c]$  ע"י פונקציה ריבועית כללית:  $q(x) = \alpha x^2 + \beta x + \gamma$  ע"י פתירת שלוש משוואות בשלושה נעלמים (הנעלמים הם מקדמי המשוואה הריבועית:  $\alpha, \beta, \gamma$ ):

$$\begin{cases} q(a) = f(a) = \alpha a^2 + \beta a + \gamma \\ q(b) = f(b) = \alpha b^2 + \beta b + \gamma \\ q(c) = f(c) = \alpha c^2 + \beta c + \gamma \end{cases}$$

כעת לאחר שיש לנו משוואה ריבועית,  $q(x) = \alpha x^2 + \beta x + \gamma$ , עם מקדמים ידועים, אפשר למצוא את המינימום של משוואה זו ע"י גזירה והשוואת הנגזרת לאפס וכו' וכך למצוא (בצורה אנליטית) את הנקודה שבה מתקבל המינימום של המשוואה הריבועית. האם זו הנקודה שבה מתקבל המינימום של משוואה  $f(x)$ ? לא בהכרח! (מדוע? כיוון שאף אחד לא אמר שמתקיים

$$f(x) \equiv q(x) \text{ לכל נקודה בקטע } [a, c], \text{ אלא אם הפונקציה } f(x) \text{ היא אכן ריבועית)}$$

נרצה להמשיך למצוא את הנקודה בה מתקבל המינימום של פונקציה  $f(x)$ . נבחין כי הנחנו כי מתקיים  $f(b) < f(c)$  וכן  $f(b) < f(a)$  ולתבנית זו אנחנו קוראים v-combination כי הצורה ששלושת הנקודות יוצרות על הגרף דומה לצורה "v":



נניח כי נקודת המינימום של המשוואה הריבועית היא  $b < m < c$  (כלומר מוכלת בקטע  $(b, c)$ ), אזי נוכל להמשיך בצורה איטרטיבית את האלגוריתם על הקטע  $[b, c]$  עם הנקודות  $b, m, c$  וזה עדיף כי אנחנו רוצים לשמר את הצורה v-combination כאשר הנקודה החדשה שמצאנו היא הנקודה הנמוכה ביותר ב-"v" כיוון ששימור זה מבטיח את יציבות והתקדמות האלגוריתם.



נבחין שבאלגוריתם זה, גודל הקטע אינו קטן באופן זהה מאיטרציה לאיטרציה, כלומר ייתכן ואנחנו נבצע הרבה איטרציות באלגוריתם אבל האלגוריתם לא יתקרב לפתרון ולכן ניתן באמצע הפעלת האלגוריתם (כאשר מגלים חוסר התקדמות), לעבור למשל לאלגוריתם אחר כדוגמת Golden section.

האלגוריתם Quadratic Interpolation:

- (1) מצא בקטע  $[a, c]$  נקודה שיוצרת v-combination, נסמנה להיות  $b$ .
- (2) מצא באופן אנליטי את המינימום של הפונקציה הריבועית שנוצרה ע"י v-combination, נסמנה להיות  $m$ .
- (3) בחר מתוך 4 הנקודות שנוצרו עד כה עוד v-combination וחזור לשלב (2). (אם ההתקדמות באלגוריתם לא מספקת ניתן לעבור לאלגוריתם Golden section).

מתי האלגוריתם Quadratic Interpolation יעיל?

אלגוריתם זה שמיש ויעיל במיוחד כאשר הפונקציה שעבורה אנחנו מחפשים את המינימום היא אכן בעלת קירוב ריבועי טוב בסביבת נקודת המינימום שאנחנו מחפשים.

קצב התכנסות אלגוריתם Quadratic Interpolation:

אם נסמן את ערך הפונקציה בנקודת מרכז ה-"v" באיטרציה ה- $k$  של האלגוריתם להיות  $f_k$  וכן נסמן את ערך הפונקציה בנקודת המינימום שאנחנו מחפשים להיות  $f^*$  אז נוכל לומר כי האלגוריתם ייתכנס אם יתקיים אי-השיוויון הבא:

$$f_k - f^* \leq c_k (f_{k-1} - f^*)$$

(כאשר ברור כי אם האלגוריתם מתכנס בהכרח מתקיים  $c_k < 1$ )

אם בשלב מסוים, הערך  $c_k = c_{k-1}$ , כלומר אנחנו מקבלים  $c_k$  קבוע (מאיטרציה לאיטרציה) אז אנחנו אומרים שקצב התכנסות האלגוריתם הוא לינארי. באלגוריתם Quadratic Interpolation מתקיים כי  $\lim_{k \rightarrow \infty} c_k = 0$  ולהתכנסות זו קוראים סופר-לינארית.

#### אלגוריתם Cubic Interpolation:

אם נניח כי עבור הפונקציה  $f(x)$ , אנחנו יודעים לחשב את הנגזרת של הפונקציה בכל נקודה שאנחנו נרצה (בין אם מישו ייתן לנו את ערך הנגזרת שאנחנו יכולים למדוד זו בצורה כלשהי) אז אפשר להשתמש באינטרפולציה חכמה אף יותר מהאינטרפולציה הריבועית. נניח ואנחנו יודעים כי  $f(x)$  מקבלת את המינימום שלה בקטע  $[a, c]$  אז במקום למצוא נקודה נוספת וליצור  $v$ -combination ולקרב אותה באמצעות משוואה ריבועית, אנחנו יכולים למצוא  $v$ -combination בצורה אחרת. אם למשל אנחנו יודעים כי  $f'(a) < 0$  וגם  $f'(b) > 0$  אז למעשה אנחנו מקבלים גם  $v$ -combination וכעת נוכל לקרב את הפונקציה  $f(x)$  באמצעות משוואה מסדר 3 עם מקדמים  $\alpha, \beta, \gamma, \delta$  שאנחנו מחפשים:

$$q(x) = \alpha x^3 + \beta x^2 + \gamma x + \delta$$

כיוון שיש לנו את 4 המשוואות הבאות (עם 4 נעלמים):

$$q(a) = f(a) = \alpha a^3 + \beta a^2 + \gamma a + \delta$$

$$q'(a) = f'(a) = 3\alpha a^2 + 2\beta a + \gamma$$

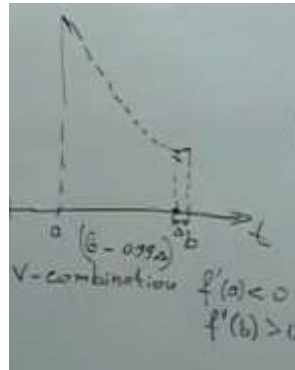
$$q(b) = f(b) = \alpha b^3 + \beta b^2 + \gamma b + \delta$$

$$q'(b) = f'(b) = 3\alpha b^2 + 2\beta b + \gamma$$

לאחר מציאת 4 המקדמים  $\alpha, \beta, \gamma, \delta$  אפשר למצוא בצורה אנליטית די פשוטה את הנקודה עבורה מתקבל המינימום של הפונקציה  $q(x)$  שננסה להיות  $b$  אז לחפש מבין 3 הנקודות  $a, b, c$  פעם נוספת אילו שתי נקודות תיתנה לנו  $v$ -combination ולהמשיך משם (זה אפשרי כי הנחנו כי גם עבור נקודה  $b$  אנחנו יכולים למצוא את ערך הנגזרת של פונקציה  $f(x)$ ).

#### קצב התכנסות אלגוריתם Cubic Interpolation:

אם ההתקדמות באלגוריתם לא מספיקה לנו אנחנו כמובן יכולים לבצע אלגוריתמים אחרים כמו Bisection. ייתכן גם מצב בו האלגוריתם, כאשר הוא מנסה להקטין את הקטע עליו הוא מחפש מינימום, מוצא נקודה חדשה  $b$  שמאד קרובה לקצה, אז ייתכן שהנקודה החדשה תהיה קרובה מידי לקצה וזה יכול לגרום לנו לעשות צעד הקטנה קטן, למשל:



נניח כי נקודת המינימום היא מאד מאד קרובה לקצה הימני של הקטע. אם הנקודה החדשה שמצאנו בצעד כלשהו של האלגוריתם, קרובה עד כדי  $\Delta$  לקצה הימני של הקטע אבל עדיין מצד שמאל לנקודת המינימום שלה אז השיפוע שלה הוא קטן מאפס, וזה טוב לנו, כי בצעד הבא אנחנו ניקח את הקטע  $[b - \Delta, b]$  שזה קטע קטן מאד והאלגוריתם מתקדם בצורה טובה. אך אם במקרה קרה שהנקודה החדשה נמצאת במצד ימין של נקודת המינימום, ולכן שיפוע הפונקציה בנקודה זו הוא חיובי, אז בצעד הבא של האלגוריתם, אשר בו אנחנו נדרשים למצוא *v-combination* אנחנו נחפש בקטע  $[a, b - \Delta]$  ואז למעשה אנחנו מתקדמים מאד לאט באלגוריתם, ולכן נגדיר כי אם הנקודה החדשה היא קרובה מידי לקצה, למשל קרובה יותר מערך  $\Delta_{\max}$  שנגדיר מראש, אז אנחנו ניקח נקודה  $b - \Delta_{\max}$  או  $a + \Delta_{\max}$  אם הנקודה החדשה קרובה לקצה השמאלי של הקטע) ובכך נוכל "לשמור" האלגוריתם בכך שהוא יבצע התקדמות טובה בכל צעד.



## הרצאה 08 – Multidimensional, Unconstrained Optimization Methods, אלגוריתמים לאופטימיזציה ללא אילוצים של פונקציות בעלות מספר משתנים

בעית אופטימיזציה כללית:

למצוא את הנקודה  $\bar{x} \in \mathbb{R}^n$  עבורה נקבל את הערך המינימלי של הפונקציה  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  שהיא פונקציה של מספר משתנים.

ישנן שיטות רבות לפתרון בעיות אופטימיזציה ושיטות אלו נהוג לחלק למחלקות.

שיטות מסוג Line Search:

שיטה זו היא שיטה איטרטיבית. באיטרציה ה- $k$  האלגוריתם נמצא בנקודה  $\bar{x}_k$  (והיא מועמדת להיות נקודת המינימום שאותה אנחנו מחפשים) והיא מחושבת ע"י הנוסחה הבאה:

$$\bar{x}_k = \bar{x}_{k-1} + \alpha_k \cdot \bar{d}_k, \quad \alpha_k \in \mathbb{R}, \bar{d}_k \in \mathbb{R}^n$$

כלומר הנקודה  $\bar{x}_k$  היא הנקודה מהאיטרציה הקודמת ( $k-1$ ) עם "הליכה" בכיוון  $\bar{d}_k$  ועם גודל צעד (step size) של  $\alpha_k$ , כלומר האלגוריתם מחפש את נקודת המינימום באמצעות הליכה על קו כלשהו (ייתכן שכיוון הקו משתנה), של ערכים בפונקציה ובכדי שזה יהיה אלגוריתם טוב נראה שלאורך קו זה ערכי הפונקציה בנקודות לאורך הקו ילכו וירדו (כי אנחנו מחפשים את מינימום הפונקציה). בהמשך נראה כי נוח מאד למשל להגדיר את כיוון ההליכה ככיוון המנוגד לכיוון הגרדיאנט, שכן הגרדיאנט מצביע בכיוון העליה החדה ביותר של ערכי הפונקציה.

בכדי שכיוון ה"הליכה" יהיה בכיוון שבו ערכי הפונקציה יורדים, בהכרח חייב להתקיים התנאי על הנגזרת הראשונה של הפונקציה, כלומר הנגזרת הראשונה חייבת להיות שלילית, ובפונקציה של מספר משתנים אנחנו אומרים כי זו הנגזרת הכיוונית, כלומר נדרוש שהנגזרת הכיוונית בכיוון ההליכה תהיה תמיד שלילית:

$$f'_{\bar{d}_k}(\bar{x}_k) = (\nabla f(\bar{x}_k))^T \cdot \bar{d}_k < 0$$

בחירת גודל הצעד  $\alpha_k$ :

בכל איטרציה אנחנו נדרשים לספק לאלגוריתם את גודל הצעד שאנחנו מעוניינים לבצע. לבחירת גודל הצעד יש כמה אפשרויות:

(1) שיטה שנקראת Exact Line Search:

בשיטה זו אנחנו בוחרים את  $\alpha_k$  להיות הערך שמביא את הביטוי  $f(\bar{x}_k + \alpha_k \cdot \bar{d}_k) = f(\bar{x}_{k+1})$  למינימום, כלומר:

$$\alpha_k = \arg \min_{\alpha} f(\bar{x}_k + \alpha \bar{d}_k)$$

למעשה זו בעית אופטימיזציה במשתנה יחיד בפני עצמה וראינו דרכים לפתור בעיה כזו – למשל ע"י אלגוריתם Bisection, או אינטרפולציה וכו'. נשים לב כי פתרון בעית האופטימיזציה הזו היא רק חלק מבעית האופטימיזציה הכללית שאנחנו מנסים לפתור ולכן רצוי לא לבזבז יותר מידי זמן בחיפוש  $\alpha_k$  אופטימלי – צריך למצוא איזון בין חיפוש  $\alpha_k$  בכל איטרציה לבין הזמן שלוקח לנו למצוא פתרון אופטימלי לבעיה הכללית שלנו.

(2) שיטת Inexact Line Search:

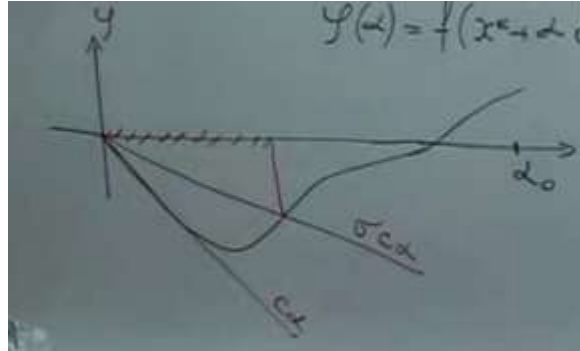
דוגמה לשיטה זו היא עבודה לפי חוק Armijo. בשיטה זו ננסה למצוא את  $\alpha_k$  באופן הבא:

$$\alpha_k = \arg \min_{\alpha} \varphi(\alpha) = \arg \min_{\alpha} \left\{ f(\bar{x}_k + \alpha \bar{d}_k) - f(\bar{x}_k) \right\}$$

חוק Armijo מציג להסתכל על הבעיה באופן הבא: נבחין כי מתקיים לפי מה שהגדרנו זה עתה  $\varphi(0) = 0$  וכן מתקיים (לפי הדרישה על הנגזרת הכיוונית):

$$\varphi'(\alpha) \Big|_{\alpha=0} = (\nabla f(\bar{x}_k))^T \bar{d}_k < 0$$

נסמן את השיפוע של  $\varphi(\alpha)$  בנקודה  $\alpha = 0$  להיות  $c < 0$ . כמו כן, החוק אומר שצריך להגדיר קו לינארי נוסף בעל שיפוע  $c\sigma < 0$  ואז נוכל לשרטט את הפונקציה  $\varphi(\alpha)$  באופן הבא (זו הקו העקום, זו פונקציה עם משתנה ממימד אחד):



(אנחנו מנסים למצוא את הנקודה  $\alpha$  עבורה נקבל את הערך המינימלי של הפונקציה  $\varphi(\alpha)$ )

כעת החוק אומר שאם נתחיל לפתור את בעיית האופטימיזציה  $\alpha_k = \arg \min_{\alpha} \varphi(\alpha)$  החל מ- $\alpha_0$  כלשהו, אז ברגע שנגיע לנקודה  $\alpha_k$  שבו אנחנו מקבלים כי:

$$\varphi(\alpha_k) < \sigma c \alpha_k$$

במצב זה אנחנו יכולים לעצור ולהחזיר את הערך  $\alpha_k$  שמצאנו. ההתקדמות בציר  $\alpha$  נעשית ע"י הכפלה בפרמטר  $0 < \beta < 1$  (למשל  $\beta = 0.2$ ) ובכך אנחנו מבטיחים להתקדם לכיוון ראשית הצירים. בשיטה זו יש כמה שיקולים שצריך לקחת בחשבון, אם נקבל גודל צעד מאד קטן ( $\alpha_k \approx 0$ ) אז אנחנו נקבל את ההתקדמות בכיוון הנכון ביותר (מבחינת הנגזרת הכיוונית) אך ההתקדמות הכוללת באלגוריתם תהיה מאד קטנה, מצד שני, אם נבחר צעדים גדולים, יש סיכוי גבוה יותר שאנחנו הולכים בכיוון שהוא לא מדויק מספיק על מנת להגיע לפתרון האופטימלי בצורה המהירה ביותר ולכן הפרמטר  $\sigma$  מאפשר לנו סוג של פשרה. באופן כללי נהוג לבחור  $0 < \sigma < 1$ . לעיתים נבחר אפילו  $\sigma = 10^{-4}$ .

(3) שיטת Constant step size:

בשיטה זו אנחנו קובעים מראש כי גודל הצעד הוא גודל (קטן יחסית) קבוע  $\alpha = \text{Const}$ . הסיכון בשיטה זו הוא שאם בוחרים בטעות קבוע גדול מידי, ייתכן והאלגוריתם יתבדר (!) וכלל לא יתכנס לפתרון האופטימלי שאנחנו מחפשים. לעומת זאת, אם נבחר גודל צעד קבוע קטן מידי, האלגוריתם יתכנס אך באופן איטי למדי.

(4) שיטת Diminishing step size:

בשיטה זו אנחנו מקטינים את גודל הצעד מאיטרציה לאיטרציה, אך צריך לדאוג כי השאיפה של גודל הצעד לאפס לא מתרחשת מהר מידי כי אז האלגוריתם יתכנס באופן איטי מאד לפתרון. נהוג לבחור קצב "הקטנה" באופן הבא:

$$\alpha_k \rightarrow 0, \quad \lim_{l \rightarrow \infty} \sum_{i=1}^l \alpha_k \rightarrow \infty$$

למשל ניתן לבחור  $\alpha_k = \frac{1}{k}$  שזה הטור ההרמוני שידוע מחדו"א 1 כי הוא לא מתכנס.

דוגמה: (לשימוש באלגוריתם Line Search לפי Steepest Descent, כלומר הולכים בכיוון הירידה החדה ביותר, או בשם אחד זה נקרא Gradient Descent)

נשתמש באלגוריתם:

$$\bar{x}_k = \bar{x}_{k-1} + \alpha_k \cdot \bar{d}_k, \quad \alpha_k \in \mathbb{R}, \bar{d}_k \in \mathbb{R}^n$$

כאשר נגדיר כי כיוון הצעד, הוא הכיוון המנוגד לכיוון הגרדיאנט (כי הגרדיאנט מצביע לכיוון העליה החדה ביותר בעוד אנחנו מחפשים את הירידה החדה ביותר), כלומר:

$$\bar{d}_k = -\nabla f(\bar{x}_k)$$

ולכן הצעד שלנו הוא:

$$\bar{x}_k = \bar{x}_{k-1} - \alpha_k \cdot \nabla f(\bar{x}_k)$$

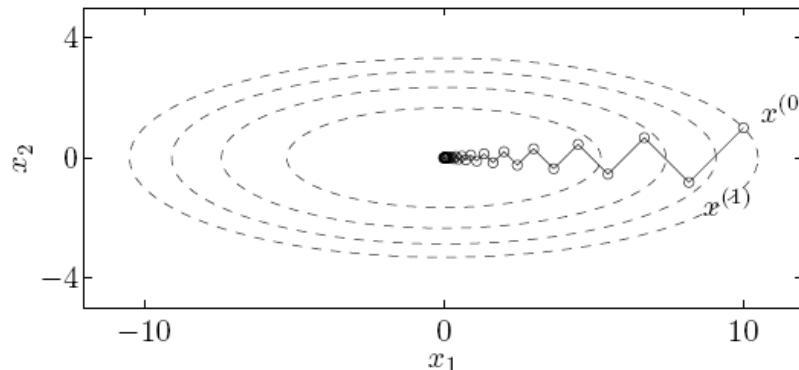
נניח כי  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$  ואנחנו נשרטט את קווי הגובה של הפונקציה במערכת הצירים  $(x_1, x_2)$ .

בשביל לחשב את גודל הצעד בכל איטרציה נשתמש בשיטת Exact Line Search ולכן אם בכל שלב נתקדם על הקו שמנוגד לכיוון הגרדיאנט, אנחנו נבחר את גודל הצעד שייתן לנו את הנקודה לאורך קו (הגרדיאנט) שהפונקציה מקבלת בו את ערכה המינימלי. באיטרציה הבאה, אנחנו שוב נבדוק את כיוון הגרדיאנט בנקודה שבה אנחנו עומדים ונבצע את אותו התהליך של מציאת גודל הצעד, שביצענו, פעם נוספת.

אם למשל קווי הגובה הן אליפסות שמאונכות לאחד הצירים, אנחנו יכולים אולי להניח כי אם הפונקציה היא מהצורה  $f(\bar{x}) = \bar{x}^T A \bar{x}$  אז יש למטריצה  $A$  ערכים עצמיים, אחד גדול ואחד קטן, כי אם למשל המטריצה היא אלכסונית (ואז ערכי האלכסון אלו הם הערכים העצמיים) אז נקבל:

$$A = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} \Rightarrow f(\bar{x}) = \bar{x}^T A \bar{x} = x_1 \lambda_1 + x_2 \lambda_2$$

ולכן אם  $\lambda_1 \ll \lambda_2$  אז אנחנו רואים כי עבור שינויים קטנים ב-  $x_2$  נקבל שינויים גדולים בערכי הפונקציה ועבור שינויים קטנים ב-  $x_1$  נקבל שינויים קטנים. ואם המטריצה  $A$  היא לא אלכסונית אז כיוון האליפסות יהיה במאונך לאחד מהווקטורים העצמיים ומקביל לשני (כי ווקטורים עצמיים הם אורתוגונליים).



(האליפסות מאונכות לציר האנכי ומקבילות לציר האופקי)

אם נמשיך באלגוריתם שלנו, אנחנו נקבל דרך שהיא סוג של "זיג-זג" וגודל הצעדים ילכו ויקטנו מאיטרציה לאיטרציה.  
קצב התכנסות האלגוריתם:

לא נוכיח בדוגמה זו את קצב ההתכנסות, אך קצב ההתכנסות של האלגוריתם הוא קצב התכנסות לינארי של הירידה החדה ביותר.

טענה (שלא נוכיח): אם נניח כי ההסיאן של הפונקציה  $f(\bar{x})$  הוא  $H(\bar{x})$  והערכים העצמיים שלו הם  $\lambda_{\min}, \lambda_{\max}$  (זו מטריצה מסדר 2 ולכן יש רק שני ערכים עצמיים) אז מתקיים:

$$\left| f(\bar{x}_{k+1}) - f(\bar{x}^*) \right| \leq c \left( f(\bar{x}_k) - f(\bar{x}^*) \right)$$

כאשר  $f(\bar{x}^*)$  זה הערך המינימלי של הפונקציה  $f(\bar{x})$  וכן  $c \in \mathbb{R}$  הוא קבוע.

לקשר זה בין שתי איטרציות באלגוריתם קוראים התכנסות לינארית, והקבוע  $c$  נקרא קצב ההתכנסות. במקרה זה מתקיים:

$$c = 1 - \frac{\lambda_{\min}}{\lambda_{\max}}$$

במקרה בו מתקיים  $\lambda_{\min} = \lambda_{\max}$  מתקיים  $c = 0$  וזה אומר:

$$f(\bar{x}_{k+1}) - f(\bar{x}^*) \leq c \left( f(\bar{x}_k) - f(\bar{x}^*) \right) = 0$$

וזה אומר כי האלגוריתם יתכנס כבר לאחר צעד אחד. נוכל לומר כי אם  $\lambda_{\min} = \lambda_{\max}$  אז קווי הגובה של הפונקציה  $f(\bar{x})$  הם עיגולים מושלמים, ולכן כבר בצעד הראשון הכיוון המנוגד לגרדיאנט הוא בדיוק לכיוון המינימום שאותו אנחנו מחפשים ולכן ע"י שיטת Exact Line אנחנו כבר אחרי צעד אחד נמצא את המינימום הנדרש.

אם במקרה אחר, נקבל למשל:  $\frac{\lambda_{\min}}{\lambda_{\max}} = 10^{-3}$  אז נקבל  $c = 0.999$  וזה אומר כי האלגוריתם יתכנס בצורה לינארית אך בצורה

איטית מאד.

הגדרה: Condition Number

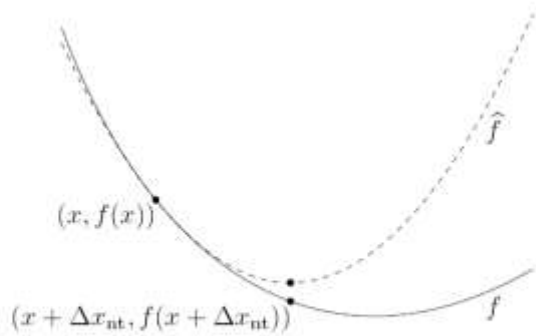
נהוג לסמן  $\frac{\lambda_{\max}}{\lambda_{\min}} = \theta$  וזה נקרא condition number (נבחין כי זה היחס ההפוך בין הערכים העצמיים). אנחנו נוהגים לומר כי

אם  $\theta$  הוא גדול אז מטריצת ההסיאן היא ill defined, וקצב ההתכנסות של האלגוריתם הוא איטי ולכן במקרים כאלה נרצה ללמוד שיטות נוספות לאופטימיזציה.

שיטת ניוטון: (Newton method):

זוהי שיטה הרבה יותר "חזקה" משיטת Line search.

נניח תחילה כי נתונה פונקציה  $f(x): \mathbb{R} \rightarrow \mathbb{R}$  וערכה המינימלי מתקבל בנקודה  $x^*$  ואנחנו נמצאים באיטרציה  $k$  של האלגוריתם בנקודה  $x_k$ . אם אנחנו נניח למשל כי הפונקציה היא דומה לפונקציה ריבועית אז נרצה לקרב אותה באמצעות פונקציה ריבועית (פולינום). לאחר מכן, נמצא את המינימום של הפולינום באופן אנליטי ונפעיל את האלגוריתם על הנקודה שמצאנו שהיא המינימום של הפולינום. למעשה נקבל את השרטוט הבא:



**Figure 9.16** The function  $f$  (shown solid) and its second-order approximation  $\hat{f}$  at  $x$  (dashed). The Newton step  $\Delta x_{nt}$  is what must be added to  $x$  to give the minimizer of  $\hat{f}$ .

השיטה: (שיטת ניוטון)

עבור פונקציה של מספר משתנים  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  ונסמן:

$$\bar{g}_k \equiv \bar{g}(\bar{x}_k) \equiv \nabla f(\bar{x}_k)$$

$$H_k \equiv H(\bar{x}_k) \equiv \nabla^2 f(\bar{x}_k)$$

וכן פיתח טיילור מסדר שני של הפונקציה  $f(\bar{x})$  סביב הנקודה  $\bar{x}_k$  יסומן להיות:

$$f(\bar{x}_k + \bar{d}_k) = f(\bar{x}_k) + \bar{g}_k^T \bar{d}_k + \frac{1}{2} \bar{d}_k^T H(\bar{x}_0) \bar{d}_k + \dots \Rightarrow q(\bar{d}_k) \equiv f(\bar{x}_k) + \bar{g}_k^T \bar{d}_k + \frac{1}{2} \bar{d}_k^T H(\bar{x}_0) \bar{d}_k$$

שיטת ניוטון ממזערת את הקירוב הריבועי  $q(\bar{d}_k)$  ובכדי למצוא מינימום של פונקציה ריבועית עם מספר משתנים נדרוש שהגרדיאנט (גוזרים לפי  $\bar{d}_k$ ) יתאפס:

$$\text{Grad}(q(\bar{d}_k)) = \nabla q(\bar{d}_k) = \bar{g}_k + H_k \bar{d}_k = 0$$

משוואה זו נקראת משוואת ניוטון. כמו כן, למדנו כי מתקיים  $H_k \bar{d}_k = -\bar{g}_k$ , ולכן על מנת לחשב את  $\bar{d}_k$  שהוא כיוון

ההתקדמות אפשר למעשה לחשב את הביטוי  $\bar{d}_k = -\frac{\bar{g}_k}{H_k}$  (ע"י מטלב למשל) ואז לבחור את גודל הצעד ע"י אלגוריתם

כדוגמת Exact Line search.

נבחין שבשביל למצוא את הנקודה בה הגרדיאנט  $\nabla q(\bar{d}_k)$  מתאפס, אנחנו למעשה פותרים בעיה אופטימיזציה במספר משתנים ולכן ניתן לפתור בעיה זו בשיטת Line search. כיוון הצעד הוא:

$$\bar{d}_k = -\frac{\bar{g}_k}{H_k} = -H_k^{-1} \cdot \bar{g}_k$$

כאשר  $\bar{d}_k$  נקרא הכיוון של ניוטון (Newton's direction) ואז באיטרציה הבא נקבל:

$$\bar{x}_{k+1} = \bar{x}_k + \alpha_k \bar{d}_k$$

ואת גודל הצעד ניתן לחשב למשל לפי Exact line search או בדרך אחרת.

את שיטת ניוטון ניתן לבצע לא רק באמצעות Line search אלא גם בשיטות האחרות שלא נדון בהן בקורס זה אך אחת מהן היא למשל Trust Region (ניתן למצוא בספר הלימוד).

קצב התכנסות האלגוריתם של שיטת ניוטון – התכנסות ריבועית אסימפטוטית:

ברור כי במידה והפונקציה  $f(\bar{x})$  היא אכן ריבועית, אז תספיק לנו איטרציה אחת בלבד של האלגוריתם על מנת להגיע לפתרון האופטימלי כיוון שאם היא ריבועית אז הקירוב שלנו (שהוא פולינום ריבועי) הוא בדיוק הפונקציה שאנחנו מחפשים עבורה את המינימום.

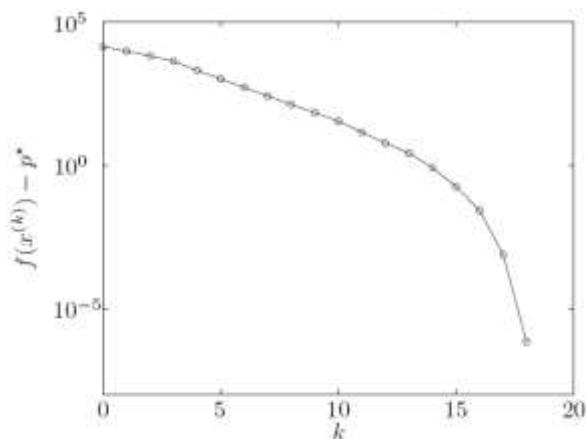
במקרה בו הפונקציה  $f(\bar{x})$  היא לא ריבועית ובנקודת האופטימום מתקיים כי  $H(\bar{x}) \succ 0$  וגם  $H(\bar{x})$  רציפה, אז ההתכנסות היא התכנסות ריבועית אסימפטוטית, כלומר אם באיטרציה מסוימת נגדיר:

$$\delta_k \triangleq \|\bar{x}_k - \bar{x}^*\| \quad or \quad \delta_k \triangleq \|f(\bar{x}_k) - f(\bar{x}^*)\|$$

אז מתקיים:

$$\delta_{k+1} = c\delta_k^2$$

ולכן אם באיטרציה מסוימת מתקיים  $\delta_k = 10^{-2}$  אז עבור  $c = 1$  נקבל כבר באיטרציה הבאה  $\delta_{k+1} = 10^{-4}$  ולאחר מכן  $\delta_{k+2} = 10^{-8}$  וזו התכנסות מאד מהירה (אקספוננציאלית). התכנסות אסימפטוטית מתרחשת רק כאשר האלגוריתם כבר קרוב לפתרון והבעיה היא שאנחנו לא יודעים מתי האלגוריתם מתקרב לפתרון ולכן רק לאחר שאנחנו רואים כי האלגוריתם קרוב לפתרון אנחנו נוכל לדעת כי ההתכנסות החל משלב זה היא מאד מהירה. אפשר לתאר את הגרף של השגיאה בצורה הבאה:



## **הרצאה 09 – Another view of Newton's Meth. Via solution of system of nonlinear equations, מבט נוסף על שיטת ניוטון באמצעות פתירת מערכת משוואות לא לינאריות**

כאשר נתונה לנו בעיית אופטימיזציה, למשל למצוא את  $\bar{x} \in \mathbb{R}^n$  אשר יביא את ערך הפונקציה  $f(\bar{x}): \mathbb{R}^n \rightarrow \mathbb{R}$  למינימום ניתן במקום זאת למצוא את  $\bar{x} \in \mathbb{R}^n$  אשר מאפס את הגרדיאנט של הפונקציה, כלומר מקיים  $\nabla f(\bar{x}) \equiv \bar{g}(\bar{x}) = 0$  כי למדנו שאלו הם תנאי מספיקים עבור פונקציה שהיא קמורה. כיוון שהגרדיאנט הוא וקטור, ואנחנו רוצים שהוא יהיה שווה לווקטור האפס, זה אומר שאנחנו למעשה צריכים לדרוש שכל הרכיבים שלו יתאפסו, ולכן למעשה יש לפנינו מערכת משוואות לא לינאריות (כי באופן כללי הגרדיאנט לא חייב להיות לינארי), כלומר אנחנו צריכים לפתור מערכת של  $n$  משוואות לא לינאריות:

$$\bar{g}(\bar{x}) = \bar{0} \Leftrightarrow g_i(\bar{x}) = 0, \forall i=1, \dots, n$$

נניח ואנחנו עומדים בנקודה מסוימת  $\bar{x}_k$ , אז אמרנו כי הרכיבים של  $\bar{g}(\bar{x}_k)$  הם לא בהכרח לינאריים, אך אם נכתוב את פיתוח טיילור עד סדר ראשון של כל רכיב נקבל  $n$  משוואות לינאריות:

$$\bar{g}(\bar{x}_k + \bar{d}_k) \cong g(\bar{x}_k) + H(\bar{x}_k)\bar{d}_k = 0$$

ולכן הפתרון למערכת המשוואות נתון ע"י:

$$\bar{d}_k = -H^{-1}(\bar{x}_k)g(\bar{x}_k)$$

וזה בדיוק צעד ניוטון שלמדנו בהקשר של שיטת ניוטון. נבחין כי אנחנו רוצים לדרוש כי כיוון הווקטור  $\bar{d}_k$  יהיה בכיוון בו ערכי הפונקציה  $f(\bar{x})$  יורדים, כלומר נדרוש כי הנגזרת הכיוונית  $f'_{\bar{d}_k}(\bar{x}) < 0$ . כמו כן:

$$f'_{\bar{d}_k}(\bar{x}) = \bar{g}_k^T \bar{d}_k \underset{\bar{d}_k = -H^{-1}(\bar{x}_k)g(\bar{x}_k)}{=} -\bar{g}_k^T H_k^{-1} \bar{g}_k < 0 \Leftrightarrow \bar{g}_k^T H_k^{-1} \bar{g}_k > 0$$

וזה דרישה שקולה לדרישה:  $H(\bar{x}_k) \succ 0$  (זה דרישה שקולה כי אם  $H^{-1}(\bar{x}_k) \succ 0$  אז זה שקול לזה שהערכים של המטריצה  $H^{-1}(\bar{x}_k)$  הם חיוביים, והערכים העצמיים של המטריצה  $H(\bar{x}_k)$  הם "1 חלקי" הערכים העצמיים של  $H^{-1}(\bar{x}_k)$ ) ולכן זה אכן דרישה שקולה). נבחין כי במידה והפונקציה  $f(\bar{x})$  לא קמורה ממש (וייתכן אפילו לא קמורה) אז אין הכרח שהתנאי יתקיים  $H(\bar{x}_k) \succ 0$  (לפי משפט שראינו: אם  $f(\bar{x})$  קמורה (ממש) אז מתקיים  $H(\bar{x}_k) \succeq 0$  ( $H(\bar{x}_k) \succ 0$ )) אבל אז זה אומר כי יש ערכים עצמיים שהם שליליים ולכן נוכל להשתמש ב"טריק" ולמעשה להוסיף להיסאן מטריצה אלכסונית  $\Delta_k$  כך שיתקיים:

$$H(\bar{x}_k) + \Delta_k \succ \varepsilon I$$

ומשמעות הביטוי האחרון זה אומר כי כל הע"ע של המטריצה באגף השמאלי, גדולים מהע"ע של המטריצה באגף הימני. אפשר באותה מידה לכתוב זאת כך (נשאיר את ההוכחה לשיעורי הבית):

$$H(\bar{x}_k) + \Delta_k - \varepsilon I \succ 0$$

וכעת נמצא את כיוון הצעד של שיטת ניוטון באמצעות:

$$\bar{d}_k = -\left(H(\bar{x}_k) + \Delta_k\right)^{-1} g(\bar{x}_k)$$

פתרון של מערכת משוואות  $A\bar{x} = \bar{b}$  כאשר המטריצה  $A$  סימטרית חיובית:

ישנן שיטות רבות לפתירת מערכת משוואות לינאריות ואנחנו צריכים לבחור את השיטה המתאימה ביותר.

שיטת Cholesky factorization (או Cholesky decomposition):

בשיטה זו, אנחנו רוצים לייצג את המטריצה הסימטרית חיובית  $A$  בתור מכפלה של מטריצה  $L$  במטריצה  $L^T$ , כאשר  $L$  היא מטריצה משולשית תחתונה, כלומר:

$$A = L \cdot L^T, \quad L = \begin{pmatrix} l_{11} & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ l_{n1} & \dots & \dots & l_{nn} \end{pmatrix}$$

לא נסביר את איך השיטה בדיוק עובדת (ניתן למצוא את ההסבר בספר הלימוד) אך נזכיר כי מספר הפעולות של שיטה זו עד שמוצאים מטריצה  $L$  בגודל  $n \times n$  הוא במידה הבעיה:

$$O\left(\frac{n^3}{6}\right)$$

נבחין כי מספר פעולות זה הוא מאד נמוך, שכן רק להכפיל שתי מטריצות ריבועיות מסדר  $n \times n$  לוקח  $n^3$  פעולות.

דוגמה:

נניח כי לפנינו מערכת המשוואות הבאה:  $L\bar{y} = \bar{b}$  (כאשר  $L$  היא מטריצה משולשית תחתונה מסדר  $3 \times 3$  ווקטור הנעלמים הוא  $\bar{y}$ ) אז באופן מפורש, מערכת המשוואות היא:

$$\begin{aligned} l_{11}y_1 &= b_1 \\ l_{21}y_1 + l_{22}y_2 &= b_2 \\ l_{31}y_1 + l_{32}y_2 + l_{33}y_3 &= b_3 \end{aligned}$$

קל לראות שנוכל לחלץ בקלות את הנעלם הראשון  $y_1$  מהמשוואה הראשונה ונציב זאת במשוואה השנייה. לאחר כן, ניתן לראות כי לאחר ההצבה, ניתן בקלות לחלץ את הנעלם השני  $y_2$  ולהמשיך באופן דומה עד אשר חילצנו את כל הנעלמים. מספר הפעולות הדרושות לביצוע פתרון זה הוא  $O(n^2)$  ודרך פתרון זו נקראת forward substitution.

דוגמה:

נניח ואנחנו רוצים לפתור את מערכת המשוואות  $A\bar{x} = \bar{b}$  (זהו ווקטור הנעלמים וכן  $A$  היא מטריצה חיובית) אז נניח כי המרנו את המערכת לפי שיטת Cholesky factorization למערכת הבאה:

$$L \cdot L^T \cdot \bar{x} = \bar{b}$$



נסמן  $\bar{y} \triangleq L^T \cdot \bar{x}$  ואז נפתור באמצעות forward substitution את המערכת  $L \cdot \bar{y} = \bar{b}$ . כעת נשאר לנו לפתור את המשוואה  $L^T \cdot \bar{x} = \bar{y}$  בכדי למצוא את ווקטור הנעלמים  $\bar{x}$  אבל נבחין כי  $L^T$  היא מטריצה משולשית עליונה ולכן נפתור את המערכת באמצעות backward substitution (באופן מאד דומה ל-forward substitution, רק שנתחיל מהמשוואה האחרונה).  
בשיטה זו, פתרנו את מערכת המשוואות  $A\bar{x} = \bar{b}$  עם סדר גודל של  $O(n^2)$  פעולות.

#### שיטת Modified Cholesky factorization:

גם אם המטריצה  $A$  לא חיובית (שזה תנאי הכרחי לביצוע השיטה) ניתן להשתמש בשיטת Cholesky רק עם שינוי, יש תחילה לחשב את המטריצה האלכסונית  $\Delta$  אשר ניתן להוסיף למטריצה  $A$  בכדי שהמטריצה  $A + \Delta$  תהיה חיובית. לא נתאר כאן את אופן החישוב אבל יש פונקציה במטלב שיודעת למצוא את המטריצה  $\Delta$  המתאימה.

#### מקרה פרטי של שיטת ניוטון: Least Squares Problem

נתאר מציאת פתרון אופטימלי למערכת משוואות עם יותר משוואות מנעלמים לפי סכום ריבועים מינימלי (Least Squares Problem) לפי השיטה של ניוטון וגאוס (שיטת ניוטון-גאוס). תחילה נתאר את בעיית ה-Least Squares Problem ולאחר מכן נציג את הפתרון לפי שיטת ניוטון-גאוס.

נניח כי מערכת המשוואות, מתוארת ע"י פונקציה ווקטורית (ניתן לחשוב על מטריצה כעל פונקציה ווקטורית, כלומר העתקה מרחב ווקטורי אחד למרחב ווקטורי אחר)  $\bar{g}(\bar{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$  ונניח כי אנחנו יודעים כי הפונקציה  $\bar{g}(\bar{x})$  כמעט מתאפסת בנקודה מסויימת ולכן המערכת שאנחנו רוצים לפתור היא:  $\bar{g}(\bar{x}) \approx \bar{0}$  (כלומר אנחנו מניחים כי לא ניתן למצוא נקודה (ווקטור)  $\bar{x}$  עברה נקבל ממש אפס אך ניתן להתקרב לזה). באופן פורמלי אנחנו מחפשים את הפתרון  $\bar{x}$  אשר יביא למינימום את סכום הריבועים של הרכיבים של  $\bar{g}(\bar{x})$  כלומר:

$$\min_{\bar{x} \in \mathbb{R}^n} f(\bar{x}) = \frac{1}{2} \|\bar{g}(\bar{x})\|_2^2 = \frac{1}{2} \sum_{i=1}^m g_i^2(\bar{x})$$

באופן כללי, אין צורך לבחור דווקא את סכום הריבועים ואפשר גם לפתור את הבעיה:

$$\min_{\bar{x} \in \mathbb{R}^n} f(\bar{x}) = \sum_{i=1}^m \varphi(g_i(\bar{x}))$$

כאשר  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$  היא פונקציה כלשהי לבחירתנו. אנחנו ננתח את המקרה הכללי הנ"ל ונרצה למשל למצוא את הגרדיאנט של הפונקציה  $\varphi$  בנקודה מסויימת וגם את ההסיאן. לפונקציות הללו, אנחנו קוראים Penalty Functions שכן הן קובעות מה השגיאה שנוצרת לנו בעקבות פתרון שמצאנו.

אם נסמן:

$$f(\bar{x}) = \sum_{i=1}^m \varphi(g_i(\bar{x}))$$

אז מתקיים (מומלץ להוכיח כשיעורי בית (הוכחנו בכיתה בהרצאה הראשונה את הבעיה הזו בדיוק עם שמות שונים לפונקציות)):

$$\begin{aligned}\nabla f(\bar{x}) &= \sum_{i=1}^m (\varphi'(g_i(\bar{x})) \nabla g_i(\bar{x})) \\ H(\bar{x}) &= \nabla^2 f(\bar{x}) = \sum_{i=1}^m (\varphi''(g_i(\bar{x})) \nabla g_i(\bar{x}) \nabla g_i^T(\bar{x}) + \varphi'(g_i(\bar{x})) \nabla^2 g_i(\bar{x}))\end{aligned}$$

נבחין כי הביטוי  $\nabla g_i(\bar{x}) \nabla g_i^T(\bar{x})$  הוא מטריצה שכן  $\nabla g_i(\bar{x})$  הוא ווקטור עמודה. מטריצה כזו נקראת Rank One Matrix וזה משום שיש לה רק ע"ע אחד שונה מאפס. נרצה לרשום את שתי הנוסחאות לעיל בכתובי מטריצי וקומפקטי.

נסמן תחילה את  $\nabla \bar{g}(\bar{x})$  להיות מטריצה שכל אחת מהעמודות שלה היא  $\nabla g_i(\bar{x})$ . כלומר:

$$\nabla \bar{g}(\bar{x}) \triangleq (\nabla g_1(\bar{x}) \quad \cdots \quad \nabla g_m(\bar{x}))$$

בנוסף, נסמן את הווקטור:

$$\varphi'(\bar{g}(\bar{x})) = \begin{pmatrix} \varphi'(g_1(\bar{x})) \\ \vdots \\ \varphi'(g_m(\bar{x})) \end{pmatrix}$$

נסמן בנוסף את המטריצה האלקסונית:

$$\mathcal{G}''(\bar{g}(\bar{x})) = \begin{pmatrix} \varphi''(g_1(\bar{x})) & & 0 \\ & \ddots & \\ 0 & & \varphi''(g_m(\bar{x})) \end{pmatrix}$$

כעת נוכל לרשום באופן מטריצי וקומפקטי את שתי הנוסחאות שרשמנו קודם באופן הבא:

$$\begin{aligned}\nabla f(\bar{x}) &= \nabla \bar{g}(\bar{x}) \varphi'(\bar{g}(\bar{x})) \\ H(\bar{x}) &= \nabla^2 f(\bar{x}) = \nabla \bar{g}(\bar{x}) \mathcal{G}''(\bar{g}(\bar{x})) (\nabla \bar{g}(\bar{x}))^T + \sum_{i=1}^m \varphi'(g_i(\bar{x})) \nabla^2 g_i(\bar{x})\end{aligned}$$

נחזור למקרה הפרטי של בעיית Least Squares ולכן הפונקציה Penalty Function היא:  $\varphi(t) = \frac{1}{2} t^2$  ולכן הסימונים שלנו מקבלים את הצורה הבאה:

$$\varphi(\bar{t}) = \frac{1}{2} \|\bar{t}\|_2^2, \quad \varphi'(\bar{t}) = \bar{t}, \quad \varphi''(\bar{t}) = 1 \quad \Rightarrow \quad \mathcal{G}''(\bar{t}) \equiv I$$

ולכן נקבל בצורה קומפקטית את הבעיה Least Squares:

$$\begin{aligned} f(\bar{x}) &= \sum_{i=1}^m \varphi(g_i(\bar{x})) = \frac{1}{2} \|\bar{g}(\bar{x})\|_2^2 = \frac{1}{2} \sum_{i=1}^m g_i^2(\bar{x}) \\ \nabla f(\bar{x}) &= \overbrace{\nabla \bar{g}(\bar{x})}^{\text{Matrix}} \cdot \overbrace{\bar{g}(\bar{x})}^{\text{Column Vector}} \\ H(\bar{x}) &= \nabla^2 f(\bar{x}) = \nabla \bar{g}(\bar{x}) (\nabla \bar{g}(\bar{x}))^T + \sum_{i=1}^m g_i(\bar{x}) \nabla^2 g_i(\bar{x}) \end{aligned}$$

שיטת ניוטון-גאוס (Newton-Gauss Method):

נניח כי אנחנו מחפשים את הנקודה  $\bar{x}$  אשר מביאה למינימום את  $f(\bar{x}) = \frac{1}{2} \|\bar{g}(\bar{x})\|_2^2 = \frac{1}{2} \sum_{i=1}^m g_i^2(\bar{x})$  כלומר:

$$\bar{x} = \arg \min f(\bar{x}) = \frac{1}{2} \|\bar{g}(\bar{x})\|_2^2 = \frac{1}{2} \sum_{i=1}^m g_i^2(\bar{x})$$

נניח כי האלגוריתם לפתרון שלנו נמצא כבר מאד קרוב לפתרון האופטימלי. נתבונן בביטוי שמצאנו בכתוב מטרצי עבור ההסיון.

נבחין כי אם אנחנו קרובים מאד לפתרון אז כל אחד מהרכיבים בסכום  $\frac{1}{2} \sum_{i=1}^m g_i^2(\bar{x})$  הוא מאד קטן, כלומר הגרדיאנט  $\bar{g}(\bar{x})$

מאד קטן (קרוב לווקטור האפס) ולכן הביטוי השני בהסיון  $\sum_{i=1}^m g_i(\bar{x}) \nabla^2 g_i(\bar{x})$  הוא מאד קטן ולכן הביטוי הדומיננטי בהסיון

הוא הביטוי:  $\nabla \bar{g}(\bar{x}) (\nabla \bar{g}(\bar{x}))^T$ . ניתן לרשום:

$$\begin{aligned} H(\bar{x}) &= \nabla^2 f(\bar{x}) = \nabla \bar{g}(\bar{x}) (\nabla \bar{g}(\bar{x}))^T + \sum_{i=1}^m g_i(\bar{x}) \nabla^2 g_i(\bar{x}) \cong \nabla \bar{g}(\bar{x}) (\nabla \bar{g}(\bar{x}))^T \\ \Rightarrow H(\bar{x}) &\cong \nabla \bar{g}(\bar{x}) (\nabla \bar{g}(\bar{x}))^T \end{aligned}$$

ושיטת ניוטון-גאוס היא למעשה האלגוריתם האיטרטיבי הבא (זו למעשה אלגוריתם ניוטון עם צעד בכיוון קצת שונה):

$$\bar{x}_{k+1} = \bar{x}_k - \alpha_k H^{-1}(\bar{x}_k) \overbrace{(\nabla \bar{g}(\bar{x}_k))}^{\text{Column Vector}} \overbrace{(\bar{g}(\bar{x}_k))}^{\text{Column Vector}}$$

אם מטריצת ההסיון היא סינגולרית (לא בעלת דרגה מלאה), כלומר אם ישנם ערכים עצמיים שהם אפס במטריצת ההסיון

$H(\bar{x}) \cong \nabla \bar{g}(\bar{x}) (\nabla \bar{g}(\bar{x}))^T$  (הם אינם יכולים להיות שלילים כי מצאנו כי המטריצה היא אי-שלילית) אז נוכל להוסיף

מטריצה אלכסונית  $\Delta$  ונקבל  $\nabla \bar{g}(\bar{x}) (\nabla \bar{g}(\bar{x}))^T + \Delta$  (ניתן להשתמש בשיטת Modified Cholesky factorization) ולהציב באלגוריתם האיטרטיבי:

$$\bar{x}_{k+1} = \bar{x}_k - \alpha_k (H(\bar{x}_k) + \Delta)^{-1} \overbrace{(\nabla \bar{g}(\bar{x}_k))}^{\text{Column Vector}} \overbrace{(\bar{g}(\bar{x}_k))}^{\text{Column Vector}}$$

באופן אלטרנטיבי ניתן להציב  $\Delta = \varepsilon_k I$  ולקבל:  $H(\bar{x}_k) + \varepsilon_k I$  ושיטה זו היא בעלת שם משל עצמה והיא נקראת

Levenberg-Marquart Method.

## הרצאה 10 – Conjugate Gradient Method. שיטת הווקטורים האורתוגונליים לגרדיאנט

למדנו שאלגוריתם ש"הולך" תמיד נגד כיוון הגרדיאנט, קצב התכנסות שלו (של האלגוריתם) תלוי בערך של condition number וכאשר הוא גדול אז קצב ההתכנסות הוא איטי ויקח הרבה זמן עד שהאלגוריתם יגיע לפתרון. לאחר מכן, מצאנו את שיטת ניוטון ששיפרה את האלגוריתם הקודם אך שיטת ניוטון צריכה לחשב את מטריצת ההסיאן בכל מהלך של האלגוריתם ואם בעית האופטימיזציה שלנו היא פונקציה של הרבה מאד משתנים אז מטריצת ההסיאן תהיה גדולה מאד ונדרש זיכרון עצום בשביל לשמור את המטריצה. כמו כן, בעיה נוספת בשיטת ניוטון זו סיבוכיות החישוב, אם הבעיה היא ב- $n$  מימדים, אז סיבוכיות

החישוב היא  $O(n^3)$  ואם נשתמש בשיטת Cholesky נגיע אמנם לסיבוכיות  $O\left(\frac{1}{6}n^3\right)$  אבל זו עדיין סיבוכיות די גדולה

ולמעשה פתרון של בעיה מסוימת יכול לקחת שעות של חישוב ע"י מחשב (כאשר  $n \approx 10,000$ ) וגם גודל מטריצת ההסיאן כל כך גדול שלא ניתן באופן מעשי לחשב את המטריצה ההופכית של ההסיאן.

אנחנו צריכים שיטה שתמצא פשרה בין פשטות האלגוריתם לבין סיבוכיות האלגוריתם. השיטה שנלמד כעת מנסה למצוא את הפשרה הנ"ל. תחילה נראה איך השיטה עובדת עבור בעיות אופטימיזציה עם פונקציות ריבועיות ולאחר מכן נעבור למקרה הכללי של פונקציות כלשהן.

השיטה Conjugate directions מתבססת בעיקרה על המכפלה הפנימית ועל התכונות של אורתוגונליות אך באופן יותר כללי מאשר מה שלמדנו במרחבים אוקלידיים. בשיטה זו אנחנו למעשה מנסים לקחת בעית אופטימיזציה ב- $n$  משתנים ולהפוך את הבעיה ל- $n$  בעיות אופטימיזציה של משתנה יחיד. נתחיל בהקדמה מתמטית.

הגדרה: מכפלה פנימית

נניח כי  $a, b, c \in \mathbb{R}^n$  ווקטורים ונניח כי  $\alpha \in \mathbb{C}$  סקלר. נגדיר את המכפלה הסקלרית בין הווקטורים  $a, b$  להיות:  $\langle a, b \rangle$ . תכונות של כל מכפלה פנימית:

- $\langle a, b \rangle = \overline{\langle b, a \rangle}$
- $\langle a, b + c \rangle = \langle a, b \rangle + \langle a, c \rangle$
- $\langle \alpha a, b \rangle = \alpha \langle a, b \rangle$
- $\langle a, a \rangle \geq 0$
- $\langle a, a \rangle = 0$  אם ורק אם  $a = \bar{0}$ .

הנורמה מוגדרת להיות:

$$\text{norm } a \triangleq \|a\| = \sqrt{\langle a, a \rangle}$$

דוגמאות למכפלה פנימית:

- במרחב  $\mathbb{R}^n$ : (המכפלה הפנימית הסטנדרטית)

$$\langle a, b \rangle \triangleq a^T \cdot b = \sum_{i=1}^n a_i b_i$$

- במרחב הפונקציות של משתנה יחיד (או רבים) שמוגדרות בתחום  $I = [a, b]$ : (המכפלה הפנימית הסטנדרטית)

$$\langle f(x), g(x) \rangle \triangleq \int_a^b f(x) \cdot \overline{g(x)} dx$$

• במרחב  $\mathbb{R}^n$ , נניח כי  $\bar{x}, \bar{y} \in \mathbb{R}^n$  וכן  $Q$  מטריצה חיובית מסדר  $n \times n$  אז ניתן להגדיר את המכפלה הפנימית הבאה:

$$\langle \bar{x}, \bar{y} \rangle_Q \triangleq \bar{x}^T Q \bar{y}$$

(יש להוכיח כי המכפלה הפנימית שהגדרנו היא אכן מכפלה פנימית באמצעות זה שנוכיח כי היא מקיימת את כל התכונות שראינו קודם שכל מכפלה פנימית חייבת לקיים – מומלץ לבדוק כשיעורי בית)

תהליך גרהם-שמידט:

תהליך גרהם-שמידט הוא תהליך שמקבל קבוצה של ווקטורים בלתי תלויים של מרחב בעל מכפלה פנימית ומחזיר קבוצה של ווקטורים אורתונורמליים של מרחב זה.

נניח כי קבוצת הווקטורים שהתהליך מקבל היא:  $\{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n\}$  אז הבסיס האורתונורמלי שנקבל ע"י התהליך הוא:

$$\begin{aligned} \bar{y}_1 &= \bar{x}_1 \\ \bar{y}_k &= \bar{x}_k - \sum_{i=1}^{k-1} \left\langle \bar{x}_i, \frac{\bar{y}_i}{\|\bar{y}_i\|} \right\rangle \cdot \frac{\bar{y}_i}{\|\bar{y}_i\|} = \bar{x}_k - \sum_{i=1}^{k-1} \frac{\langle \bar{x}_i, \bar{y}_i \rangle}{\|\bar{y}_i\|^2} \cdot \bar{y}_i \end{aligned}$$

ווקטורים אורתוגונלים ביחס למכפלה הפנימית עם המטריצה Q (Q-conjugate (Q-orthogonal) directions):

נניח כי  $\bar{d}_i, \bar{d}_j \in \mathbb{R}^n$  ווקטורים, אז נאמר כי הם אורתוגונליים לפי המכפלה הפנימית לפי המטריצה Q אם מתקיים:

$$\langle \bar{d}_i, \bar{d}_j \rangle_Q = \bar{d}_i^T Q \bar{d}_j = 0$$

אם נניח כי נתונה לנו קבוצה של ווקטורים בלתי תלויים לינארית  $\{\bar{\xi}_1, \bar{\xi}_2, \dots, \bar{\xi}_n\}$  אז ע"י תהליך גרהם-שמידט אפשר למצוא קבוצת ווקטורים אורתוגונליים  $\{\bar{d}_1, \bar{d}_2, \dots, \bar{d}_n\}$  באופן הבא:

$$\begin{aligned} \bar{d}_1 &= \bar{\xi}_1 \\ \bar{d}_k &= \bar{\xi}_k - \sum_{i=1}^{k-1} \left\langle \bar{\xi}_i, \frac{\bar{d}_i}{\|\bar{d}_i\|} \right\rangle \cdot \frac{\bar{d}_i}{\|\bar{d}_i\|} = \bar{\xi}_k - \sum_{i=1}^{k-1} \frac{\langle \bar{\xi}_i, \bar{d}_i \rangle}{\|\bar{d}_i\|^2} \cdot \bar{d}_i = \bar{\xi}_k - \sum_{i=1}^{k-1} \frac{\bar{\xi}_i^T Q \bar{d}_i}{\bar{d}_i^T Q \bar{d}_i} \cdot \bar{d}_i \end{aligned}$$

נראה בקרוב כי הווקטורים האורתוגונליים הללו מאפשרים לנו באמצעות שיטה פשוטה למדי למצוא את המינימום של פונקציה ריבועית.

מציאת מינימום של פונקציה ריבועית:

נניח כי  $Q$  מטריצה חיובית מסדר  $n \times n$  וכן ווקטור  $\bar{b} \in \mathbb{R}^n$  אז נוכל להגדיר פונקציה ריבועית  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  (התנאי שהמטריצה תהיה חיובית אומר שבהכרח לפונקציה הריבועית יש מינימום ולא מקסימום, בדיוק כמו פונקציה ריבועית ממימד אחד אשר יהיה לה מקסימום אם הנגזרת השניה תהיה חיובית) באופן הבא:

$$f(\bar{x}) = \bar{x}^T Q \bar{x} + \bar{b}^T \cdot \bar{x} = \frac{1}{2} \|\bar{x}\|_Q^2 + \bar{b}^T \cdot \bar{x}$$

בנוסף, נניח ובידינו יש קבוצה של  $n$  ווקטורים אורתוגונליים (בנינו ע"י תהליך גרהם-שמידט) לפי המטריצה  $Q$ :

$$\{\bar{d}_i\}_{i=1}^n = \{\bar{d}_1, \dots, \bar{d}_n\}$$

נוכל למצוא צירוף לינארי של  $\bar{x}$  באמצעות הבסיס שבנינו ונקבל:

$$\bar{x} = \sum_{i=1}^n \alpha_i \bar{d}_i, \quad \forall 1 \leq i \leq n: \alpha_i \in \mathbb{R}$$

(כלומר  $\bar{\alpha} \in \mathbb{R}^n$  הוא ווקטור ואז נוכל לומר כי  $\bar{x}$  הוא פונקציה של הווקטור  $(\bar{\alpha})$ )

וכעת נוכל לכתוב את הפונקציה הריבועית באופן הבא:

$$f(\bar{x}(\bar{\alpha})) = \frac{1}{2} \|\bar{x}(\bar{\alpha})\|_Q^2 + \bar{b}^T \cdot \bar{x} = \sum_{i=1}^n \left( \frac{1}{2} \alpha_i^2 \|\bar{d}_i\| + \alpha_i \cdot \bar{b}^T \bar{d}_i \right)$$

ואם נסמן  $\varphi_i(\alpha_i) \triangleq \frac{1}{2} \alpha_i^2 \|\bar{d}_i\| + \alpha_i \cdot \bar{b}^T \bar{d}_i$  אז נקבל את בעית האופטימיזציה הבאה:

$$\min_{\bar{\alpha} \in \mathbb{R}^n} f(\bar{x}(\bar{\alpha})) = \sum_{i=1}^n \min_{\alpha_i \in \mathbb{R}} \varphi_i(\alpha_i) = \min_{\alpha_1 \in \mathbb{R}} \varphi_1(\alpha_1) + \dots + \min_{\alpha_n \in \mathbb{R}} \varphi_n(\alpha_n)$$

כלומר הצלחנו לבטא את הפונקציה הריבועית שלנו בתור פונקציה פריקה (separable) לפי  $\alpha_i$  וזה אומר שהפכנו בעית אופטימיזציה אחת ב- $n$  משתנים ל- $n$  בעיות אופטימיזציה של משתנה אחד וזה טוב לנו, כי אנחנו רואים כי מציאת המינימום של בעיה אחת, בהכרח לא משפיע על בעיה אחרת (כי המטרה היא למזער את הסכום כולו בסופו של דבר). בשביל לפתור כל אחת מ- $n$  בעיות האופטימיזציה נוכל לבצע אלגוריתם איטרטיבי של אופטימיזציה במשתנה יחיד (למשל Line search).

תכונת היריעה של מרחבים אפיניים (Expanding Manifold Property):

באמצעות השיטה Conjugate directions method שראינו לעיל, ניתן לפתור בעיות גם במרחבים אפיניים, כלומר הנקודה  $\bar{x}$  יכולה להיות מתוארת גם באופן הבא:

$$\bar{x} = \bar{x}_0 + \sum_{i=1}^n \alpha_i \bar{d}_i, \quad \forall 1 \leq i \leq n: \alpha_i \in \mathbb{R}$$

ותהליך האופטימיזציה הוא למה שראינו לעיל.

לאחר שבידינו  $n$  בעיות האופטימיזציה, אם למשל נסתכל על מצב שבו פתרנו רק  $k < n$  בעיות אז אנחנו מקבלים את הפתרון החלקי הבא:

$$\bar{x} = \bar{x}_0 + \sum_{i=1}^k \alpha_i \bar{d}_i$$

הגדרה:

תת-מרחב אפיני  $G_k$  הוא המרחב:

$$G_k = \left\{ \bar{x} \mid \bar{x} = \bar{x}_0 + \sum_{i=1}^k \alpha_i \bar{d}_i \right\}, \quad \forall 1 \leq i \leq k: \alpha_i \in \mathbb{R}$$

כלומר  $G$  הוא מרחב לינארי מוזז ע"י הוקטור  $\sum_{i=1}^k \alpha_i \bar{d}_i$ .

תכונת היריעה של מרחבים אפיניים אומרת כי לאחר הצעד ה- $k$  בשיטה Conjugate directions method (כלומר פתרנו  $k$  בעיות אופטימיזציה מתוך ה- $n$ ) אנחנו נקבל את הפתרון החלקי האופטימלי הבא:

$$\bar{x}_k = \arg \min_{\bar{x} \in G_k} f(\bar{x})$$

כלומר בנקודה  $\bar{x}_k$  נקבל את הערך האופטימלי האפשרי לאחר  $k$  איטרציות באלגוריתם.

## הרצאה 11 – Conjugate Gradient Method 2, שיטת הווקטורים האורתוגונליים לגרדיאנט - המשך

ראינו בהרצאה הקודמת כי ניתן בצורה פשוטה למצוא את המינימום של הפונקציה הריבועית  $f(\bar{x}) = \bar{x}^T Q \bar{x} + \bar{b}^T \bar{x}$  באמצעות זה שנמצא את הווקטורים האורתוגונליים לפי המטריצה החיובית  $Q$ .

נסמן את הגרדיאנט של הפונקציה הריבועית בנקודה  $\bar{x}_k$  להיות:

$$\bar{g}_k \triangleq \bar{g}(\bar{x}_k) = \nabla f(\bar{x}_k) = Q\bar{x}_k + \bar{b}$$

אם בהרצאה הקודמת הנחנו שבידינו קבוצת ווקטורים אורתוגונליים  $\{\bar{d}_i\}_{i=0}^n$  שיצרנו באמצעות תהליך גרהם-שמידט, כעת נרצה

לבטא את הווקטורים האורתוגונליים  $\{\bar{d}_i\}_{i=0}^n$  באמצעות הגרדיאנטים  $\bar{g}_k$  ולראות איך השיטה Conjugate Gradient נראית באופן הזה. נבחר את הווקטור הראשון:

$$\bar{d}_0 = -\bar{g}_0 = -(Q\bar{x}_0 + \bar{b})$$

הצעד ה- $k$ -י באלגוריתם הוא:

$$\bar{x}_{k+1} = \bar{x}_k + \gamma_k \bar{d}_k, \quad \gamma_k \in \mathbb{R}$$

ואת הערך  $\gamma_k$  נמצא בשיטת Exact line Search (בשיטת Conjugate direction חייבים למצוא את  $\gamma_k$  בשיטה זו בלבד כיוון שזו השיטה הכי מדויקת ובמקרה זה השיטה Exact Line Search הכרחית). ואז נמצא את הווקטור האורתוגונלי הבא  $\bar{d}_{k+1}$  באמצעות תהליך גרהם-שמידט:

$$\bar{d}_{k+1} = -\bar{g}_{k+1} + \sum_{j=0}^k \frac{\bar{g}_{k+1}^T Q \bar{d}_j}{\bar{d}_j^T Q \bar{d}_j} \bar{d}_j$$

אחת היתרונות העיקריים של השיטה Conjugate Gradient על פני השיטה הכללית יותר שהיא Conjugate Direction נעוץ ברישום של הווקטורים האורתוגונליים  $\bar{d}_k$ , כאשר רושמים את הווקטורים באמצעות הגרדיאנט זה נהיה הרבה יותר פשוט וקומפקטי.

נבחין כי הצעד באלגוריתם הוא  $\bar{x}_{k+1} = \bar{x}_k + \gamma_k \bar{d}_k$  ולכן מתקיים:

$$\bar{d}_j = \frac{1}{\gamma_j} (\bar{x}_{j+1} - \bar{x}_j)$$

ומכאן נקבל:

$$Q \bar{d}_j = \frac{1}{\gamma_j} Q (\bar{x}_{j+1} - \bar{x}_j) = \frac{1}{\gamma_j} (\bar{g}_{j+1} - \bar{b} - (\bar{g}_j - \bar{b})) = \frac{1}{\gamma_j} (\bar{g}_{j+1} - \bar{g}_j)$$



בנוסף ניזכר כי לפי תכונת היריעה של מרחבים אפיניים אנחנו יודעים כי בצעד ה- $k$ -י באלגוריתם, אנחנו מגיעים לפתרון האופטימלי האפשרי  $\bar{x}_k$  בתת המרחב האפיני  $G_k$  (המרחב שנפרש ע"י כל הווקטורים שמצאנו עד לאותו שלב באלגוריתם) וכיוון שבשיטה Conjugate Gradient אנחנו גורמים לזה שכל ווקטור שאנחנו יוצרים הוא צירוף של גרדיאנטים קודמים אז זה אומר שבכל שלב באלגוריתם אנחנו מוציאים את הפתרון האופטימלי במרחב נפרש ע"י הגרדיאנטים הקודמים שמצאנו, וכמו שראינו בהרצאה הקודמת, כל ווקטור חדש שאנחנו מוצאים הוא אורתוגונלי לכל הקודמים ולכן גם הגרדיאנט החדש שנמצא בצעד הבא הוא אורתוגונלי לכל הגרדיאנטים האחרים ולמרחב כולו שנפרש על ידם, כלומר באופן פורמלי נקבל:

$$\bar{g}_{k+1} \perp \{\bar{g}_k, \bar{g}_{k-1}, \dots, \bar{g}_0\}$$

כעת מהחישוב שעשינו עבור  $Q\bar{d}_j$  ומהעובדה שמצאנו כי מתקיים  $\bar{g}_{k+1} \perp \{\bar{g}_k, \bar{g}_{k-1}, \dots, \bar{g}_0\}$  אנחנו מסיקים כי מתוך כל הסכום בביטוי:

$$\bar{d}_{k+1} = -\bar{g}_{k+1} + \sum_{j=0}^k \frac{\bar{g}_{k+1}^T Q\bar{d}_j}{\bar{d}_j^T Q\bar{d}_j} \bar{d}_j$$

נשאר רק האיבר האחרון (של הסכום) כי הגרדיאנט  $\bar{g}_{k+1}$  אורתוגונלי לכל הגרדיאנטים הקודמים (לאחר הצבת החישוב שמצאנו עבור  $Q\bar{d}_j$ ) ולכן המונה מתאפס כי הוא מכפלה סקלרית בין ווקטורים אורתוגנליים ולכן נקבל:

$$\bar{d}_{k+1} = -\bar{g}_{k+1} + \frac{\bar{g}_{k+1}^T Q\bar{d}_k}{\bar{d}_k^T Q\bar{d}_k} \bar{d}_k = -\bar{g}_{k+1} + \frac{\bar{g}_{k+1}^T (\bar{g}_{k+1} - \bar{g}_k)}{\bar{d}_k^T (\bar{g}_{k+1} - \bar{g}_k)} \bar{d}_k$$

$Q\bar{d}_k = \frac{1}{\gamma_k}(\bar{g}_{k+1} - \bar{g}_k)$

ואם נסמן את הסקלר שמוכפל בווקטור  $\bar{d}_k$  להיות  $\beta_k$  נקבל את הביטוי הפשוט הבא:

$$\bar{d}_{k+1} = -\bar{g}_{k+1} + \beta_k \bar{d}_k$$

ננסה לפשט את הביטוי עבור  $\beta_k$ . במכנה נבחין כי  $\bar{g}_{k+1} \perp \bar{d}_k^T$  ולכן המכפלה הסקלרית ביניהם מתאפסת ולכן המכנה נותר רק  $(-\bar{g}_k) \cdot \bar{d}_k^T$  ואם נציב במפורש את  $\bar{d}_k$  מתוך האיטרציה הקודמת באלגוריתם נקבל:

$$\bar{d}_k = -\bar{g}_k + \underbrace{\frac{\bar{g}_k^T (\bar{g}_j - \bar{g}_{j-1})}{\bar{d}_{k-1}^T (\bar{g}_j - \bar{g}_{j-1})}}_{\beta_{k-1}} \bar{d}_{k-1} = -\bar{g}_k + \beta_{k-1} \bar{d}_{k-1}$$

(נציב את הביטוי שקיבלנו זה עתה עבור  $\bar{d}_k$  במכנה) ולכן המכנה הוא:

$$(-\bar{g}_k + \beta_{k-1} \bar{d}_{k-1})^T (-\bar{g}_k)$$

וכיוון שמתקיים  $\bar{d}_{k-1} \perp \bar{g}_k$  (כי מתקיים  $\bar{d}_{k-1}^T \cdot \bar{g}_k = 0$ ) נשארו במכנה עם הסקלר:

$$(-\bar{g}_k^T)(-\bar{g}_k) = \|\bar{g}_k\|_2^2$$

כלומר הצלחנו לפשט את המכנה של  $\beta_k$  וקיבלנו:

$$\beta_k = \frac{\bar{g}_{k+1}^T (\bar{g}_{k+1} - \bar{g}_k)}{\|\bar{g}_k\|_2^2}$$

הנוסחה האחרונה פותחה ע"י Polak-Ribiere ונקראת על שמם, וכן אם נבחין כי מתקיים גם  $\bar{g}_k \perp \bar{g}_{k+1}$  אז נוכל לרשום את הנוסחה האחרונה גם באופן הבא:

$$\beta_k = \frac{\|\bar{g}_{k+1}\|_2^2}{\|\bar{g}_k\|_2^2}$$

ונוסחה זו קרויה על שם Fletcher-Reeves.

שתי הנוסחאות האלה במקרה שלנו מתלכדות אך כאשר נשתמש בשיטה Conjugate Gradient עבור פונקציה כללית (ולא דווקא ריבועית) אנחנו נראה כי העובדה  $\bar{g}_k \perp \bar{g}_{k+1}$  לא מתקיימת ולכן שתי הנוסחאות האחרונות שונות.

כמו כן, עבור חישובים נומריים, דווקא הנוסחה הראשונה מתנהגת טוב יותר ולכן עבור חישובים נומריים לרוב משתמשים בנוסחה של Polak-Ribiere.

סיכום השיטה Conjugate Gradient:

עבור הפונקציה הריבועית  $f(\bar{x}) = \frac{1}{2} \bar{x}^T Q \bar{x} + \bar{b}^T \bar{x}$  השיטה היא כלהלן:

(1) מתחילים לבנות את קבוצת הווקטורים האורתוגונליים לפי הווקטור הראשון:

$$\bar{d}_0 = -\bar{g}_0 = \nabla f(\bar{x}_0) = Q\bar{x}_0 + \bar{b}$$

(2) בכל שלב באלגוריתם מחשבים את הנקודה הבאה לפי הנוסחה:  $\bar{x}_{k+1} = \bar{x}_k + \gamma_k \bar{d}_k$  ואת הסקלר  $\gamma_k$  מוצאים בצורה אנליטית (כיוון שזו פונקציה ריבועית אז ניתן למצוא באופן אנליטי) לפי הנוסחה:

$$\gamma_k = -\frac{\bar{g}_k^T \bar{d}_k}{\bar{d}_k^T Q \bar{d}_k} = -\frac{f'_{\bar{d}_k}(\bar{x}_k)}{f''_{\bar{d}_k \bar{d}_k}(\bar{x}_k)}$$

(זו נוסחה דומה לנוסחה של מציאת מינימום של פונקציה ריבועית במשתנה יחיד – הוכיחו זאת לעצמכם) ניתן גם להגיע לפתרון זה באמצעות איטרציה אחת (כיוון שמדובר בפונקציה ריבועית) בשיטת ניוטון.

(3) את הגרדיאנט הבא מוצאים לפי  $\bar{g}_{k+1} = Q\bar{x}_{k+1} + \bar{b}$  ונבחין כי ניתן לחשב  $\bar{g}_{k+1}$  ללא חישוב נוסף כיוון

שמתקיים  $Q\bar{x}_{k+1} = Q\bar{x}_k + \gamma_k Q\bar{d}_k$  ונבחין כי את הביטוי  $Q\bar{d}_k$  כבר חישבנו כאשר חישבנו את  $\gamma_k$ .

(4) מוצאים את ווקטור הכיוון הבא, הווקטור האורתוגנולי הבא לפי הנוסחה  $\bar{d}_{k+1} = -\bar{g}_{k+1} + \beta_k \bar{d}_k$  כאשר מתקיים (לא רק עבור פונקציה ריבועית):

$$\beta_k = \frac{\bar{g}_{k+1}^T (\bar{g}_{k+1} - \bar{g}_k)}{\|\bar{g}_k\|_2^2}, \quad \text{Polak - Ribiere}$$

$$\beta_k = \frac{\|\bar{g}_{k+1}\|_2^2}{\|\bar{g}_k\|_2^2}, \quad \text{Fletcher - Reeves}$$

עבור הפונקציה הכללית  $f(\bar{x})$  השיטה היא כלהלן:

(1) מתחילים לבנות את קבוצת הווקטורים האורתוגונליים לפי הווקטור הראשון:

$$\bar{d}_0 = -\bar{g}_0 = \nabla f(\bar{x}_0)$$

(2) בכל שלב באלגוריתם מחשבים את הנקודה הבאה לפי הנוסחה:  $\bar{x}_{k+1} = \bar{x}_k + \gamma_k \bar{d}_k$  ואת הסקלר  $\gamma_k$  מוצאים ע"י שיטת Line Search Exact (הסברנו כי בשיטת Conjugate Gradient חייבים להשתמש בשיטה זו).

(3) מוצאים את הגרדיאנט הבא:  $\bar{g}_{k+1} = \nabla f(\bar{x}_{k+1})$ .

(4) מוצאים את ווקטור הכיוון הבא, הווקטור האורתוגונלי הבא לפי הנוסחה  $\bar{d}_{k+1} = -\bar{g}_{k+1} + \beta_k \bar{d}_k$  כאשר מתקיים (לא רק עבור פונקציה ריבועית):

$$\beta_k = \frac{\bar{g}_{k+1}^T (\bar{g}_{k+1} - \bar{g}_k)}{\|\bar{g}_k\|_2^2}, \quad \text{Polak - Ribiere}$$

$$\beta_k = \frac{\|\bar{g}_{k+1}\|_2^2}{\|\bar{g}_k\|_2^2}, \quad \text{Fletcher - Reeves}$$

#### קצב ההתכנסות של השיטה Conjugate Gradient:

כמו שהראנו כי בשיטת Conjugate Direction, לאחר בדיוק  $n$  איטרציות של האלגוריתם, כאשר  $n$  זה מספר המשתנים של בעיית האופטימיזציה, נגיע לפתרון האופטימלי אז גם בשיטת Conjugate Gradient עבור פונקציה ריבועית האלגוריתם יתכנס לפתרון האופטימלי לאחר  $n$  איטרציות כאשר  $n$  זה מספר המשתנים של בעיית האופטימיזציה. בבעיות מסוימות, אשר בהן יש מאות אלפי משתנים, לפעמים עוצרים את האלגוריתם עוד לפני שהוא ביצע את כל האיטרציות הדרושות.

נסמן את הערך העצמי המקסימלי של מטריצת ההסיון להיות  $M = \lambda_{\max}$  ואת הערך העצמי המינימלי להיות  $m = \lambda_{\min}$ .

נסמן את  $\bar{x}^*$  להיות הערך עבורו הפונקציה שעבורה אנחנו מחפשים את המינימום מקבלת אכן את המינימום אזי ניתן לתאר את קצב ההתכנסות של האלגוריתם Conjugate Gradient ע"י קצב ההתקדמות הבא:

$$\|\bar{x}_{k+1} - \bar{x}^*\| = c \|\bar{x}_k - \bar{x}^*\| = \frac{\sqrt{M} - \sqrt{m}}{\sqrt{M} + \sqrt{m}} \|\bar{x}_k - \bar{x}^*\|$$

בנוסף, כאשר מתקיים  $M \gg m$  אפשר להניח כי מתקיים:  $c = 1 - 2\sqrt{\frac{m}{M}}$  (נוכל גם להיזכר כי היחס  $\theta = \frac{\lambda_{\max}}{\lambda_{\min}}$

נקרא גם Condition Number ובסימונים כעת מתקיים:  $c = 1 - 2\sqrt{\frac{1}{\theta}}$  (ואם נשווה את קצב ההתכנסות של

האלגוריתם Conjugate Gradient לקצב ההתכנסות של האלגוריתם Steepest Descent (או בשמו האחר Gradient Descent) אז מתקיים קצב ההתכנסות של Steepest Descent הוא:

$$f(\bar{x}_{k+1}) - f(\bar{x}^*) \leq c (f(\bar{x}_k) - f(\bar{x}^*))$$

ומתקיים:

$$c = 1 - 2\sqrt{\frac{m}{M}}, \quad M = \lambda_{\max}, \quad m = \lambda_{\min}$$

וכמובן נזכור כי ככל שהקצב  $c$  קרוב יותר לאפס אז קצב ההתכנסות מהיר יותר.

דוגמה:

אם נניח כי  $\theta = \frac{M}{m} = 10^6$  (כאשר  $\theta$  זה ה-condition number) אז נקבל:

$$c_{\text{Conjugate Gradient}} = 1 - 2\sqrt{\frac{m}{M}} = 0.998$$

$$c_{\text{Gradient Descent}} = 1 - 2\frac{m}{M} = 0.999998$$

כלומר צריך לבצע בערך 1000 צעדים של שיטת Gradient Descent על מנת לקבל את התוצאות של שיטת Conjugate Gradient – כלומר שיטת Conjugate Gradient בדוגמה זו יעילה פי 1000 מהשיטה Gradient Descent!

שיטה להקטנת הפרמטר (Preconditioning) Condition number:

ראינו כי גם בשיטת Conjugate Gradient וגם בשיטת Gradient Descent (Steepest Descent) קצב ההתכנסות מושפע

מהפרמטר  $\theta = \frac{\lambda_{\max}}{\lambda_{\min}}$  (condition number) וראינו כי כאשר הפרמטר  $\theta$  הוא קטן, אז מקבלים קצב התכנסות מהיר ולכן נרצה למצוא דרכים כיצד להקטין את הפרמטר  $\theta$ .

מקרה פרטי של שיטת Preconditioning:

נניח ואנחנו מחפשים את המינימום של פונקציה ריבועית:

$$f(\bar{x}) = \frac{1}{2} \bar{x}^T Q \bar{x} + \bar{b}^T \bar{x}$$

כאשר מניחים כי מטריצה  $Q$  היא חיובית  $Q \succ 0$  ולכן גם סימטרית.

נחליף משתנים בפונקציה באופן הבא: נניח כי ברשותנו מטריצה הפיכה  $S$  אז נסמן:

$$\bar{x} = S\bar{y}$$

ולכן קיבלנו את החלפת המשתנים:

$$\varphi(\bar{y}) = f(S\bar{y}) = \frac{1}{2} \bar{y}^T S^T Q S \bar{y} + \bar{b}^T S \bar{y}$$

כלומר במקום המטריצה  $Q$  יש לנו את המטריצה  $S^T Q S$ . החלפת משתנים זו יעילה לנו כאשר אנחנו יודעים כי ה-condition number של המטריצה  $S^T Q S$  הוא קטן מה-condition number של המטריצה  $Q$  ועכשיו אפשר להשתמש בשיטה Conjugate Gradient בצורה יעילה ומהירה יותר כי מובטח לנו קצב התכנסות מהיר יותר מאשר אם היינו מפעילים את השיטה על הפונקציה המקורית עם המטריצה  $Q$ .

באופן מעשי, פעמים רבות אנחנו נתקלים במטריצה  $Q$  ישנם ערכים גדולים מאד וקטנים מאד על האלכסון הראשי (מה שמצביע גם על ערכים עצמיים גדולים מאד וקטנים מאד) ולכן לאחר שנכפול במטריצה  $S$  מתאימה נוכל להקטין את ההפרש בין הערכים על האלכסון ואז הערכים העצמיים יהיו קרובים יותר אחד לשני וכתוצאה מזה גם ה-condition number יקטן ואז קצב ההתכנסות יהיה מהיר יותר.

דוגמה:

נניח כי המטריצה  $Q$  היא כמעט מטריצה אלכסונית (כלומר רוב הערכים נמצאים על האלכסון וערכים מעטים נמצאים במקומות אחרים) ונסמן כי המטריצה  $Q' = \text{diagonal } Q$  היא המטריצה  $Q$  כאשר מאפסים את כל הערכים שהם לא על האלכסון הראשי אז המטריצה  $S$  הדרושה היא:

$$S = (\sqrt{Q'})^{-1} = \begin{pmatrix} \frac{1}{\sqrt{q_{11}}} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{\sqrt{q_{nn}}} \end{pmatrix} \Rightarrow S = S^T$$

ולכן נקבל:

$$S^T Q' S = I \Rightarrow \lambda_{\max} = \lambda_{\min} = 1 \Rightarrow \theta = 1$$

הערה: במקרה הכללי ביותר, אם מניחים כי  $Q$  היא סימטרית אז אם נוכל לחשב את  $S = (\sqrt{Q})^{-1}$  ואז שוב נקבל  $\theta = 1$ . שיטה זו, שיטה של Preconditioning יעילה מאד אך ניתן להשתמש בה רק במקרים יחסית ספציפיים.

הכללה של שיטת Preconditioning:

ביטאנו את השיטה (Preconditioning) באמצעות החלפת משתנים, וכעת נבטא את השיטה בצורתה המקורית: נסמן  $W = S^T S$ . בשיטת Conjugate Gradient הכיוון הראשון אליו מתקדמים באלגוריתם עבור פונקציה ריבועית נקבע ע"י:

$$\bar{d}_0 = W \cdot \bar{g}_0$$

כאשר הגרדיאנט של פונקציה ריבועית הוא כפי שראינו כבר בעבר:

$$f(\bar{x}) = \frac{1}{2} \bar{x}^T Q \bar{x} + \bar{b}^T \bar{x} \Rightarrow \bar{g}(\bar{x}) = Q \bar{x} + \bar{b}$$

הצעד הבא הוא:

$$\bar{x}_{k+1} = \bar{x}_k + \gamma_k \bar{d}_k$$

הכיוון הבא הוא:

$$\bar{d}_{k+1} = W \cdot \bar{g}_{k+1} + \beta_k \bar{d}_k$$

כאשר:

$$\beta_k = \frac{\bar{g}_{k+1}^T \cdot W \cdot \bar{g}_{k+1}}{\bar{g}_k^T \cdot W \cdot \bar{g}_k}$$

בשיטה זו אנחנו לא משתמשים בצורה מפורשת במטריצה  $S$  אלא במטריצה  $W = S^T S$  שלעיתים חוסך חישובים מיותרים.

שיטת הקטימה של ניוטון (Truncated Newton's Method):

בשיטת ניוטון למדנו כיצד למצוא את המינימום של פונקציות ריבועיות (או פונקציות שאנחנו רוצים לקרב באמצעות פונקציה ריבועית) במספר רב של משתנים בצורה יעילה. ניזכר כי את כיוון הצעד הבא מצאנו לפי מערכת ניוטון (משוואת ניוטון):

$$\begin{cases} H(\bar{x}_k) \bar{d}_k = -\bar{g}(\bar{x}_k) \\ \min_{\bar{d}_k} \nabla q(\bar{d}_k) = \min_{\bar{d}_k} \bar{g}^T(\bar{x}_k) \bar{d}_k + \frac{1}{2} \bar{d}_k^T H(\bar{x}_k) \bar{d}_k \end{cases}$$

(ניזכר כי  $q(\bar{d}_k)$  הוא פיתוח טיילור מסדר שני (ולכן היא פונקציה ריבועית) של הפונקציה המקורית  $f(\bar{x})$  שאנחנו מנסים למזער)

וראינו כי בעיית מציאת המינימום של פונקציה ריבועית שקולה לפתירת מערכת משוואות כפי שנראה שוב מיד.

טענה:

המשוואה  $H(\bar{x}_k) \bar{d}_k = -\bar{g}(\bar{x}_k)$  (שנותנת לנו פתרון עבור  $\bar{d}_k$ , כי  $\bar{d}_k = -\frac{\bar{g}(\bar{x}_k)}{H(\bar{x}_k)}$ ) היא תנאי מספיק למינימום של

הביטוי  $\nabla q(\bar{d}_k)$ . כלומר אם נציב  $\bar{d}_k = -\frac{\bar{g}(\bar{x}_k)}{H(\bar{x}_k)}$  נקבל את המינימום של הביטוי  $\nabla q(\bar{d}_k)$ .

הוכחה:

אם נחשב את הגרדיאנט של הביטוי  $\nabla q(\bar{d}_k) = \frac{1}{2} \bar{d}_k^T H(\bar{x}_k) \bar{d}_k + \bar{g}^T(\bar{x}_k) \bar{d}_k$  ביחס לווקטור  $\bar{d}_k$  נקבל כי הגרדיאנט מתאפס עבור  $\bar{d}_k = -\frac{\bar{g}(\bar{x}_k)}{H(\bar{x}_k)}$  ולכן הוכחנו את הנדרש. מש"ל.

נסביר את הצורך בשיטת הקטימה. לעיתים אנחנו נדרשים לפתור בעיית אופטימיזציה ובמהלכה אנחנו נדרשים למצוא את כיוון ההתקדמות, כלומר לפתור את משוואת ניוטון:  $H(\bar{x}_k) \bar{d}_k = -\bar{g}(\bar{x}_k)$ , אך פתירת משוואה זו היא יקרה (מבחינת זמן) ולכן לעיתים נעדיף לפתור את הבעיה השקולה:

$$\min_{\bar{d}_k} \frac{1}{2} \bar{d}_k^T H(\bar{x}_k) \bar{d}_k + \bar{g}^T(\bar{x}_k) \bar{d}_k$$

ואת הבעיה הזו נפתור באמצעות שיטת Conjugate Gradient עד השלב שבו נקבל קירוב מספיק טוב ולא באופן מלא כדי לחסוך זמן בחישובים ולכן שיטה זו נקראת Truncated Newton's Method.

אנחנו למעשה רוצים למצוא  $\bar{d}_k$  אשר יקיים:  $H(\bar{x}_k) \bar{d}_k = -\bar{g}(\bar{x}_k)$  או במילים אחרות יקיים  $H(\bar{x}_k) \bar{d}_k + \bar{g}(\bar{x}_k) = 0$ . כפי שאמרנו, אנחנו לא נמצא את  $\bar{d}_k$  שאכן מביא למינימום את הביטוי  $\frac{1}{2} \bar{d}_k^T H(\bar{x}_k) \bar{d}_k + \bar{g}^T(\bar{x}_k) \bar{d}_k$  (ולכן גם פותר את משוואת ניוטון) אלא נמצא בקירוב, אך מה זה קירוב טוב? אנחנו נעצור כאשר יתקיים:

$$\|H(\bar{x}_k) \bar{d}_k + \bar{g}(\bar{x}_k)\| \leq \sigma \|H(\bar{x}_0) \bar{d}_0 + \bar{g}(\bar{x}_0)\|$$

כאשר  $\sigma$  הוא קבוע שנקבע לעצמנו. באופן, הזה, אנחנו מצאנו את כיוון הצעד המקורב שיש לעשות באלגוריתם על מנת לפתור את בעיית האופטימיזציה הכללית יותר שאנחנו מנסים לפתור.

הערה (זמן החישוב של אלגוריתמי אופטימיזציה):

בזמן פתירת בעיית אופטימיזציה למעשה אנחנו נדרשים לבצע הרבה פעולות חישוב. למשל, בשיטת ניטון אנחנו נדרשים הרבה פעמים לכפול את מטריצת ההסיאן בוקטור הכיוון ומתברר כי אם אנחנו נמצאים בשלב כלשהו באלגוריתם שפותר בעיית אופטימיזציה עם סדר גודל גדול מאד של משתנים (למשל מאות אלפי משתנים), אז כל החישובים הבאים לוקחים אותו סדר גודל של מאמץ חישובי:

- חישוב של מטריצת ההסיאן כפול ווקטור הכיוון
- חישוב ערך הגרדיאנט של הפונקציה המקורית
- חישוב ערך הפונקציה עצמה בנקודה שבה אנחנו עומדים

כל החישובים לעיל לוקחים את אותה כמות הזמן לערך (עד כדי קבוע כפלי קטן יחסית).

בנוסף, בשביל ליעל את החישוב (למשל במטלב) אם למשל  $H = A^T B A$  (ההסיאן הוא תוצאה של מכפלה של מטריצות) ואנחנו רוצים לחשב  $H \bar{x}$  רצוי לרשום במטלב את השורה הבאה:

$$Answer = A^T (B(A\bar{x}))$$

במקום:

$$Answer = A^T B A \bar{x}$$

כי באפשרות השניה מטלב יחשב תחילה כפל מטריצות ולא יזהה שהוא יכול לבצע רק כפל ווקטור במטריצה במקום.

בנוסף, נבחין כי למדנו כי מתקיים בצורה הדיפרנציאלית  $H(\bar{x}) d\bar{x} = d\bar{g}(\bar{x})$ , אז אם נסמן  $\varepsilon \bar{x}$  אז נוכל לומר כי מתקיים:

$$H(\bar{x}) \bar{x} \approx \frac{\bar{g}(\bar{x} + \varepsilon \bar{x}) - \bar{g}(\bar{x})}{\varepsilon}$$

כלומר מצאנו שיטה לחשב באופן מקורב את התוצאה של מכפלת הסיאן בוקטור כלשהו ע"י הגרדיאנט ולכן זה יקח פחות כוח חישוב.

## הרצאה 12 – Sequential Subspace Optimization (SESOP) and Quasi-Newton's Method

גם שיטת SESOP (שהיא היא הרחבה של שיטת Conjugate Gradient) וגם השיטה Quasi-Newton's Method שתיהן שיטות שמיועדות לבעיות אופטימיזציה שהפונקציה בהן היא פונקציה של מספר רב של משתנים.

שיטת Sequential Subspace Optimization (SESOP):

שיטה זו, היא כאמור הרחבה של שיטת Conjugate Gradient שהיא שיטה איטרטיבית למציאת הנקודה  $\bar{x}$  עבור הפונקציה שאנחנו מנסים למצוא את המינימום שלה,  $f(\bar{x})$ , מקבלת את המינימום הזה. כעת, נניח כי  $f(\bar{x})$  היא חלקה (לפחות גזירה ברציפות) ובעית האופטימיזציה היא כרגיל:

$$\bar{x} = \arg \min_{x \in \mathbb{R}^n} f(\bar{x}), \quad \bar{x} \in \mathbb{R}^n$$

כאשר הגרדיאנט של הפונקציה  $f(\bar{x})$  בזמן האיטרציה ה- $k$  הוא  $\bar{g}_k(\bar{x}_k) = \nabla f(\bar{x}_k)$ .

בנוסף נזכיר כי הכיוון (שהלכנו בו על מנת להגיע לנקודה  $\bar{x}_k$ ) (ללא התחשבות בגודל הצעד) הוא:

$$\bar{d}_{k-1} = \bar{x}_k - \bar{x}_{k-1}$$

(למה זה הביטוי עבור  $\bar{d}_{k-1}$ ? זה הביטוי לכיוון שהכיוון הקודם ניתן לביטוי גם כחיבור, או במקרה שלנו כחיסור של שני הווקטורים של הצעדים שביצענו באלגוריתם)

בשיטה SESOP אנחנו מצמצמים את מרחב החיפוש של הנקודה הבאה  $\bar{x}_{k+1}$  באמצעות שני הווקטורים:  $\bar{d}_{k-1}$  ו- $\bar{g}_k(\bar{x}_k)$ , כלומר אנחנו מחפשים את הנקודה  $\bar{x}_{k+1}$  בתת המרחב האפייני שנפרש ע"י הגרדיאנט בנקודה הנוכחית שבה אנחנו נמצאים וע"י ווקטור הכיוון שעל ידו הגענו לנקודה הנוכחית.

תחילה נבנה את המטריצה (באיטרציה ה- $k$  של האלגוריתם Conjugate Gradient) שעמודותיה הן:

$$P_k = (\bar{g}_k(\bar{x}_k) \quad \bar{d}_{k-1} \quad \bar{d}_{k-2} \quad \dots)$$

(כלומר עמודה ראשונה זה ווקטור הגרדיאנט  $\bar{g}_k(\bar{x}_k)$ , עמודה שניה זה ווקטור הכיוון  $\bar{d}_{k-1}$  וניתן (אך לא חובה) להשתמש גם בווקטורי הכיוון או בגרדיאנטים שמצאנו עד השלב ה- $k$  שאנחנו נמצאים בו עכשיו, המטריצה  $P_k$  חייבת להכיל לכל הפחות שתי עמודות).

האלגוריתם SESOP באיטרציה ה- $k$ :

(1) נמצא את כיוון וגודל ההתקדמות  $\bar{\alpha}_k$  מהנקודה  $\bar{x}_k$ , ע"י לפתור את תת-בעית האופטימיזציה:

$$\alpha_k = \arg \min_{\bar{\alpha}} f(\bar{x}_k + P \cdot \bar{\alpha})$$

(2) והצעד הבא יהיה:

$$\bar{x}_{k+1} = \bar{x}_k + P \cdot \alpha_k$$



## הערה:

ניתן לחשוב כי כמות הפעולות שמתבצעת באיטרציה יחידה של אלגוריתם SESOP היא משמעותית גדולה מזו של Conjugate Gradient אך לרוב מבצעים את האלגוריתם הזה על פונקציות מסוימות (שנראה מיד) ובפונקציות אלו (שהן די נפוצות) אלגוריתם זה (שמשמש בעיקרון ה-Subspace) בכל זאת מאד יעיל.

## שימוש בשיטת Fast Subspace Optimization:

ביישומים רבים הפונקציה שעבורה אנחנו מחפשים את המינימום,  $f(\bar{x})$ , היא מהצורה:

$$f(\bar{x}) = \varphi(A\bar{x})$$

כאשר ניתן לחשב את הפונקציה  $\varphi$  יחסית מהר ("זולה לחישוב") והמטריצה  $A$  זו מטריצה שמתפקדת למעשה כהעתקה לינארית והכפל שלה בווקטור זו הפעולה האיטית ("יקרה לחישוב"). נראה כי בשיטת Subspace Optimization אין צורך בחישובים מיותרים הכוללים את המטריצה  $A$  ולכן זו שיטה מהירה במקרה זה.

בשיטת Subspace Optimization אנחנו למעשה מחפשים את גדול וכיוון הצעד הבא,  $\alpha_k$ , באופן הבא:

$$\alpha_k = \arg \min_{\bar{\alpha}} f(\bar{x}_k + P \cdot \bar{\alpha}) = \arg \min_{\bar{\alpha}} \varphi(A(\bar{x}_k + P \cdot \bar{\alpha})) = \arg \min_{\bar{\alpha}} \varphi(A\bar{x}_k + AP \cdot \bar{\alpha})$$

ובשביל לפתור בעיה זו אנחנו צריכים לחשב את הפונקציה  $\varphi$  בהרבה נקודות שונות כתלות ב- $\alpha$ . נסמן  $R = AP$  ונבחין שמטריצה  $R$  היא בעלת מימדים יחסית קטנים (ביחס ל- $A$ ) (כי המימדים תלויים במטריצה  $P$  וראינו כי זו יכולה להיות גם מטריצה עם שתי עמודות בלבד) ולכן ניתן לחשב אותה במהירות, כלומר בזמן החישוב אנחנו נחשב פעם אחת את  $A\bar{x}_k$  ופעם אחת את  $R = AP$  ולאחר מכן נבצע רק חישוב מסוג  $R \cdot \bar{\alpha}$ . בנוסף המטריצה  $P$  משנה לכל היותר עמודה אחת מאיטרציה לאיטרציה ואין צורך לחשב את כולה כל פעם מחדש ולכן בכל איטרציה אנחנו מחשבים עמודה אחת במטריצה  $R$ . בנוסף נבחין כי באיטרציה הבאה אנחנו צריכים לדעת את ערכו של הביטוי  $A \cdot \bar{x}_{k+1}$  ולכאורה זה יגזול מאיתנו זמן חישוב גדול אך נבחין כי מתקיים:

$$A \cdot \bar{x}_{k+1} = A\bar{x}_k + R \cdot \bar{\alpha}_k$$

(ע"י כפל של מטריצה  $A$  משמאל במשוואה שראינו קודם:  $\bar{x}_{k+1} = \bar{x}_k + P \cdot \alpha_k$ )

ולסיכום, בכל איטרציה של האלגוריתם לפי Subspace Optimization אנחנו צריכים לחשב את הפונקציה  $\varphi$  מספר פעמים אך בכל פעם אנחנו צריכים לחשב בסך הכול עמודה אחת במטריצה  $AP$ .

נתבונן בנוסף בגרדיאנט ובהסיאן של הפונקציה  $\varphi$  ביחס לווקטור  $\bar{\alpha}$ : (נסמן  $\bar{u} \triangleq A\bar{x}_k + R\bar{\alpha}$ )

$$\begin{aligned}\nabla_{\bar{\alpha}} \varphi(\bar{u}) &= R^T \nabla_{\bar{u}} \varphi(\bar{u}) \\ \nabla_{\bar{\alpha}\bar{\alpha}}^2 \varphi(\bar{u}) &= R^T \nabla_{\bar{u}\bar{u}}^2 \varphi(\bar{u}) R\end{aligned}$$

## שיטות מסוג Quasi-Newton's Method:

שיטות מסוג זה שימושיות כאשר בעיית האופטימיזציה שלנו היא ממימד לא גבוה יותר מידי, למשל עד כמה אלפי משתנים. בשיטת ניוטון המקורית, ראינו כי מספר הפעולות הנדרשות בכפל מטריצות (או פתירת מערכת משוואות) הוא סדר גודל של  $O(n^3)$  ובשיטות מסוג זה אנחנו נראה כי מספר הפעולות הוא בסדר גודל  $O(n^2)$  כאשר עיקר החיסכון בחישובים נעוץ

בעובדה שאנחנו לא נחשב את מטריצת ההסיאן בצורה מדויקת אלא רק בצורה מקורבת ולכן שיטה זו שימושית מאד כאשר החישוב של מטריצת ההסיאן הוא קשה עד בלתי אפשרי מבחינה מעשית.

נניח כי  $f(\bar{x})$  היא "חלקה" (לפחות גזירה ברציפות) ובעית האופטימיזציה היא כרגיל:

$$\bar{x} = \arg \min_{\bar{x} \in \mathbb{R}^n} f(\bar{x}), \quad \bar{x} \in \mathbb{R}^n$$

כאשר הגרדיאנט של הפונקציה  $f(\bar{x})$  בזמן האיטרציה ה- $k$  הוא  $\bar{g}_k(\bar{x}_k) = \nabla f(\bar{x}_k)$ .

נסמן את המטריצה ההופכית של מטריצת ההסיאן בצעד ה- $k$  של האלגוריתם להיות:

$$B_k = (H(\bar{x}_k))^{-1} = (\nabla^2 f(\bar{x}_k))^{-1}$$

ונדאג לכך שהמטריצה  $B_k$  תהיה מטריצה סימטרית וכן חיובית ( $B_k \succ 0$ ) כך שכיוון ההתקדמות יהיה תמיד בכיוון הירידה של הפונקציה  $f(\bar{x})$  (כיוון שאם  $B_k \succ 0$  אז באזור הנקודה שאנחנו נמצאים הפונקציה  $f(\bar{x})$  תהיה קמורה כי אם  $B_k \succ 0$  אז גם  $(H(\bar{x})) \succ 0$ ).

האלגוריתם של שיטות מסוג Quasi-Newton's Method בכל איטרציה  $k$ :

(1) אם תנאי העצירה מתקיים, עצור. אחרת, חשב את הכיוון הניוטוני המקורב:

$$\bar{d}_k = -B_k \cdot \bar{g}(\bar{x}_k)$$

(2) מצא את גודל הצעד בכיוון ההתקדמות,  $\alpha_k$ , ע"י אלגוריתם Line Search (ולאו דווקא Exact Line Search אלא

אפשר למשל להשתמש בחוק Armijo וזה יתרון על פני שימוש בשיטת Conjugate Gradient ששם אנחנו מחוייבים להשתמש ב-Exact Line Search אשר מצריך יותר חישובים).

(3) בצע את הצעד הבא:

$$\bar{x}_{k+1} = \bar{x}_k + \alpha_k \bar{d}_k$$

(4) חשב את הגרדיאנט בנקודה החדשה  $\bar{x}_{k+1}$ , כלומר את  $\bar{g}(\bar{x}_{k+1})$ .

(5) חשב את המטריצה  $B_{k+1}$  באמצעות הפרמטרים הבאים:  $\{B_k, \bar{x}_{k+1} - \bar{x}_k, \bar{g}(\bar{x}_{k+1}) - \bar{g}(\bar{x}_k)\}$ . (מיד נראה כיצד לחשב זאת בצורה יעילה) וחזור לשלב 1 באלגוריתם.

כיצד לחשב את המטריצה ההופכית של ההסיאן המקורב בכל איטרציה  $k$  באלגוריתם – המטריצה  $B_k$ :

נתחיל דווקא בהסבר כיצד לחשב את ההסיאן המקורב (ולא את המטריצה ההופכית של ההסיאן) בכל איטרציה.

נסמן תחילה:

$$\bar{p}_k \triangleq \bar{x}_{k+1} - \bar{x}_k$$

$$\bar{q}_k \triangleq \bar{g}(\bar{x}_{k+1}) - \bar{g}(\bar{x}_k)$$

ניזכר כי לפי ההגדרה הדיפרנציאלית של מטריצת ההסיאן מתקיים:

$$H(\bar{x}) d\bar{x} = d\bar{g}(\bar{x})$$

והביטוי הזה מתאר את השינוי בגרדיאנט אם נזוז בכיוון אינפניטיסימלי  $d\bar{x}$ . בפונקציה  $f(\bar{x})$  כללית כל שינוי קטן ב-  $d\bar{x}$  ייתן לנו (לרוב) שינוי קטן ב-  $d\bar{g}(\bar{x})$  ושינוי גדול ב-  $d\bar{x}$  יכול לתת לנו שינוי גדול ב-  $d\bar{g}(\bar{x})$ . נבחין כי בפונקציה  $f(\bar{x})$  ריבועית, המצב אחר, גם עבור שינוי קטן וגם עבור שינוי גדול (לא יותר מידוי) ב-  $d\bar{x}$  נקבל שינוי יחסית קטן ב-  $d\bar{g}(\bar{x})$  כי הגרדיאנט בפונקציה ריבועית הוא לינארי). תחילה נתאר את החישוב עבור פונקציה  $f(\bar{x})$  ריבועית.

חישוב מטריצת ההסיאן המקורבת לפונקציה ריבועית:

מתקיים (לפי ההגדרות עבור  $\bar{p}_k, \bar{q}_k$  שראינו זה עתה):

$$H(\bar{x}_{k+1})\bar{p}_k = \bar{q}_k$$

(המשוואה האחרונה נקראת משוואת המיתר (secant equation))

נניח כי בידינו החישוב המקורב הקודם של מטריצת ההסיאן,  $H(\bar{x}_k)$ , ואנחנו רוצים להשתמש במידע זה על מנת לחשב את המטריצה המקורבת בצעד הבא. חישוב/עדכון המטריצה בצעד הבא יכול להיות מחושב בשיטת Rank one update, אשר שיטה זו אומרת כי אנחנו מעדכנים את המטריצה ע"י הוספת התוצאה של מכפלה פנימית בין ווקטור עמודה כלשהו לבין ווקטור שורה כלשהו (שתוצאה זו היא מטריצה) למטריצה בצעד הקודם. כלומר מתקיים:

$$\text{Rank One update: } H(\bar{x}_{k+1}) = H(\bar{x}_k) + \bar{u} \cdot \bar{v}^T \quad (\bar{u} - \text{Column Vector})$$

וממשוואת המיתר נקבל:

$$H(\bar{x}_{k+1})\bar{p}_k = (H(\bar{x}_k) + \bar{u} \cdot \bar{v}^T)\bar{p}_k = \bar{q}_k$$

נבחין כי מטריצת ההסיאן היא מטריצה סימטרית  $f(\bar{x})$  גזירה ברציפות ומשפט שוורץ מתקיים ואפשר להחליף סדר גזירה) אז אנחנו רוצים שגם המטריצה המקורבת תהיה סימטרית ולכן אנחנו צריכים לדאוג שמטריצה  $\bar{u} \cdot \bar{v}^T$  תהיה סימטרית (כי אנחנו מוסיפים אותה למטריצה סימטרית ולכן בשביל להישאר עם מטריצה סימטרית אנחנו חייבים שהמטריצה  $\bar{u} \cdot \bar{v}^T$  תהיה גם כן סימטרית), אך זה אפשרי אם ורק אם  $\bar{u} = \alpha \cdot \bar{v}$  כאשר  $\alpha \in \mathbb{R}$ , כלומר  $\bar{u} \propto \bar{v}$ .  
 וע"י פישוט קל של המשוואה האחרונה נקבל:

$$\bar{u}(\bar{v}^T \bar{p}_k) = \bar{q}_k - H(\bar{x}_k)\bar{p}_k$$

ולכן עבור כל ווקטור  $\bar{v}^T$  שאינו אורתוגנלי לווקטור  $\bar{p}_k$ , כלומר  $\bar{v}^T \bar{p}_k \neq 0$  (ולכן  $\bar{v}^T \bar{p}_k = \langle \bar{v}, \bar{p}_k \rangle \neq 0$ ), נקבל:

$$\bar{u} = \frac{1}{\bar{v}^T \bar{p}_k} (\bar{q}_k - H(\bar{x}_k)\bar{p}_k)$$

נבחין כי  $\frac{1}{\bar{v}^T \bar{p}_k}$  זה סקלר ולכן נוכל לבחור למשל:

$$\bar{v} = \bar{q}_k - H(\bar{x}_k)\bar{p}_k$$

כלומר קיבלנו:

$$\bar{u} = \frac{1}{\bar{v}^T \bar{p}_k} \bar{v}$$

ובסך הכול נקבל כי מטריצת ההסיאן בצעד הבא של האלגוריתם:

$$H(\bar{x}_{k+1}) = H(\bar{x}_k) + \frac{1}{\bar{v}^T \bar{p}_k} \bar{v} \cdot \bar{v}^T$$

אך למעשה אנחנו מעוניינים במטריצה ההופכית של מטריצת ההסיאן. ישנן נוסחאות שמוצאות את המטריצה ההופכית של המטריצה החדשה בהינתן זו שהמטריצה החדשה חושבה ע"י שיטת Rank one update, למשל נוסחת Sherman-Morrison אשר יודעת לחשב את מטריצת ההסיאן ההופכית,  $B_{k+1} = (H(\bar{x}_{k+1}))^{-1}$ , בהינתן שהמטריצה  $H(\bar{x}_{k+1})$  חושבה בשיטת Rank one update.

כעת, בהינתן זה שאנחנו יודעים את המטריצה ההופכית, ממשוואת המיתר שראינו:  $H(\bar{x}_{k+1}) \bar{p}_k = \bar{q}_k$ , נכפול מצד שמאל במטריצה ההופכית,  $B_{k+1} = (H(\bar{x}_{k+1}))^{-1}$ , ונקבל:

$$B_{k+1} \bar{q}_k = \bar{p}_k$$

ונבחין כי המשוואה האחרונה דומה מאד למשוואת המיתר:  $H(\bar{x}_{k+1}) \bar{p}_k = \bar{q}_k$ , ולכן נוכל לרשום ע"י החלפת המקומות של הווקטורים  $\bar{p}_k, \bar{q}_k$  ונקבל בדומה לביטוי שקיבלנו עבור  $H(\bar{x}_{k+1})$  את הביטוי עבור  $B_{k+1} = (H(\bar{x}_{k+1}))^{-1}$ :

$$\begin{aligned} \bar{v} &= \bar{p}_k - B_k \bar{q}_k \\ B_{k+1} &= B_k + \frac{1}{\bar{v}^T \bar{q}_k} \bar{v} \cdot \bar{v}^T \end{aligned}$$

(לא הגדרנו מהו הווקטור  $\bar{v}$  ורק אמרנו שהוא ווקטור כלשהו!)

ובזאת מצאנו שיטה לחשב את המטריצה  $B_{k+1} = (H(\bar{x}_{k+1}))^{-1}$  אך זו לא השיטה היחידה שקיימת ויש שיטות אחרות ויעילות יותר ויש עוד פרטים שיש לדאוג להם על מנת שהמטריצה  $B_{k+1}$  תהיה חיובית והאלגוריתם הבא שנראה למציאת  $B_{k+1}$  נחשב לאלגוריתם הטוב ביותר שהומצא עד היום:

שיטות מסוג Broyden family of Quasi-Newton's Method:

נזכיר תחילה את הסימון:

$$\begin{aligned} \bar{p}_k &\triangleq \bar{x}_{k+1} - \bar{x}_k \\ \bar{q}_k &\triangleq \bar{g}(\bar{x}_{k+1}) - \bar{g}(\bar{x}_k) \end{aligned}$$

כעת נסמן בנוסף:

$$\begin{aligned}\bar{s}_k &\triangleq B_k \bar{q}_k \\ \tau_k &\triangleq \bar{s}_k^T \bar{q}_k, \quad \tau_k \in \mathbb{R} \\ \mu_k &\triangleq \bar{p}_k^T \bar{q}_k, \quad \mu_k \in \mathbb{R} \\ \bar{v}_k &= \frac{\bar{p}_k}{\mu_k} - \frac{\bar{s}_k}{\tau_k}\end{aligned}$$

כאשר נבחין כי  $\tau_k, \mu_k$  הם סקלרים. לכן בשיטות מסוג Broyden family of Quasi-Newton's Method מתקיים:

$$B_{k+1} = B_k + \frac{\bar{p}_k \cdot \bar{p}_k^T}{\mu_k} - \frac{\bar{s}_k \cdot \bar{s}_k^T}{\tau_k} + \xi_k \cdot \tau_k \cdot \bar{v}_k \cdot \bar{v}_k^T$$

כאשר החישוב האחרון הוא עדכון שמתבסס על Rank two update. כמו כן,  $\xi_k \in \mathbb{R}$  הוא פרמטר.

שיטת BFGS ושיטת DFP:

אחת השיטות החשובות מסוג Broyden family of Quasi-Newton's Method נקראת BFGS על שם: Broyden, Shanno ו-Fletcher, Goldfrab ושיטה מוכרת נוספת מסוג זה נקראת DFP על שם Davidon, Fletcher ו-Powell. כמו כן, בכל שיטה הפרמטר  $\xi_k \in \mathbb{R}$  הוא שונה למשל:

$$\begin{aligned}BGFS &\Rightarrow \xi_k = 1 \\ DFP &\Rightarrow \xi_k = 0\end{aligned}$$

נבחין כי בנוסחה האחרונה עבור  $B_{k+1}$  מספר פעולות החישוב הנדרש הוא  $O(n^2)$  כאשר  $n \times n$  הוא המימד של מטריצה  $B_{k+1}$ .

הערות:

- השיטה BFGS היא מאד שימושית ואפשר למצוא אותה גם במטלב (Matlab) ב-Matlab Optimization Toolbox וזה נתון לבחירתנו באיזו מן השיטות לעיל להשתמש לחישובים. נזכיר כי שיטות אלו שימושיות כאשר חישוב ערך הפונקציה והגרדיאנט בנקודה הן יקרות יחסית, כי אם הפעולות האלה זולות אולי כדאי דווקא לחזור לשיטות שלמדנו לפני שהן Conjugate Gradient או Subspace Optimization.

- הנוסחה  $B_{k+1} = B_k + \frac{\bar{p}_k \cdot \bar{p}_k^T}{\mu_k} - \frac{\bar{s}_k \cdot \bar{s}_k^T}{\tau_k} + \xi_k \cdot \tau_k \cdot \bar{v}_k \cdot \bar{v}_k^T$  מקיימת את משוואת המיתר  $H(\bar{x}_{k+1}) \bar{p}_k = \bar{q}_k$ .

והמטריצה  $B_{k+1}$  היא סימטרית ומקיימת  $B_{k+1} \succ 0$  וכמו כן מתקיים כי המטריצה  $B_{k+1}$  היא בעלת הנורמה המינימלית במובן של:  $\min \|B_{k+1} - B_k\|$ .

- במידה והפונקציה  $f(\bar{x})$  היא לא קמורה, אז בשביל שנקבל  $B_{k+1} \succ 0$  ע"י הנוסחה

$$B_{k+1} = B_k + \frac{\bar{p}_k \cdot \bar{p}_k^T}{\mu_k} - \frac{\bar{s}_k \cdot \bar{s}_k^T}{\tau_k} + \xi_k \cdot \tau_k \cdot \bar{v}_k \cdot \bar{v}_k^T$$

יש לדרוש כי ה-Line Search יקיים:

$$f'_{\bar{d}_k}(\bar{x}_k) < f'_{\bar{d}_k}(\bar{x}_{k+1})$$

שלב האתחול בשיטות מסוג Broyden family of Quasi-Newton's Method:

אם איננו יודעים את מטריצת ההסיאן בשלב ההתחלה אז נוכל להתחיל את השיטה ע"י:  $B_0 = I_{[n \times n]}$

ואם אנחנו יודעים את מטריצת ההסיאן באופן כזה או אחר אז נקרב את המטריצה  $B_0$  למטריצה ההופכית של הסיאן כמה שנוכל:

$$B_0 \approx H(\bar{x}_0)$$

הערה:

- אם נשתמש בשיטת BFGS Quasi-Newton's Method על פונקציה ריבועית, ונשתמש ב-Exact Line Search (שעבור פונקציה ריבועית הוא למעשה חישוב אנליטי) אז באופן מפתיע אנחנו נקבל את אות המסלול באופטימיזציה שהיינו מקבלים לו היינו משתמשים בשיטת Conjugate Gradient ולכן נפתור את בעיית האופטימיזציה (שהיא במימד  $n$ ) בדיוק ב- $n$  צעדים ולאחר  $n$  הצעדים הללו תהיה בידינו מטריצת ההסיאן המדויקת של הפונקציה המקורית (אפילו אם בשלב האתחול של השיטה נתחיל עם מטריצה מקורבת להסיאן  $B_0 = I_{[n \times n]}$ ).
- אם נשתמש בשיטת BFGS Quasi-Newton's Method על פונקציה שאינה ריבועית, בדר"כ נקבל פתרון לאחר פחות איטרציות מאשר אם היינו משתמשים בשיטת Conjugate Gradient וכן שיטת BFGS אינה דורשת את שיטת Exact Line Search ולכן היא חסכונית במידה משמעותית משיטת Conjugate Gradient מבחינת כמות החישובים.
- קצב ההתכנסות של השיטה: ראינו שבשיטת Gradient Descent עבור פונקציה שאינה לינארית, קצב ההתכנסות הוא לינארי ובשיטה ניוטון על פונקציה לא לינארית קצב ההתכנסות הוא ריבועי. בשיטת Quasi-Newton קצב ההתכנסות הוא סופר-לינארי. נזכיר כי אם נסמן את נקודת המינימום להיות  $\bar{x}^*$  אז קצב התכנסות לינארי אומר כי מתקיים:

$$Linear\ Convergence\ Rate: \forall k: \frac{\|\bar{x}_{k+1} - \bar{x}^*\|}{\|\bar{x}_k - \bar{x}^*\|} = Const \in \mathbb{R}$$

וקצב התכנסות סופר לינארי מקיים:

$$Super - Linear\ Convergence\ Rate: \lim_{k \rightarrow \infty} \frac{\|\bar{x}_{k+1} - \bar{x}^*\|}{\|\bar{x}_k - \bar{x}^*\|} \rightarrow 0$$

## **הרצאה 13 – Summary Of Unconstrained Optimization And Intro To Optimization With Constraints**

סיכום השיטות לאופטימיזציה ללא אילוצים:

1. שיטות אופטימיזציה של בעיות במימד אחד:
  - a. שיטת Golden Section – שיטה זו שימושית בבעיות אופטימיזציה מממד אחד שבהן אנחנו יכולים רק לחשב את ערכי הפונקציה אך לא יכולים לחשב את ערכי הנגזרת (מסיבות כאלה ואחרות).
  - b. שיטת Bisection – שיטה זו הייתה שימושית כאשר היינו יכולים לחשב את הנגזרת של הפונקציה בכל נקודה, ותנאי הכרחי שנקודה תהווה מינימום הוא התאפסות הנגזרת.
  - c. שיטת Quadratic/Cubic Interpolation – בשיטה זו אנחנו מניחים כי בסביבת נקודת המינימום של הפונקציה, הפונקציה דומה לפולינום מדרגה שנייה או שלישית (וזה אכן נכון לרוב הפונקציות החלקות בסביבת נקודת המינימום). שימושים לשיטה זו:
    - i. ניתן גם להשתמש בשיטה זו יחד עם שיטת Golden Section ובכל פעם שיש לנו 3 נקודות (בשיטת Golden Section) לחשב את המינימום באמצעים אנליטיים כי באמצעות 3 נקודות אפשר לקרב פונקציה ריבועית.
    - ii. ניתן גם להשתמש בשיטת זו יחד עם שיטת Bisection ולקרב את הפונקציה שלנו בעזרת ערכי הנגזרת ולא רק ע"י ערכי הפונקציה.
  - d. שיטת Inexact Line Search – למדנו ליישם שיטה זו באמצעות Backtracking Method יחד עם חוק Armijo.
2. שיטות אופטימיזציה של בעיות במספר מימדים:
  - a. שיטת Steepest Descent / Gradient Descent – בכל צעד באלגוריתם הולכים בכיוון שמנוגד לכיוון הגרדיאנט (כי הגרדיאנט מצביע על כיוון העליה החדה ביותר) ולאחר כל בחירת כיוון מחפשים את גודל הצעד המתאים, כאשר בעיה זו היא למעשה בעיה במימד אחד ולכן משתמשים באחת משיטות האופטימיזציה במימד אחד. אחד החסרונות הבולטים של שיטה זו, זה כאשר נקודת המינימום נמצא במעיין "ואדי צר" (ill condition function) ואז כיווני הגרדיאנט מתנהגים כמו "זיג-זג" והתכנסות האלגוריתם היא מאד איטית.
  - b. שיטת Newton – שיטה טובה עבור ill condition function. עבור בעיות מסוג Least Square Problems השתמשנו ב"תת-שיטה" של שיטת Newton וקראנו לשיטה זו Newton-Gauss Method. שיטת Newton דורשת בכל צעד של האלגוריתם לפתור מערכת גדולה של משוואות לינאריות (בגודל  $n \times n$  עבור בעיות אופטימיזציה מממד  $n$ ). אם  $n$  גדול מידי (סדר גודל של אלפים/עשרות אלפים) שיטה זו לא ישימה במחשבים של ימינו.
  - c. שיטת Conjugate Gradient – זו שיטה שראינו למזעור של פונקציות ריבועיות וכן פונקציות לינאריות אחרות (לא ריבועיות אך "חלקות").
    - i. שיטת Truncated Newton – בשיטה זו ראינו שאפשר לפתור בעית אופטימיזציה בשיטת Newton אך את מערכת המשוואות בנוצרת בשיטת Newton נפתור לא באופן מלא אלא באופן מקורב ע"י שיטת Conjugate Gradient.
  - d. שיטות Quasi-Newton – בכל איטרציה בשיטה זו, אנחנו בונים מטריצה אשר מקרבת את מטריצת ההסיאן האמיתית או למעשה מקרבת את הפונקציה ההופכית של ההסיאן. שיטה מסוג Quasi-Newton שלמדנו היא שיטת BFGS. שיטה מסוג Quasi-Newton היא ככל הנראה השיטה היעילה ביותר לכל בעית אופטימיזציה עד מימד של כמה אלפים בודדים של משתנים.
  - e. שיטת Sequential Subspace Optimization (SESOP) – שיטה זו היא למעשה עוד הכללה של שיטת Conjugate Gradient למזעור של פונקציות לא לינאריות ולא ריבועיות וכל איטרציה של שיטה זו, אנחנו

למעשה משתמשים בגרדיאנט של הנקודה בה אנחנו עומדים ובכל הכיוונים שהשתמנו בהם עד לאיטרציה הנוכחית על מנת לחשב את הכיוון וגודל הצעד באיטרציה הנוכחית. נבחין כי בכדי לחשב כל צעד אנחנו פותרים בעיית אופטימיזציה אבל במספר יחסית קטן של משתנים, כי המשתנים הם כל האיטרציות הקודמות ולא כמספר המשתנים המקורי של הבעיה. במקרה של פונקציות ריבועיות השיטה הזו תתהג באופן זהה כמו שיטת Conjugate Gradient. שיטה זו יעילה מאד כאשר הפונקציה שאותה אנחנו ממזערים היא מהצורה:  

$$f(\bar{x}) = \varphi(A\bar{x})$$
, כלומר הפונקציה היא פונקציה לא לינארית ( $\varphi$ ) של אופרטור לינארי כלשהו ( $A$ ) וכן הקושי החישוב נעוץ דווקא בחישוב של האופרטור הלינארי ( $A\bar{x}$ ) ולא בחישוב הפונקציה  $\varphi$  ואז כאשר אנחנו מחפשים את גודל וכיוון הצעד הבא, אנחנו מחפשים אותו ב-Subspace ולא צריכים לבצע את החישוב  $A\bar{x}$  כל פעם מחדש (שהוא, כאמור, החישוב היקר מבחינת זמן). שיטה זו יכולה להיות יעילה מאד בבעיות אופטימיזציה אפילו עם מאות אלפי משתנים.  
 f. שיטת Nelder-Mead Simplex Method – שיטה זו לא נלמדה בקורס. שיטה זו טובה לבעיות במספר משתנים שלא עולה על עשרות בודדות. שיטה זו לא דורשת לדעת את הנגזרות של הפונקציה אותה ממזערים אלא רק את ערכיה בנקודות שנרצה.

#### אופטימיזציה עם אילוצים:

אופטימיזציה עם אילוצים זה החלק המרכזי והעיקרי של החלק השני של הקורס. תחילה נלמד מה הם התנאים ההכרחיים עבור נקודה להיות הנקודה בה מתקבל המינימום של הפונקציה אותה אנחנו ממזערים ולאחר מכן נלמד שיטות ואלגוריתמים יעילים למציאת הנקודות הללו.

בעיות אופטימיזציה עם אילוצים הן מהצורה:

$$\begin{aligned} \min_{\bar{x} \in \mathbb{R}^n} f(\bar{x}) \\ \text{Subject To } (S.T): \quad & g_i(\bar{x}) \leq 0, \quad i = 1, \dots, m \\ & h_j(\bar{x}) = 0, \quad j = 1, \dots, k \end{aligned}$$

כאשר אנחנו נוהגים לקרוא לפונקציה שאותה אנחנו ממזערים  $f$  להיות פונקציה המטרה (Objective function או Cost function). הפונקציות  $g_i: \mathbb{R}^n \rightarrow \mathbb{R}$  נקראות Inequality Constraints והפונקציות  $h_j: \mathbb{R}^n \rightarrow \mathbb{R}$  נקראות Equality Constraints. בבעיות מסוג זה אנחנו מחפשים את הנקודה  $\bar{x} \in \mathbb{R}^n$  אשר מביאה למינימום את הפונקציה  $f$  אך גם מקיימת את כל אי-השוויונות בהצבה בפונקציות  $g_i$  ומקיימת את כל השוויונות בהצבה בפונקציות  $h_j$ .

אנחנו נתחיל בשיטות לפתירת בעיות אופטימיזציה עם פונקציות מסוג Inequality Constraints בלבד ולאחר מכן נכליל לבעיות אופטימיזציה גם עם פונקציות Equality Constraints.

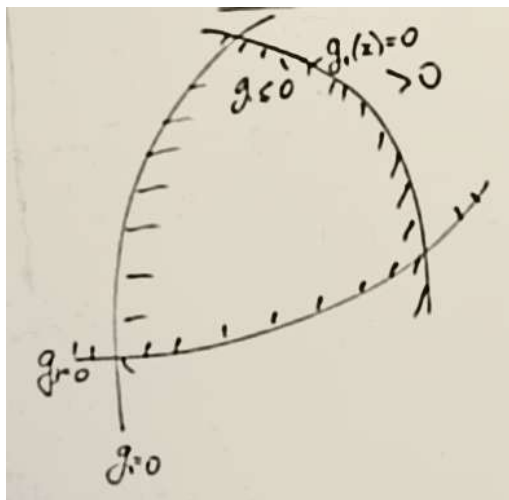
#### בעיות אופטימיזציה עם פונקציות מסוג Inequality Constraints (בלבד):

בעיות אלו באופן כללי מיוצגות באופן הבא:

$$\begin{aligned} \min_{\bar{x} \in \mathbb{R}^n} f(\bar{x}) \\ \text{Subject To } (S.T): \quad & g_i(\bar{x}) \leq 0, \quad i = 1, \dots, m \end{aligned}$$

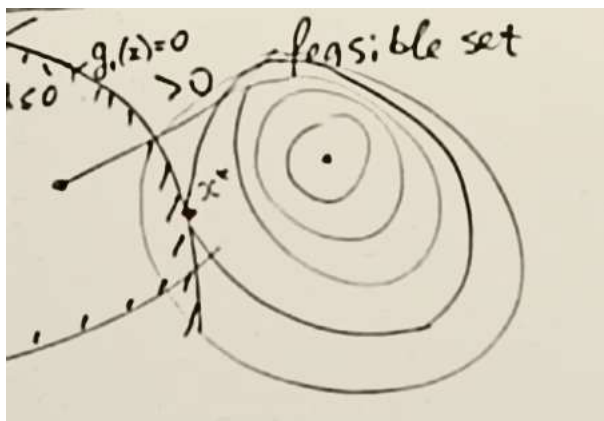
ננסה תחילה להבין את הגיאומטריה של הבעיה. אנחנו נצייר את קווי הגובה של הפונקציה  $f(\bar{x})$  ואת קווי הגובה של הפונקציות  $g_i(\bar{x})$ . נניח ומתקיים  $m = 3$  ומתקיים "קווי הגובה" של הפונקציות  $g_i(\bar{x})$  הן:





כלומר הפונקציות  $g_i(\bar{x})$  בעלות תחום משותף שבו שלושתן מקבלות ערכים קטנים מאפס בכל נקודה בתחום זה. לתחום זה קוראים Feasible Set. כל נקודה  $\bar{x} \in \text{Feasible Set}$  מקיימת את כל אילוצי הבעיה.

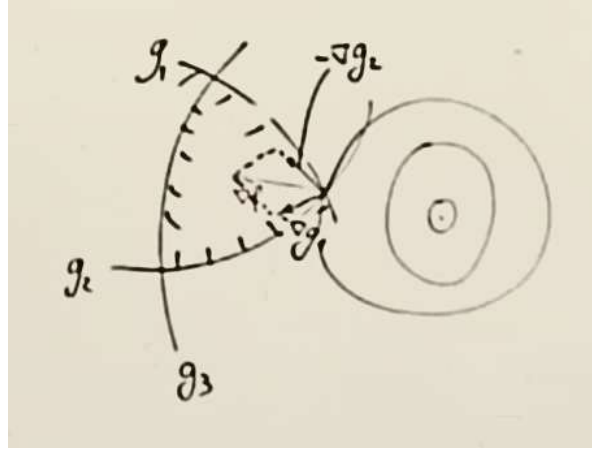
כעת נצייר את קווי הגובה של הפונקציה,  $f(\bar{x})$ , אותה אנחנו רוצים למזער:



במרכז יש את נקודת המינימום (מודגשת), אך מהצירור ברור כי רק הנקודה המודגשת משמאל  $(\bar{x}^*)$  מקיימת את האילוצים של הפונקציות  $g_i(\bar{x})$  והיא הנקודה בה מתקבל המינימום החוקי של הפונקציה  $f(\bar{x})$  בבעיה הנתונה. אנחנו יודעים כי כיוון הגרדיאנט (למשל של  $f(\bar{x})$ ) הוא תמיד מאונך לקווי הגובה של הפונקציה ולכן גם בנקודה  $\bar{x}^*$  יש גרדיאנט שמאונך לקו הגובה המתאים. בנוסף, נביט על הגרדיאנט של הפונקציה  $g_1(\bar{x})$  בנקודה  $\bar{x}^*$ , ברור כי מתקיים  $g_1(\bar{x}^*) = 0$  (ניתן לראות בצירור כי מצד אחד של הנקודה הפונקציה היא חיובית ומצד שני היא שלילית, וכיוון שהיא רציפה אז לפני משפט ערך הביניים על הקו הנ"ל הפונקציה מתאפסת) אך כיוון הגרדיאנט הוא בדיוק כיוון ההפוך של כיוון הגרדיאנט של הפונקציה  $f(\bar{x})$ , ובאופן פורמלי:

$$\nabla f(\bar{x}^*) = -\lambda \nabla g_1(\bar{x}^*) \quad \Rightarrow \quad \nabla f(\bar{x}^*) + \lambda \nabla g_1(\bar{x}^*) = 0, \quad \lambda \in \mathbb{R}$$

במקרה הנ"ל ראינו כי רק אילוץ אחד "נוגע" בנקודה  $\bar{x}^*$ , אך תיתכנה בעיות בהן כמה אילוצים יגעו בנקודה האופטימלית שאנחנו מחפשים, למשל שני אילוצים:



כלומר במקרה זה מתקיים (ובאופן כללי):

$$\nabla f(\bar{x}^*) + \sum_{i \in I_a} \lambda_i \nabla g_i(\bar{x}^*) = 0, \quad I_a \equiv \{\text{Active Constraints}\}, \quad \lambda_i \in \mathbb{R}$$

כאשר  $I_a$  מייצגת את קבוצת האילוצים הפעילים (Active Constraints) אשר באמת נוגעים בנקודה (בניגוד לפונקצית האילוץ  $g_3(\bar{x})$  בדוגמה לעיל).

נגדיר את פונקציית הלגרנז'אן:

$$L(\bar{x}, \bar{\lambda}) \triangleq f(\bar{x}) + \sum_{i \in I} \lambda_i g_i(\bar{x}^*)$$

ולכן נוכל למעשה לבטא את המשוואה  $\nabla f(\bar{x}^*) + \sum_{i \in I_a} \lambda_i \nabla g_i(\bar{x}^*) = 0$  (שראינו זה עתה) באמצעות הלגרנז'אן:

$$\nabla_{\bar{x}} L(\bar{x}^*, \bar{\lambda}^*) = 0 \quad \Leftrightarrow \quad \nabla f(\bar{x}^*) + \sum_{i \in I_a} \lambda_i \nabla g_i(\bar{x}^*) = 0$$

התנאים ההכרחיים מסדר ראשון לאופטימליות של פתרון של Karush-Kuhn-Tucker (תנאי KKT):

נניח כי הנקודה  $\bar{x}^*$  היא הנקודה הפותרת את בעיית האופטימיזציה הבאה:

$$\min_{\bar{x} \in \mathbb{R}^n} f(\bar{x})$$

$$\text{Subject To (S.T):} \quad g_i(\bar{x}) \leq 0, \quad i = 1, \dots, m$$

ונסמן את הגרדיאנטים  $\nabla g_i(\bar{x}^*)$ , כאשר  $i \in I_a$  (כאשר  $I_a \equiv \{\text{Active Constraints}\}$ ), וכן נניח כי הם בלתי תלויים לינארית אחד בשני. במצב זה, קיים  $\bar{\lambda}^*$  (ווקטור כופלי לגרנז') אשר עבורו נקבל:

$$\nabla_{\bar{x}} L(\bar{x}^*, \bar{\lambda}^*) = 0, \quad \bar{g}(\bar{x}^*) \leq 0$$

(כאשר  $\bar{g}(\bar{x})$  זה סימון ווקטורי לכל פונקציות האילוצים)

ועבור אילוצים מסוג Inequality Constraints מתקיים כי  $\lambda_i^* \geq 0$ . בנוסף, עבור אילוצים שהם לא בקבוצה  $I_a$  (אילוצים לא פעילים) מתקיים  $\lambda_i^* = 0$ .

אפשר לסכם את כל התנאים הללו:

$$\begin{aligned} \nabla_{\bar{x}} L(\bar{x}^*, \bar{\lambda}^*) &= 0, & \bar{g}(\bar{x}^*) &\leq 0 \\ \lambda_i^* &\geq 0, & i &\in I_a \\ \lambda_i^* &= 0, & i &\notin I_a \end{aligned}$$

לתנאי אחד בלבד (תנאי Complimentary Slackness):

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

(כי עבור כל אילוץ לא פעיל, מתקיים  $\lambda_i^* = 0$  ועבור כל אילוץ פעיל מתקיים  $g_i(\bar{x}^*) = 0$ )

שיטת פונקצית עונשין (Penalty function method) עבור אילוצים מסוג Inequality Constraints:

את בעיית האופטימיזציה:

$$\begin{aligned} \min_{\bar{x} \in \mathbb{R}^n} f(\bar{x}) \\ \text{Subject To (S.T): } g_i(\bar{x}) &\leq 0, \quad i = 1, \dots, m \end{aligned}$$

ניתן גם לרשום באופן הבא:

$$\begin{aligned} \min_{\bar{x} \in G} f(\bar{x}) \\ G = \{ \bar{x} \in \mathbb{R}^n \mid g_i(\bar{x}) \leq 0, i = 1, \dots, m \} \end{aligned}$$

אנחנו נגדיר פונקצית עונשין ( $\varphi: \mathbb{R} \rightarrow \mathbb{R}$ ) לכל אחת מפונקציות האילוצים  $g_i(\bar{x})$ , כך שעבור כל  $g_i(\bar{x})$  (אשר לא מקיים את האילוץ) תקבל ערך ("עונשין") שונה, וככל שערך זה גדול יותר כך זה אומר כי הנקודה הנוכחית שאנחנו נמצאים בה פחות קרובה לפתרון האופטימלי. כלומר יהיו אילוצים שנרצה לעמוד בהם באופן נוקשה ויהיו אילוצים שנאפשר הפרה שלהם (אם העונש לא גדול מדי). נגדיר פרמטר  $p \in \mathbb{R}$  באופן הבא:

עבור  $p \rightarrow 0$  נקבל כי פונקצית העונשין  $\varphi_p$ , נותנת עונש "קל" עבור חריגה של  $g_i(\bar{x}) = t \geq 0$  ועבור  $p \rightarrow \infty$  העונש עבור חריגה הוא אינסופי. למשל אפשר לחשוב על פונקצית העונשין הבאה:

$$\varphi_p(t) = e^{pt} - 1, \quad t, p \in \mathbb{R}$$

בנוסף, עבור  $t \leq 0$  נרצה שהעונש (ערך הפונקציה  $\varphi(t)$ ) יהיה מינימלי (ולכן הפונקציה שהוגדרה לעיל לא מתאימה לחלוטין לדרישות שלנו). הדרישות שלנו מפונקצית העונשין הן:

- (1) פונקציה קמורה.
- (2) פונקציה מונוטונית עולה.
- (3) פונקציה "חלקה".

$$(4) \quad \lim_{t \rightarrow \infty} \varphi'(t) = \infty \quad \text{וכן} \quad \lim_{t \rightarrow \infty} \varphi(t) = 0$$

$$(5) \quad \text{תקיים: } \varphi(0) = 0 \quad \text{וגם} \quad \varphi'(0) = 1 \quad (\text{דרישה זו היא לשם נוחות})$$

נגדיר את פרמטר העונשין (Penalty Parameter):  $p \in \mathbb{R}^+ \setminus \{0\}$

נגדיר גם את פונקציית העונשין  $\varphi_p(t)$  להיות:

$$\varphi_p(t) = \frac{1}{p} \varphi(pt)$$

נראה שעבור  $t > 0$  ועבור ערך  $p$  גדול אנחנו מקבלים עונש "כבד", נבדוק את הנגזרת הראשונה ונקבל:

$$\frac{\partial}{\partial t}(\varphi_p(t)) = \frac{1}{p} \varphi'(pt) \cdot p = \varphi'(pt)$$

ולכן עבור ערך  $p$  גדול, הארגומנט  $pt$  הוא מאד גדול ולכן לפי הדרישה שלנו מתקיים:  $\lim_{t \rightarrow \infty} \varphi'(t) = \infty$ , כלומר אפילו חריגה קלה מהתחום המותר עבור  $t$  תגרום לעונש "כבד". גם במקרה ההפוך, אם  $p$  הוא מאד גדול, וכן  $t < 0$  אז פונקציית העונשין שואפת לאפס.

דוגמאות לפונקציות עונשין:

$$(1) \quad \varphi(t) = e^t - 1$$

$$(2) \quad \varphi(t) = \begin{cases} \frac{1}{2}t^2 + t & , t \geq -\frac{1}{2} \\ -\frac{1}{4}\log(-2t) - \frac{3}{8} & , t < -\frac{1}{2} \end{cases}$$

עונשין כפי שהגדרנו אותן

פונקציית העונשין האופטימלית/אידיאלית (Ideal Penalty Function):

באופן כללי אפשר לומר כי פונקציית העונשין האופטימלית/אידיאלית (Ideal Penalty Function) מקיימת:

$$\lim_{p \rightarrow \infty} \varphi_p(t) = \varphi_\infty(t) = \begin{cases} 0 & , t \leq 0 \\ \infty & , t > 0 \end{cases} = \text{Ideal Penalty Function}$$

פונקציית העונשין המצרפית (Penalty Aggregate):

נוכל להגדיר את הפונקציה שאותה אנחנו ממזערים בבעיית האופטימיזציה באופן הבא:

$$F_p(\bar{x}) = f(\bar{x}) + \sum_{i=1}^m \varphi_p(g_i(\bar{x}))$$

הרעיון הכללי: בעיה זו ניתן לפתור ע"י בחירת ערך התחלתי  $p$  שהוא יחסית נמוך ואז בעיית אופטימיזציה זו הופכת לבעיית אופטימיזציה ללא אילוצים, ואז מתחילים להעלות את הפרמטר  $p$  ומקבלים בעיית אופטימיזציה עם יותר ויותר אילוצים. אם נצליח למצוא פתרון עבור  $p \rightarrow \infty$  אז נוכל לומר שפתרנו את בעיית האופטימיזציה באופן מוחלט.

בעית האופטימיזציה המצרפית האידיאלית (Ideal Penalty Aggregate):

$$F_{\infty}(\bar{x}) = f(\bar{x}) + \sum_{i=1}^m \varphi_{\infty}(g_i(\bar{x})) = \begin{cases} f(\bar{x}) & , \quad \bar{x} \in G = \{Feasible Set\} \\ \infty & , \quad \bar{x} \notin G = \{Feasible Set\} \end{cases}$$

כלומר שתי בעיות האופטימיזציה הבאות הן שקולות:

$$\min_{\bar{x} \in \mathbb{R}} F_{\infty}(\bar{x}) \quad \Leftrightarrow \quad \min_{\bar{x} \in G} f(\bar{x})$$

אלגוריתם שיטת פונקציית העונשין (Penalty Function algorithm):

(1) התחל עם פרמטר עונשין  $p$  קטן יחסית בנקודה  $\bar{x}_0$ .

(2) באיטרציה ה- $k$  בצע את השלבים הבאים:

a. הנקודה הבאה היא  $\bar{x}_{k+1} \approx \arg \min_{\bar{x}} F_{p_k}(\bar{x})$ , כלומר פותרים בעית אופטימיזציה ללא אילוצים עבור

פרמטר  $p_k$  מסויים.

b. הגדל את פרמטר העונשין  $p$ , באופן הבא:  $p_{k+1} = \alpha \cdot p_k$  (כאשר  $2 \leq \alpha \leq 10$ ) עד אשר

$10^5 \leq p \leq 10^6$  (לא בהכרח המספרים הללו).

הערה:

בשיטת פונקציית העונשין אנחנו אכן יכולים לקבל פתרון מדוייק יחסית, אך עבור ערכים גדולים של פרמטר השגיאה, אנחנו נאלצים לפתור בעית אופטימיזציה ללא האילוצים קשה יותר ויותר ולכן אנחנו מאלצים את הדיוק של הפתרון האמיתי של בעית האופטימיזציה המקורית שניסינו לפתור (הבעיה עם האילוצים). בהרצאה הבאה נראה איך להתמודד עם מצב שכזה.

## הרצאה – Penalty Function Method and Augmented Lagrangian Method For Constrained Optimization

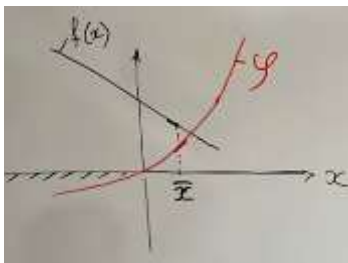
נתבונן בבעיית אופטימיזציה פשוטה (ממימד אחד):

$$\min_{x \in \mathbb{R}} f(x)$$

$$\text{Subject To (S.T): } x \leq 0$$

$$\text{Penalty Function: } \varphi(x) = e^x - 1$$

ולשם המחשה:



ונפתור את הבעיה באמצעות פונקציות עונשין ואנחנו רואים כי פונקציות העונשין מקבלת ערכים חיוביים עבור  $x > 0$ .

ובעית האופטימיזציה היא למעשה:

$$\bar{x} \triangleq \arg \min_x \{F(x) = f(x) + \varphi(x)\}$$

טענה (ללא הוכחה):

הנקודה  $\bar{x}$  (הנקודה בה היא הנקודה בה הפונקציה  $F(x) = f(x) + \varphi(x)$  מקבלת את המינימום שלה) מתקיים

$$f'(\bar{x}) = -\varphi'(\bar{x})$$

הערה: הטענה האחרונה היא למעשה טענה שקולה שבנקודת המינימום הלגרנז'יאן מתאפס.

נבחין כי הפתרון במצב זה לא מדויק, שכן, אנחנו יודעים כי הפתרון לבעיית האופטימיזציה כפי שמתוארת בשרטוט:

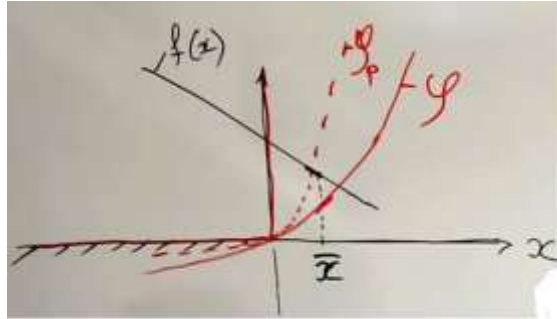
$$\min_{x \in \mathbb{R}} f(x)$$

$$\text{Subject To (S.T): } x \leq 0$$

היא למעשה  $x = 0$ . נגדיר כעת פונקציות העונשין עם פרמטר עונשין:

$$\varphi_p(t) = \frac{1}{p} \varphi(pt)$$

ולכן כאשר  $p \rightarrow \infty$  אנחנו מקבלים כי  $\varphi_p \rightarrow \varphi_\infty$  ואז מקבלים פונקציות עונשין אידיאליות:



$$\varphi_{\infty} = \begin{cases} 0 & , x \in \text{Feasible Set} \\ \infty & , x \notin \text{Feasible Set} \end{cases}$$

ולכן אנחנו למעשה מנסים לפתור את בעיית האופטימיזציה הבאה:

$$\bar{x}_p \triangleq \arg \min_x \{F_p(x) = f(x) + \varphi_p(x)\}$$

כאשר ערך הפרמטר  $p$  הולך וגדל.

נבחין כי כאשר  $p$  הולך וגדל, פתרון בעיית האופטימיזציה ללא האילוצים הולכת ונהיית קשה יותר (מדוע?), כיוון שפונקציית העונשין נהיית יותר ויותר חדה. בכדי להתגבר על הבעיה הזו קיימת השיטה הבאה:

שיטת Augmented Lagrangian Penalty – Multiplier method:

פתרון בעיית האופטימיזציה שנוצרת בעקבות שימוש בשיטת פונקציית העונשין נהיה יותר ויותר קשה כאשר שיפוע פונקציית העונשין הולך וגדל (זוהי קורה כאשר  $p$  הולך וגדל), ולכן אנחנו רוצים להכניס פרמטר נוסף לפונקציית העונשין, פרמטר אשר "ישלוט" על שיפוע פונקציית העונשין. פרמטר זה נהוג לסמן  $\mu$  והוא נקרא המכפיל. פונקציית העונשין הנשלטת נקראת

Penalty-multiplier function ומסומנת להיות  $\varphi_{p\mu}(x)$  והיא מקיימת:

$$\varphi'_{p\mu}(0) = 0$$

דוגמה לפונקציית עונשין שנשלטת ע"י מכפיל כזה היא:

$$\varphi_{p\mu} = \mu \varphi_p(x)$$

בעיית האופטימיזציה שהצגנו קודמת:

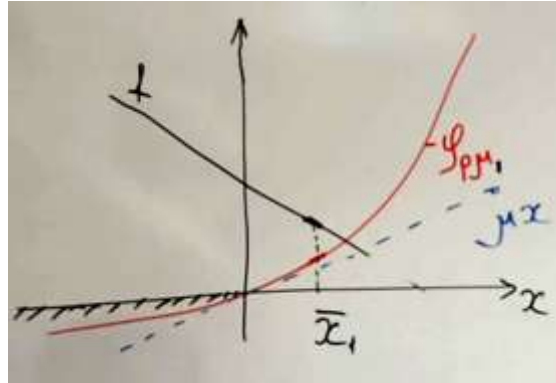
$$\min_{x \in \mathbb{R}} f(x)$$

$$\text{Subject To (S.T): } x \leq 0$$

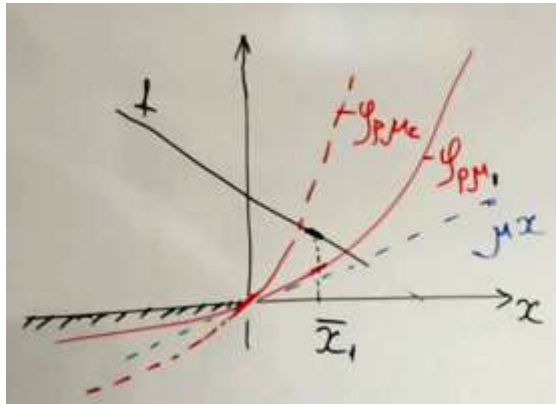
$$\text{Penalty Function: } \varphi(x) = e^x - 1$$

כאשר בכל צעד אנחנו בוחרים מכפיל  $\mu$  חדש (גדול יותר מהקודם), כאשר באיטרציה הראשונה נסמנו להיות  $\mu_1$  והבעיה שאנחנו פותרים היא:

$$\bar{x}_1 \triangleq \arg \min_x \{F_{p\mu_1}(x) = f(x) + \varphi_{p\mu_1}(x)\}$$



ובאיטרציה הבאה אנחנו בוחרים את המכפיל הבא להיות  $\mu_2 = \varphi'_{p\mu_1}(\bar{x}_1)$ , כלומר נקבל באיטרציה הבאה פונקציה עונשין שקרובה יותר לפונקציה עונשין אידיאלית כאשר השיפוע שלה בראשית הוא בדיוק  $\mu_2 = \varphi'_{p\mu_1}(\bar{x}_1)$ :



ועכשיו באיטרציה הבאה אנחנו מקבלים את הבעיה:

$$\bar{x}_2 \triangleq \arg \min_x \{F_{p\mu_2}(x) = f(x) + \varphi_{p\mu_2}(x)\}$$

ובדוגמה זו, למעשה באיטרציה השנייה כבר נקבל את הפתרון  $x_2 = 0$  (שזה הפתרון האופטימלי של הבעיה המקורית שהצגנו) וזה קרה כיוון שלפי הטענה, בנקודה האופטימלית מתקיים:  $f'(\bar{x}) = -\varphi'_{p\mu}(\bar{x})$ , וכיוון שבחרנו את  $\mu_2 = \varphi'_{p\mu_1}(\bar{x}_1)$  אז זה אומר כי פונקציה העונשין תהיה בעלת שיפוע זהה לשיפוע הפונקציה  $f(x)$  בדיוק בראשית.

#### האלגוריתם לשיטת Augmented Lagrangian Penalty – Multiplier method

נניח ובעיה האופטימיזציה המקורית היא:

$$\begin{aligned} &\min_{\bar{x} \in \mathbb{R}^n} f(\bar{x}) \\ &\text{Subject To (S.T):} \quad g_i(\bar{x}) \leq 0, \quad i = 1, \dots, m \end{aligned}$$

נגדיר את פונקציה העונשין המצרפית:

$$F_{p\mu}(\bar{x}) = f(\bar{x}) + \sum_{i=1}^m \varphi_{p\mu_i}(g_i(\bar{x}))$$



בכל איטרציה של השיטה נמצא את הצעד הבא באמצעות פתירת בעיית האופטימיזציה הבאה:

$$\bar{x}_{k+1} = \arg \min_{\bar{x}} F_{p_k \bar{\mu}_k}(\bar{x})$$

ונעדכן את המכפיל בכל איטרציה להיות:

$$\mu_{i,k+1} = \varphi'_{p_k \mu_{i,k}}(g_i(\bar{x}_{k+1}))$$

(למעשה יש לנו ווקטור של מכפילים ולעיל אנחנו רואים כיצד מעדכנים כל רכיב בווקטור המכפילים)

ובכל איטרציה אנחנו גם מגדילים את פרמטר העונשין  $p_k$  באופן הבא:  $p_{k+1} = \alpha \cdot p_k$  כאשר  $2 \leq \alpha \leq 10$ .

בעיות אופטימיזציה עם אילוצים מסוג Equality Constraints (בלבד):

נניח ובעיית האופטימיזציה היא (במימד אחד):

$$\min_{x \in \mathbb{R}} f(x)$$

$$\text{Subject To } (S.T): \quad h_j(x) = 0, \quad j = 1, \dots, k$$

במקרה זה, אנחנו נרצה פונקציית עונשין שתקבל ערכים חיוביים גם עבור ערכי נקודות  $x$  שחורגות מהכיוון החיובי וגם מהכיוון השלילי של האילוץ  $h_j(x) = 0$ , ולכן הכי מתבקש להגדיר פונקציית עונשין מהצורה:

$$\varphi: \mathbb{R} \rightarrow \mathbb{R}, \quad \varphi(t) = \frac{1}{2} t^2$$

ואם נשתמש בשיטת Augmented Lagrangian Penalty – Multiplier method אז נקבל את פונקציית העונשין:

$$\varphi_{p\mu}(t) = \mu t + \frac{p}{2} t^2$$

ופונקציית העונשין המצרפית במקרה זה היא:

$$F_{p\mu}(x) = f(x) + \sum_{i=1}^m \varphi_{p\mu_i}(h_i(x)) = f(x) + \mu^T \bar{h}(x) + \frac{p}{2} \|\bar{h}(x)\|^2$$

כמו כן, בדוגמה זו, המכפיל (הבעיה היא ממימד אחד בלבד) משתנים מאיטרציה לאיטרציה בצורה הבאה:

$$\varphi'_{p\mu}(t) = \mu + pt \quad \Rightarrow \quad \mu_{k+1} = \varphi'_{p_k \mu_k}(h(x_{k+1})) = \mu_k + p \cdot h(x_{k+1})$$

## **הרצאה 14 – Lagrange Multipliers and Penalty Function Method. Augmented Lagrangian**

בהרצאה הקודמת הגדרנו את פונקצית העונשין המצרפית  $(F_p(\bar{x}))$ . ניזכר כי בעית האופטימיזציה הכללית (עם אילוצים) היא:

$$\arg \min_{\bar{x} \in \mathbb{R}^n} F_p(\bar{x}) = \arg \min_{\bar{x} \in \mathbb{R}^n} \left( f(\bar{x}) + \sum_{i=1}^m \varphi_p(g_i(\bar{x})) \right)$$

כאשר נסמן את הפתרון של בעית אופטימיזציה כללית זו להיות  $\bar{x}_p$ . נגדיר תנאים הכרחיים לאופטימליות של פתרון שנציע.

תנאי מסדר ראשון (על הנגזרת הראשונה) לפתרון אופטימלי (Karush-Kuhn-Tucker) (תנאי KKT):

כמו שראינו בעבר עבור בעית אופטימיזציה ללא אילוצים, תנאי KKT הוא התאפסות פגרדיאנט פונקצית המטרה בנקודה האופטימלית. כיוון שהמרנו את בעית האופטימיזציה שהייתה לנו עם אילוצים לבעית אופטימיזציה ללא אילוצים, אז כעת בעיה האופטימיזציה ללא אילוצים שיש לנו היא פונקצית העונשין המצרפית ולכן תנאי KKT הוא שוב התאפסות הגרדיאנט בנקודת הפתרון:

$$\nabla F_p(\bar{x}_p) = \nabla f(\bar{x}_p) + \sum_{i=1}^m \varphi'_p(g_i(\bar{x}_p)) \cdot \nabla g_i(\bar{x}_p) = 0$$

נסמן את הסקלר:  $\lambda_{i,p}(\bar{x}_p) \triangleq \varphi'_p(g_i(\bar{x}_p))$  ואז נקבל:

$$\nabla F_p(\bar{x}_p) = \nabla f(\bar{x}_p) + \sum_{i=1}^m \lambda_{i,p}(\bar{x}_p) \cdot \nabla g_i(\bar{x}_p) = 0$$

נבחין כי הביטוי האחרון הוא בדיוק הגרדיאנט של הלגרנג'יאן שראינו בהרצאה הקודמת. בנוסף, ראינו כי כאשר  $p \rightarrow \infty$  אז אנחנו מקבלים כי הפתרון שאנחנו מציעים שואף לפתרון האופטימלי, כלומר  $\bar{x}_p \rightarrow \bar{x}^*$ , כלומר הפתרון  $\bar{x}_p$  מתכנס לפתרון האמיתי של בעית האופטימיזציה המקורית עם אילוצים.

ננתח מה הם הסקלרים  $\lambda_{i,p}(\bar{x}_p) = \varphi'_p(g_i(\bar{x}_p))$  (לכל  $1 \leq i \leq m$ ) כאשר  $p \rightarrow \infty$ . עבור אילוץ לא פעיל (Non-Active Constrain) מתקיים (לפי הדרישה על פונקצית העונשין):

$$\lim_{p \rightarrow \infty} \lambda_{i,p}(\bar{x}_p) = \varphi'_p(g_i(\bar{x}_p)) = 0$$

ועבור אילוצים פעילים (Active Constrains) נביט על הגרדיאנט של הלגרנג'יאן בצורה מטריצית. נביט על המטריצה  $\nabla G(\bar{x})$  שעמודותיה הן הגרדיאנטים של כל אחד מהאילוצים:

$$\nabla G(\bar{x}) = \begin{pmatrix} \nabla g_{i_1}(\bar{x}) & \nabla g_{i_2}(\bar{x}) & \cdots & \cdots & \nabla g_{i_l}(\bar{x}) & \nabla g_{i_{l+1}}(\bar{x}) \\ \underbrace{\qquad\qquad\qquad}_{\substack{\nabla G_A(\bar{x}) \\ A \equiv \text{Active}}} & \underbrace{\qquad\qquad\qquad}_{\substack{\nabla G_N(\bar{x}) \\ N \equiv \text{Non-Active}}} \end{pmatrix}$$

ואז מהדרישה של התנאי מסדר ראשון:  $\nabla f(\bar{x}_p) + \sum_{i=1}^m \lambda_{i,p}(\bar{x}_p) \cdot \nabla g_i(\bar{x}_p) = 0$  נקבל בצורה מטריצית:

$$-\nabla f(\bar{x}_p) = \sum_{i=1}^m (\nabla G_A(\bar{x}_p) \lambda_{i,p}(\bar{x}_p) + \nabla G_N(\bar{x}_p) \lambda_{i,p}(\bar{x}_p))$$

עבור אילוצים לא פעילים, אנחנו מקבלים כי נגזרת פונקציית העונשין  $\varphi'_p(\bar{x}_p)$  היא שואפת לאפס כיוון שאין חריגה בתחום זה (אילוץ לא פעיל) אך עבור אילוצים פעילים, אנחנו נמצאים בנקודות  $\bar{x}_p$  שיכולות לחרוג (לפחות לפני שאנחנו מגיעים לפתרון האופטימלי) ולכן מעניינת אותנו הנגזרת  $\varphi'_p(\bar{x}_p)$ , והיא אף יותר מעניינת כאשר  $p \rightarrow \infty$  כי אז פונקציית העונשין מתקרבת לפונקציית עונשין אידיאלית והנגזרת שואפת לאינסוף או אפילו נהיית לא מוגדרת. נראה כי הנגזרת תלויה בהאם המטריצה  $\nabla G(\bar{x})$  היא מטריצה מדרגה מלאה, או שבמילים אחרות אפשר לומר כי המטריצה היא Full Column Rank.

עבור  $p \rightarrow \infty$  ראינו כי עבור אילוצים לא פעילים מתקיים  $\lim_{p \rightarrow \infty} \lambda_{i,p}(\bar{x}_p) = 0$  ולכן  $\lim_{p \rightarrow \infty} \nabla G_N(\bar{x}_p) \lambda_{i,p}(\bar{x}_p) = 0$  כלומר מהמשוואה האחרונה אנחנו מקבלים כעת:

$$-\nabla f(\bar{x}_p) = \sum_{i=1}^m \nabla G_A(\bar{x}_p) \lambda_{i,p}(\bar{x}_p)$$

אם המטריצה  $\nabla G(\bar{x})$  היא בעלת דרגה מלאה (ואז עמודותיה בלתי תלויות לינארית) חייב להיות פתרון יחיד למשוואה האחרונה וכן עבור שינויים "קטנים" במטריצה  $\nabla G(\bar{x})$  נקבל כי השינויים בפתרון  $(\lambda_{i,p}(\bar{x}_p))$  יהיו "קטנים" גם כן. כמו כן, כאשר  $p \rightarrow \infty$  נקבל  $\lambda_{i,p}(\bar{x}_p) \rightarrow \lambda_{i,p}^*(\bar{x}_p)$ , כלומר נקבל את הפתרון האופטימלי.

שיטת פונקציות עונשין (Penalty function method) עבור אילוצים מסוג Equality Constraints:

ניזכר כי בעית האופטימיזציה עם אילוצים מסוג Equality Constraints היא מהצורה:

$$\begin{aligned} \min_{\bar{x} \in \mathbb{R}^n} f(\bar{x}) \\ \text{Subject To (S.T):} \quad h_j(\bar{x}) \leq 0, \quad j = 1, \dots, l \end{aligned}$$

במקרה זה, אנחנו נרצה פונקציית עונשין שתקבל ערכים חיוביים גם עבור ערכי נקודות  $x$  שחורגות מהכיוון החיובי וגם מהכיוון השלילי של האילוץ  $h_j(x) = 0$ , ולכן הכי מתבקש להגדיר פונקציית עונשין מהצורה:

$$\psi: \mathbb{R} \rightarrow \mathbb{R}, \quad \psi(t) = \frac{1}{2} t^2$$

בנוסף, נרצה גם כאן להגדיר פרמטר עונשין,  $p$ , אשר יקיים:

$$\lim_{p \rightarrow \infty} \psi_p(t) = \begin{cases} 0 & , t = 0 \\ \infty & , t \neq 0 \end{cases}$$

למשל ניקח את פונקציית העונשין:

$$\psi(t) = t^2$$

ונוסיף לה פרמטר עונשין בצורה הבאה:

$$\psi_p(t) = p\psi(t)$$

או למשל:

$$\psi_p(t) = \frac{1}{p}\psi(pt)$$

ונבחין כי עבור  $\psi(t) = t^2$ , שתי הדוגמאות עבור  $\psi_p(t)$  תהיינה זהות ונקבל:

$$\psi_p(t) = pt^2$$

ניתן כעת גם להגדיר את פונקציית העונשין המצרפית:

$$F_p(\bar{x}) = f(\bar{x}) + \sum_{j=1}^l \psi_p(h_j(\bar{x}))$$

בצורה דומה למקודם, נוכל להגדיר תנאי מסדר ראשון הכרחי לאופטימליות של פתרון. אם נניח כי הפתרון האופטימלי הוא  $\bar{x}^*$  אז קיים פתרון  $\lambda^*(\bar{x}^*) \in \mathbb{R}^l$  כך שמתקיים:

$$\nabla F_p(\bar{x}^*) = f'(\bar{x}^*) + \sum_{j=1}^l \lambda_j^*(\bar{x}^*) \cdot h_j'(\bar{x}^*) = 0$$

ראינו שני תנאים מסדר ראשון, כביכול שונים עבור מקרה של אילוצים מסוג Inequality ומסוג Equality אך למעשה אפשר להכליל את שני התנאים באמצעות זה שנאמר כי התנאי הוא כי הגרדיאנטים של האילוצים הפעילים הם בלתי תלויים לינארית וכופלי לגרנז' עבור אילוצים מסוג Inequality יהיו אי-שלייליים ועבור אילוצים מסוג Equality כופלי לגרנז' יכולים להיות חיוביים או שליליים.

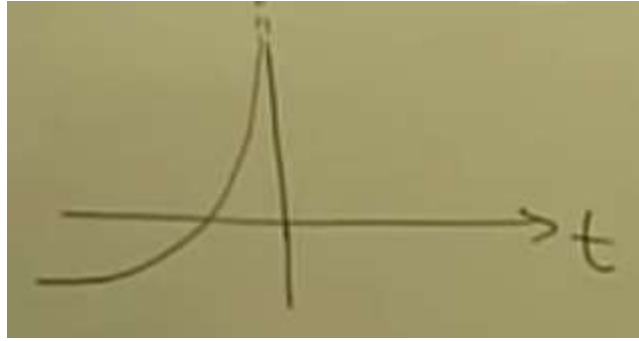
מקרה פרטי של פונקציית עונשין, פונקציית Barrier, (שיטת Barrier, שיטת המחסום):

נניח ובעית האופטימיזציה היא:

$$\begin{aligned} \min_{x \in \mathbb{R}^n} f(\bar{x}) \\ \text{Subject To } (S.T): \quad g_i(\bar{x}) \leq 0, \quad i = 1, \dots, m \end{aligned}$$

נניח ועבור נקודות  $\bar{x}$  שלא עומדות באילוצים  $g_i(\bar{x}) \leq 0$ , הפונקציות (האילוצים)  $g_i(\bar{x})$  כלל לא ניתנים להגדרה.

נגדיר את פונקציית המחסום (Barrier) במשתנה יחיד:



כלומר באפס יש לה אסימפטוטה ששואפת לאינסוף. ומשפחה פרמטרית של פונקציות מסוג כזה תהיה בעלת פרמטר  $p$ . עבור  $p \rightarrow \infty$  נרצה כי הפונקציה תקיים:

$$\lim_{\substack{p \rightarrow \infty \\ t \rightarrow 0^-}} \varphi_p(t) = \begin{cases} \infty & , t = 0 \\ 0 & , t < 0 \end{cases}$$

ניתן גם להגדיר באופן הבא:

$$\begin{aligned} \lim_{t \rightarrow 0^-} \varphi(t) &= \infty \\ \varphi_p(t) &= \frac{1}{p} \varphi(t) \end{aligned}$$

דוגמה לפונקצית מחסום:

$$\varphi(t) = -\log(-t)$$

פונקצית העונשין המצרפית, כמו קודם, היא:

$$F_p(\bar{x}) = f(\bar{x}) + \sum_{i=1}^m \varphi_p(g_i(\bar{x}))$$

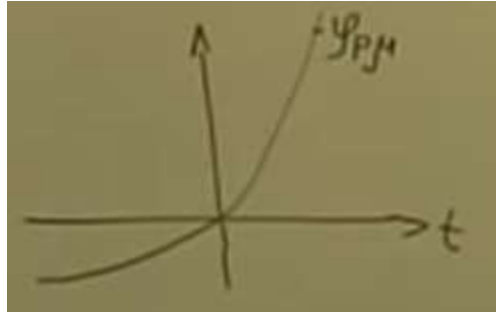
היתרון העיקרי בשיטה זו (Barrier Method) הוא שהנקודה  $\bar{x}_p$  תמיד נמצאת בתחום, כלומר תמיד עומדת בדרישת האילון, כלומר תמיד נמצאת ב-Feasible Set.

הערה:

כיוון שפונקצית העונשין אינה מוגדרת עבור  $t > 0$  אסור לחרוג מתחום זה בעת פתרון בעית האופטימיזציה, ולכן למשל כאשר מבצעים Line Search צריך לוודא כי לא חורגים מהתחום המותר.

הרחבה של שיטת פונקצית העונשין, שיטת Augmented Lagrangian method:

המטרה בשיטה זו היא למצוא פתרון מדויק לבעית האופטימיזציה (עם אילוצים) ללא השאפת פרמטר העונשין  $p$  לאינסוף. נגדיר פונקצית עונשין נשלטת הנקראת Penalty-multiplier function ומסומנת להיות  $\varphi_{p\mu}(t)$ . פונקציה זו מקיימת את הצורה:



מטרת הפרמטר הנוסף  $\mu$  הוא לדרוש כי הנגזרת הראשונה של פונקציית העונשין בראשית תהיה בדיוק:

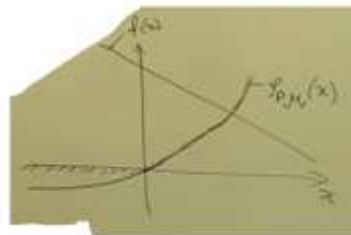
$$\varphi'_{p\mu}(0) = \mu$$

בקרב נראה איך דרישה זו מתקשרת לשיטת האופטימיזציה באמצעות כופלי לגרנז'.

תחילה נבין את השיטה לבעיית אופטימיזציה ממימד אחד. נניח כי בעיית האופטימיזציה היא:

$$\begin{aligned} \arg \min_{x \in \mathbb{R}} f(x) \quad , \quad x \in \mathbb{R} \\ \text{Subject To (S.T):} \quad x \leq 0 \end{aligned}$$

ובנוסף נניח כי  $f(x)$  היא פונקציית לינארית יורדת. ברור כי הפתרון האופטימלי הוא  $x^* = 0$  אך נרצה להגיע לפתרון זה באמצעות פונקציית העונשין שהגדרנו. בצורה גרפית נקבל: (באיטרציה הראשונה של האלגוריתם פונקציית העונשין היא  $\varphi_{p\mu_1}$ )



נגדיר את פונקציית העונשין המצרפית:

$$F_p(x, \mu) = f(x) + \varphi_{p\mu}(x)$$

כעת אנחנו מחפשים את הנקודה  $x_1$  אשר בה מתקיים  $-f'(x) = \varphi'_{p\mu_1}(x)$  (לפי התנאי מסדר ראשון, כלומר זה התנאי שגרדיאנט הפונקציה המצרפית יתאפס). כעת, מטרנו היא למצוא פונקציית עונשין חדשה, לאיטרציה הבאה באלגוריתם,  $\varphi_{p\mu_2}$  אשר תקרב אותנו לפתרון האופטימלי שהוא כאמור  $x^* = 0$ . כיוון שאנחנו יודעים כי הפונקציה אותה אנחנו ממזערים היא לינארית, אז אם נבחר פונקציית עונשין אשר בראשית שלה יש שיפוע זהה והפוך בסימן לנגזרת של הפונקציה  $f(x)$  נקבל את הפתרון האופטימלי כבר באיטרציה הבאה. לפי האמור לעיל, נגדיר באיטרציה הבאה כי פונקציית העונשין בראשית, היא בעלת שיפוע שזהה לשיפוע של פונקציית העונשין בנקודה  $x_1$ . באופן פורמלי הגדרנו את הפרמטר  $\mu$  באיטרציה הבאה לפי:

$$\mu_2 = \varphi'_{p\mu_2}(0) = \varphi'_{p\mu_1}(x_1)$$

כעת נעבור לבעיית אופטימיזציה במספר מימדים. נניח ובעיית האופטימיזציה הבא:

$$\min_{\bar{x} \in \mathbb{R}^n} f(\bar{x})$$

$$\text{Subject To (S.T):} \quad g_i(\bar{x}) \leq 0, \quad i = 1, \dots, m$$

ופונקצית העונשין המצרפית היא:

$$F_p(\bar{x}, \bar{\lambda}) = f(\bar{x}) + \sum_{i=1}^m \varphi_{p\lambda_i}(g_i(\bar{x}))$$

האלגוריתם Augmented Lagrangian:

באיטרציה ה- $k$  אנחנו פותרים את בעיית האופטימיזציה:

$$\bar{x}_{k+1} = \arg \min_{\bar{x} \in \mathbb{R}^n} F_p(\bar{x}, \bar{\lambda}_k)$$

מעדכנים את כופלי לגרנז' באופן הבא:

$$\lambda_{i,k+1} = \varphi'_{p\lambda_{i,k}}(g_i(\bar{x}_k))$$

מעדכנים את פרמטר העונשין באופן הבא:

$$p_{k+1} = \min\{\alpha p_k, p_{\max}\}, \quad \alpha > 1, \quad 2 \leq \alpha \leq 10, \quad p_{\max} \in [100, 1000]$$

(את פרמטר העונשין אנחנו משנים באופן מתון. אמנם ככל שהפרמטר גדול יותר כך אנחנו מקבלים פתרון מדויק יותר אך לוקח גם יותר זמן להגיע אליו בכל איטרציה ולכן אנחנו עושים פשרה בין דיוק הפתרון לבין המהירות שמגיעים אליו).

הערה:

על מנת שהאלגוריתם יעבוד בצורה יעילה, רצוי כי הפרמטר  $\lambda_{i,k+1}$  לא ישתנה בצורה חדה מידי ולכן כדאי להגביל את השינוי בו בתחום:

$$\frac{1}{3} \leq \frac{\lambda_{i,k+1}}{\lambda_{i,k}} \leq 3$$

הערה:

השיטה Augmented Lagrangian עובדת עם אילוצים מסוג Inequality, אך מה נעשה כאשר בעיית האופטימיזציה היא עם אילוצים מסוג Equality? נוכל לומר כי מתקיים:

$$h(x) = 0 \quad \Leftrightarrow \quad h(x) \leq 0 \text{ and } h(x) \geq 0$$

ועכשיו קיבלנו אילוצים מסוג Inequality, כלומר נפתור שתי בעיות אופטימיזציה עם אילוצים מסוג Inequality. אפשרות נוספת לפתור בעיית אופטימיזציה עם אילוצים מסוג Equality נוכל להגדיר Penalty-multiplier function מהצורה:

$$\varphi_{p\mu}(t) = pt^2 + \mu t$$

וקל לראות כי בראשית נקבל כי השיפוע הוא תמיד  $\mu$  כפי שנדרש בשיטת Augmented Lagrangian:

$$\varphi'_{p\mu}(0) = \mu$$

סיכום שיטת Augmented Lagrangian:

בשיטת Augmented Lagrangian ראינו ניתן למזער את פונקציית העונשין המצרפית  $F_p(\bar{x}, \bar{\lambda}^*)$ , ולהגיע לפתרון בעית האופטימיזציה ללא השאפת פרמטר העונשין,  $p$  לאינסוף. התנאי מסדר ראשון של פונקציית העונשין המצרפית:

$$\nabla_{\bar{x}} F_p(\bar{x}, \bar{\lambda}) = \nabla_{\bar{x}} f(\bar{x}) + \sum_{i=1}^m \varphi'_{p\lambda_i}(g_i(\bar{x})) \cdot \nabla_{\bar{x}} g_i(\bar{x}) = 0$$

ובפתרון האופטימלי הוא למעשה:

$$\nabla_{\bar{x}} F_p(\bar{x}^*, \bar{\lambda}^*) = \nabla_{\bar{x}} L(\bar{x}^*, \bar{\lambda}^*) = 0$$

וזה אכן נכון כי עבור אילוצים פעילים נקבל:

$$g_i(\bar{x}^*) = 0 \quad \Rightarrow \quad \varphi'_{p\lambda_i^*}(g_i(\bar{x}^*)) = \varphi'_{p\lambda_i^*}(0) \triangleq \lambda_i^*$$

ועבור אילוצים לא פעילים, מתקיים  $\lambda_i^* = 0$  ולכן נקבל:

$$g_i(\bar{x}^*) < 0 \quad \Rightarrow \quad \varphi'_{p\lambda_i^*}(g_i(\bar{x}^*)) = 0 \triangleq \lambda_i^*$$



## הרצאה 15 – Minmax, Game Theory, Lagrangian Duality

נרצה למצוא ביטוי שיבטא את דואליות הלגרנז'יאן נותנת לנו מבט נוסף על בעיות אופטימיזציה רבות ותעזור לנו לפשט את הבעיות. דואליות הלגרנז'יאן תאפשר לנו לדעת עד כמה אנחנו קרובים לפתרון וכן לתת גבול תחתון לפונקציה שאותה אנחנו מנסים למזער. דואליות הלגרנז'יאן היא הבסיס להרבה שיטות אופטימיזציה יעילות.

משפט Minmax:

נניח כי פונקציית המטרה היא בעלת שני משתנים, כאשר כל משתנה הוא ווקטור, כלומר  $f: (\mathbb{R}^n, \mathbb{R}^n) \rightarrow \mathbb{R}$ . מתקיים:

$$\min_{\bar{z}} \left( \max_{\bar{\omega}} f(\bar{z}, \bar{\omega}) \right) \geq \max_{\bar{\omega}} \left( \min_{\bar{z}} f(\bar{z}, \bar{\omega}) \right)$$

פירוש למשפט ה-Minmax בתורת המשחקים:

נתאר מהלך משחק:

נניח כי שחקן  $z$  בחר נקודה  $\bar{z} \in Z$  ושחקן  $\omega$  בחר נקודה  $\bar{\omega} \in \Omega$ . לאחר מכן, המשחק אומר כי שחקן  $z$  משלם לשחקן  $\omega$  סכום של  $f(\bar{z}, \bar{\omega})$  שקלים (ורק שחקן  $z$  משלם ל- $\omega$ ). המטרה של שחקן  $\omega$  לקבל הכי הרבה כסף שהוא יכול והמטרה של שחקן  $z$  לשלם כמה שפחות לשחקן  $\omega$ .

במשחק יכולות להיות שתי אפשרויות לשחקן שמשחק ראשון, או שחקן  $\omega$  או שחקן  $z$ .

- שחקן  $\omega$  משחק שני: אם שחקן  $\omega$  משחק שני (כלומר קודם שחקן  $z$  בחר נקודה  $\bar{z}$  ורק לאחר מכן שחקן  $\omega$  בחר נקודה  $\bar{\omega}$ ) אז שחקן  $z$  (שבוחר ראשון) יבחר נקודה  $\bar{z}$  כך יתקיים:  $\min_{\bar{z}} f(\bar{z}, \bar{\omega})$  (שחקן  $z$  למעשה חוזה כי שחקן  $\omega$  יבחר את נקודה  $\bar{\omega}$  כך שיתקיים  $\max_{\bar{\omega}} \left( \min_{\bar{z}} f(\bar{z}, \bar{\omega}) \right)$ ). כעת שחקן  $\omega$  כבר יודע מה שחקן  $z$  בחר, ולכן רוצה למקסם את  $f(\bar{z}, \bar{\omega})$  (כלומר למקסם את הסכום שהוא יקבל) ולכן יבחר  $\max_{\bar{\omega}} \left( \min_{\bar{z}} f(\bar{z}, \bar{\omega}) \right)$  (בדיוק כפי ששחקן  $z$  "חזה").
- שחקן  $\omega$  משחק ראשון: אם שחקן  $\omega$  משחק ראשון, אז הוא כמובן יבחר את הנקודה  $\max_{\bar{\omega}} f(\bar{z}, \bar{\omega})$  (בכדי לקבל את כמות הכסף הגדולה ביותר האפשרית מבחינתו), ואז שחקן  $z$ , ירצה לבחור נקודה כך שכמות הכסף שהוא יצטרך לשלם תקטן ולכן יבחר  $\min_{\bar{z}} \left( \max_{\bar{\omega}} f(\bar{z}, \bar{\omega}) \right)$  ולכן שחקן  $z$  ישלם סך הכול  $\min_{\bar{z}} \left( \max_{\bar{\omega}} f(\bar{z}, \bar{\omega}) \right)$ .

לפי משפט Minmax:

$$\min_{\bar{z}} \left( \max_{\bar{\omega}} f(\bar{z}, \bar{\omega}) \right) \geq \max_{\bar{\omega}} \left( \min_{\bar{z}} f(\bar{z}, \bar{\omega}) \right)$$

ולכן אם שחקן  $\omega$  משחק ראשון אז הוא יקבל סכום גדול יותר. (צריך לוודא שזה אכן הניתוח הנכון למשחק)

משפט נקודת האוכף (Saddle point):

אם קיימת נקודה  $(\bar{z}^*, \bar{\omega}^*)$  כך שמתקיים:

$$\forall \bar{z}, \bar{\omega} \in \mathbb{R}^n: f(\bar{z}^*, \bar{\omega}) \leq f(\bar{z}^*, \bar{\omega}^*) \leq f(\bar{z}, \bar{\omega}^*)$$

אז מתקיים:

$$\min_{\bar{z}} \left( \max_{\bar{\omega}} f(\bar{z}, \bar{\omega}) \right) = \max_{\bar{\omega}} \left( \min_{\bar{z}} f(\bar{z}, \bar{\omega}) \right)$$

(דוגמה: לפונקציה  $f(z, \omega) = z^2 - \omega^2$  יש נקודת אוכף)

הוכחה:

נתון לנו כי קיימת נקודה  $(\bar{z}^*, \bar{\omega}^*)$  כך שמתקיים:

$$f(\bar{z}^*, \bar{\omega}) \underset{(left)}{\leq} f(\bar{z}^*, \bar{\omega}^*) \underset{(right)}{\leq} f(\bar{z}, \bar{\omega}^*)$$

מתקיים:

$$\min_{\bar{z}} \left( \max_{\bar{\omega}} f(\bar{z}, \bar{\omega}) \right) \leq \max_{\bar{\omega}} f(\bar{z}^*, \bar{\omega}) \underset{(left)}{=} f(\bar{z}^*, \bar{\omega}^*) \underset{(right)}{=} \min_{\bar{z}} f(\bar{z}^*, \bar{\omega}) \leq \max_{\bar{\omega}} \left( \min_{\bar{z}} f(\bar{z}^*, \bar{\omega}) \right)$$

כלומר קיבלנו כי:

$$\min_{\bar{z}} \left( \max_{\bar{\omega}} f(\bar{z}, \bar{\omega}) \right) \leq \max_{\bar{\omega}} \left( \min_{\bar{z}} f(\bar{z}^*, \bar{\omega}) \right)$$

אך המשפט Minmax אומר כי מתקיים:

$$\min_{\bar{z}} \left( \max_{\bar{\omega}} f(\bar{z}, \bar{\omega}) \right) \geq \max_{\bar{\omega}} \left( \min_{\bar{z}} f(\bar{z}, \bar{\omega}) \right)$$

ולכן בהכרח מתקיים:

$$\min_{\bar{z}} \left( \max_{\bar{\omega}} f(\bar{z}, \bar{\omega}) \right) = \max_{\bar{\omega}} \left( \min_{\bar{z}} f(\bar{z}, \bar{\omega}) \right)$$

מש"ל.

אי-שיוויון Minmax של הלגרנז'יאן, Minmax of Lagrangian:

כעת נראה את הקשר בין משפט ה-Minmax לבין פתרון בעיות אופטימיזציה עם אילוצים. נניח ובעית האופטימיזציה היא:

$$\min_{x \in \mathbb{R}^n} f(\bar{x})$$

$$Subject To (S.T): \quad g_i(\bar{x}) \leq 0, \quad i = 1, \dots, m$$

נזכיר כי הלגרנז'יאן הוא:

$$L(\bar{x}, \bar{\lambda}) = f(\bar{x}) + \sum_{i=1}^m \lambda_i g_i(\bar{x})$$

נביט בביטוי הבא:

$$\max_{\substack{\lambda_i \geq 0 \\ \forall i \in [1, m]}} L(\bar{x}, \bar{\lambda}) = \max_{\substack{\lambda_i \geq 0 \\ \forall i \in [1, m]}} f(\bar{x}) + \sum_{i=1}^m \lambda_i g_i(\bar{x})$$

אם  $\bar{x} \in Feasible Set$  אז כל האילוצים חלים עליו ולכן  $\forall 1 \leq i \leq m: g_i(\bar{x}) \leq 0$  ולכן  $\sum_{i=1}^m \lambda_i g_i(\bar{x}) \leq 0$  (כאשר נזכור כי  $\lambda_i \geq 0$ ).

ולכן למעשה קיבלנו כי אם  $\sum_{i=1}^m \lambda_i g_i(\bar{x}) \leq 0$  אז  $L(\bar{x}, \bar{\lambda})$  הוא חסם תחתון של פונקציית המטרה  $f(\bar{x})$ . מכל האמור לעיל נוכל לקבל (ע"י בחירת  $\lambda_i$  מתאימים):

$$\max_{\substack{\lambda_i \geq 0 \\ \forall i \in [1, m]}} L(\bar{x}, \bar{\lambda}) = \begin{cases} f(\bar{x}) & , \quad x \in Feasible Set \\ \infty & , \quad x \notin Feasible Set \end{cases}$$

למעשה קיבלנו את פונקציית העונשין המצרפית האידיאלית שסימנו בעבר להיות  $F_\infty(\bar{x})$ . מכאן, אם נניח כי הפתרון האופטימלי הוא  $\bar{x}^*$  אז נקבל כי ערך הפונקציה האופטימלי הוא:

$$f(\bar{x}^*) = \min_{\bar{x}} \max_{\substack{\lambda_i \geq 0 \\ \forall i \in [1, m]}} L(\bar{x}, \bar{\lambda})$$

לפי משפט ה-Minmax מתקיים:

$$f(\bar{x}^*) = \min_{\bar{x}} \max_{\substack{\lambda_i \geq 0 \\ \forall i \in [1, m]}} L(\bar{x}, \bar{\lambda}) \geq \max_{\substack{\lambda_i \geq 0 \\ \forall i \in [1, m]}} \min_{\bar{x}} L(\bar{x}, \bar{\lambda})$$

נסמן  $\eta(\bar{\lambda}) = \min_{\bar{x}} L(\bar{x}, \bar{\lambda})$ , ופונקציה זו נקראת Dual function. לסיכום קיבלנו:

$$f(\bar{x}^*) \geq \max_{\substack{\lambda_i \geq 0 \\ \forall i \in [1, m]}} \eta(\bar{\lambda})$$

הבעיה של למקסם את הפונקציה הדואלית  $\eta(\bar{\lambda})$ , נקראת הבעיה הדואלית של מיזעור הלגרנז'יאן. אי-השיוויון

$$f(\bar{x}^*) \geq \max_{\substack{\lambda_i \geq 0 \\ \forall i \in [1, m]}} \eta(\bar{\lambda})$$

נקרא משפט הדואליות החלש (בדואליות חזקה נקבל אי-שיוויון ממש).

משפט הדואליות ה"חזקה":

נניח ובעית האופטימיזציה היא:

$$\min_{x \in \mathbb{R}^n} f(\bar{x})$$

$$Subject To (S.T): \quad g_i(\bar{x}) \leq 0, \quad i = 1, \dots, m$$

אם מתקיים כי  $f(\bar{x}), g_i(\bar{x})$  (לכל  $i = 1, \dots, m$ ) הן פונקציות קמורות וגם תנאי מסדר ראשון (Karush-Kuhn-Tucker) (KKT) מתקיים, אז מתקיים אי-שוויון ממש, כלומר:

$$f(\bar{x}^*) > \max_{\substack{\lambda_i \geq 0 \\ \forall i \in [1, m]}} \eta(\bar{\lambda})$$

הוכחה:

נביט על הלגרנז'יאן:

$$L(\bar{x}, \bar{\lambda}) = f(\bar{x}) + \sum_{i=1}^m \lambda_i g_i(\bar{x})$$

הלגרנז'יאן הוא פונקציה קמורה כאשר  $\lambda_i \geq 0$  לכל  $1 \leq i \leq m$ .

בנוסף, לפי התנאי ההכרחי מסדר ראשון מתקיים:

$$\nabla_{\bar{x}} L(\bar{x}^*, \bar{\lambda}^*) = 0$$

אבל עבור פונקציה שהיא קמורה, התנאי ההכרחי מסדר ראשון הוא גם תנאי מספיק להיותה של הנקודה  $\bar{x}^*$  להיות הפתרון האופטימלי ולכן מתקיים:

$$\bar{x}^* = \arg \min_{\bar{x} \in \mathbb{R}^n} L(\bar{x}, \bar{\lambda}^*)$$

ולכן מתקיים:

$$L(\bar{x}, \bar{\lambda}^*) \geq L(\bar{x}^*, \bar{\lambda}^*)$$

וכן לפי התנאי מסדר ראשון, KKT, ראינו כי מתקיים Complimentary Slackness, כלומר:

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

בנוסף, כיוון שדרשנו  $\lambda_i \geq 0$  לכל  $1 \leq i \leq m$ , אז בנקודה  $\bar{x}^* \in Feasible Set$ , האילוצים מתקיימים, כלומר  $g_i(\bar{x}^*) \leq 0$  לכל  $1 \leq i \leq m$  ולכן למעשה מתקיים:

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

ולכן מתקיים:

$$L(\bar{x}^*, \bar{\lambda}^*) \geq L(\bar{x}^*, \bar{\lambda})$$

ומכאן נוכל לראות כי הנקודה  $(\bar{x}^*, \bar{\lambda}^*)$  היא נקודת אוקף של הלגרנז'יאן כי הראנו כי מתקיים:

$$L(\bar{x}, \bar{\lambda}^*) \geq L(\bar{x}^*, \bar{\lambda}^*) \geq L(\bar{x}^*, \bar{\lambda})$$

כעת, כיוון שהנקודה  $(\bar{x}^*, \bar{\lambda}^*)$  היא נקודת אוכף של הלגרנז'יאן,  $L(\bar{x}, \bar{\lambda})$ , אז לפי משפט נקודת האוכף מתקיים:

$$f(\bar{x}^*) = \min_{\bar{x}} \left( \max_{\substack{\lambda_i \geq 0 \\ i \in [1, m]}} L(\bar{x}, \bar{\lambda}) \right) = \max_{\substack{\lambda_i \geq 0 \\ i \in [1, m]}} \left( \underbrace{\min_{\bar{x}} L(\bar{x}, \bar{\lambda})}_{\eta(\bar{\lambda})} \right)$$

מש"ל.

הערה: (Slater Condition)

קיים תנאי מספיק נוסף (גם התנאים שהוכחנו במשפט הדואליות החזקה היו תנאים מספיקים). אם קיים  $\bar{x} \in Feasible Set$  אשר מקיים  $g_i(\bar{x}) < 0$  לכל אילוץ לא לינארי  $g_i$  (עבור אילוץ לינארי אין שום צורך ש-  $\bar{x} \in Feasible Set$  יקיים משהו) אז מתקיימת דואליות חזקה.

דוגמאות של בעיות דואליות:

בעיה מספר 1: בעיה ריבועית

נתונה בעית האופטימיזציה הבאה:

$$\min_{\bar{x} \in \mathbb{R}^n} f(\bar{x}) = \frac{1}{2} \bar{x}^T x$$

$$Subject To (S.T): \quad A\bar{x} - \bar{b} \leq 0$$

הלגרנז'יאן עבור בעיה זו הוא:

$$L(\bar{x}, \bar{\lambda}) = \frac{1}{2} \bar{x}^T \bar{x} + \bar{\lambda}^T (A\bar{x} - \bar{b})$$

(ניזכר כי הלגרנז'יאן נתון ע"י:  $L(\bar{x}, \bar{\lambda}) \triangleq f(\bar{x}) + \sum_{i \in I} \lambda_i g_i(\bar{x}^*)$  כאשר  $g_i(\bar{x})$  הם האילוצים)

ככלל, הפתרון לבעיה נתון ע"י:

$$\min_{\bar{x}} \left( \max_{\bar{\lambda}} L(\bar{x}, \bar{\lambda}) \right)$$

אך לפי משפט Minmax אנחנו יודעים כי מתקיים:

$$\min_{\bar{x}} \left( \max_{\bar{\lambda}} L(\bar{x}, \bar{\lambda}) \right) \geq \max_{\bar{\lambda}} \left( \min_{\bar{x}} L(\bar{x}, \bar{\lambda}) \right)$$

ולכן אם נפתור את הבעיה  $\max_{\bar{\lambda}} \left( \min_{\bar{x}} L(\bar{x}, \bar{\lambda}) \right)$  אז בוודאות נקבל פתרון שאפילו ייתכן שהוא טוב יותר מהפתרון של

$\min_{\bar{x}} \left( \max_{\bar{\lambda}} L(\bar{x}, \bar{\lambda}) \right)$ . למעשה הבעיה שלנו כעת שקולה למציאת הפונקציה הדואלית, ובשביל זה אנחנו צריכים למזער את

הלגרנז'יאן לפי  $\bar{x}$ . נסמן כי הפתרון האופטימלי למיזעור הלגרנז'יאן מתקבל בנקודה  $\bar{x}_{\bar{\lambda}}$  ולכן בנקודה זו מתקיים תנאי הכרחי מסדר ראשון, KKT (הנגזרת/גרדיאנט הראשונה מתאפסת):

$$\nabla_{\bar{x}} L(\bar{x}_{\bar{\lambda}}, \bar{\lambda}) = \bar{x}_{\bar{\lambda}} + A^T \bar{\lambda} = 0$$

מהשיויון האחרון נקבל:

$$\bar{x}_{\bar{\lambda}} = -A^T \bar{\lambda}$$

כאשר נציב את  $\bar{x}_{\bar{\lambda}}$  בפונקצית המטרה (כאשר נזכיר כי המטרה הראשונה היא למזער את הלגרנז'יאן), כלומר בלגרנז'יאן ואז נישאר עם פונקציה שתלויה רק במשתנה  $\bar{\lambda}$  כפי שרצינו. נקבל:

$$\eta(\bar{\lambda}) \triangleq \min_{\bar{x}} L(\bar{x}, \bar{\lambda}) = L(\bar{x}_{\bar{\lambda}}, \bar{\lambda}) = \frac{1}{2} \bar{\lambda}^T A A^T \bar{\lambda} - \bar{\lambda}^T A A^T \bar{\lambda} - \bar{\lambda}^T \bar{b} = -\frac{1}{2} \bar{\lambda}^T A A^T \bar{\lambda} - \bar{\lambda}^T \bar{b}$$

ולכן הבעיה הדואלית היא למקסם:

$$\max_{\bar{\lambda} \geq 0} \eta(\bar{\lambda}) = \max_{\bar{\lambda} \geq 0} \left( -\frac{1}{2} \bar{\lambda}^T A A^T \bar{\lambda} - \bar{\lambda}^T \bar{b} \right)$$

(כאשר  $\bar{\lambda} \geq 0$  משמעו:  $\forall i: \lambda_i \geq 0$ )

נגזור לפי  $\bar{\lambda}$  ונדרוש שהנגזרת תתאפס (זו למעשה פונקציה ריבועית שאנחנו מחפשים את המקסימום שלה ע"י התאפסות הנגזרת/גרדיאנט). נקבל:

$$\nabla_{\bar{\lambda}} \left( -\frac{1}{2} \bar{\lambda}^T A A^T \bar{\lambda} - \bar{\lambda}^T \bar{b} \right) = -\frac{1}{2} \left( A A^T + (A A^T)^T \right) \bar{\lambda} - \bar{b} = 0$$

ולכן קיבלנו כי נקודת המינימום מתקבלת עבור:

$$\bar{\lambda} = -2 \left( A A^T + (A A^T)^T \right)^{-1} \bar{b} = -2 \left( A A^T + A A^T \right)^{-1} \bar{b} = - \left( A A^T \right)^{-1} \bar{b}$$

ולכן קיבלנו כי המינימום (של בעית האופטימיזציה המקורית) מתקבל בנקודה:

$$\bar{x}_{\bar{\lambda}} = -A^T \bar{\lambda} = (-A^T) \cdot \left( - \left( A A^T \right)^{-1} \bar{b} \right) = A^T \left( A A^T \right)^{-1} \bar{b}$$

אם מטריצה  $A$  היא הפיכה אז ניתן להמשיך ולקבל:

$$\bar{x}_{\bar{\lambda}} = A^T \left( A A^T \right)^{-1} \bar{b} = A^T \left( A^T \right)^{-1} A^{-1} = A^{-1} \bar{b}$$

בעיה מספר 2: בעיה לינארית, תכנות לינארי

נתונה בעית האופטימיזציה הבאה:

$$\min_{\bar{x} \in \mathbb{R}^n} f(\bar{x}) = \bar{c}^T \bar{x}$$

$$\text{Subject To (S.T.):} \quad - \left( A \bar{x} - \bar{b} \right) \leq 0$$

הלגרנז'יאן עבור בעיה זו הוא:

$$L(\bar{x}, \bar{\lambda}) = \bar{c}^T \bar{x} - \bar{\lambda}^T (A\bar{x} - \bar{b})$$

(נזכר כי הלגרנז'יאן נתון ע"י:  $L(\bar{x}, \bar{\lambda}) \triangleq f(\bar{x}) + \sum_{i \in I} \lambda_i g_i(\bar{x}^*)$  כאשר  $g_i(\bar{x})$  הם האילוצים)

נרצה למצוא את הפונקציה הדואלית, ובשביל זה אנחנו צריכים למזער את הלגרנז'יאן לפי  $\bar{x}$ . נסמן כי הפתרון האופטימלי למיזעור הלגרנז'יאן מתקבל בנקודה  $\bar{x}_{\bar{\lambda}}$  ולכן בנקודה זו מתקיים תנאי הכרחי מסדר ראשון, KKT, (הנגזרת/גרדיאנט הראשונה מתאפסת):

$$\nabla_{\bar{x}} L(\bar{x}, \bar{\lambda}) = \bar{c} - A^T \bar{\lambda} = 0$$

נבחין כי לפי תנאי זה אנחנו רואים כי קיים תחום שלם (שלא תלוי ב- $\bar{x}$ ) שבו מתקבלת נקודת המינימום של הפונקציה  $f(\bar{x})$ . מהשיוויון האחרון, אם קיים מינימום ללגרנז'יאן אז ניתן לדרוש:

$$\bar{c} - A^T \bar{\lambda} = 0$$

ולכן הבעיה הדואלית היא למקסם:

$$\eta(\bar{\lambda}) \triangleq \min_{\bar{x}} L(\bar{x}, \bar{\lambda}) = \min_{\bar{x}} \left( (\bar{c} - A^T \bar{\lambda})^T \bar{x} + \bar{\lambda}^T \bar{b} \right) = \begin{cases} \bar{\lambda}^T \bar{b} & , A^T \bar{\lambda} = \bar{c} \\ -\infty & , otherwise \end{cases}$$

או בצורה אחרת:

$$\max_{\bar{\lambda} \geq 0} \eta(\bar{\lambda}) = \max_{\bar{\lambda} \geq 0} \bar{\lambda}^T \bar{b} = \max_{\bar{\lambda} \geq 0} \bar{b}^T \bar{\lambda}$$

*Subject To (S.T):*  $A^T \bar{\lambda} = \bar{c}$

(כאשר  $\bar{\lambda} \geq 0$  משמעו  $\forall i: \lambda_i \geq 0$ )

כאשר את הבעיה הזו אפשר לפתור באמצעות פונקצית עונשין. כיוון שזו בעית אופטימיזציה שצריך למצוא בה את המקסימום (ולא את המינימום כפי שאנחנו רגילים) אז פונקציית העונשין תהיה שווה למינוס אינסוף כאשר אנחנו לא נמצאים בתחום המותר לפי האילוץ ( $A^T \bar{\lambda} = \bar{c}$ ) וערכה יהיה אפס כאשר אנחנו נמצאים בתחום המותר.

בזאת סיימנו את החלק העיקרי השני בקורס: אופטימיזציה של פונקציות כלליות לא לינאריות. כעת מה שנשאר זה ללמוד שיטות אופטימיזציה מודרניות שהתפתחו במהלך שנות ה-90 ומתפתחות גם בשנות האלפיים. שיטות אלו ישתמשו בתכנות קוני (מהמילה: קונוס).

## הרצאה 16 – Conic Programming, תכנות קוני

נושא התכנות הקוני הוא הנושא המסכם של הקורס והוא מתייחס לשיטות אופטימיזציה מודרניות. תכנות קוני זו למעשה הכללה של תכנות לינארי. נניח ויש לנו בעיית תכנות לינארי (LP (Linear Programming):

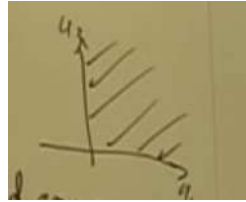
$$LP: \min_{\bar{x} \in \mathbb{R}^m} \langle \bar{c}, \bar{x} \rangle = \min_{\bar{x} \in \mathbb{R}^m} \bar{c}^T \bar{x}, \quad \bar{c} \in \mathbb{R}^m$$

$$S.T (Subject To): A\bar{x} - \bar{b} \geq \bar{0}, \quad \bar{b} \in \mathbb{R}^n, A_{[n \times m]}$$

ניתן להכליל את ניסוח הבעיה הזו למחלקה רחבה יותר של בעיות שניתנות לתכנות. מה שלמדנו עד עכשיו זה תכנות לא לינארי ואת זה עשינו כאשר החלפנו את פונקציית המטרה מפונקציה לינארית לפונקציה אחרת (פונקציה עונשין מצרפית למשל) כלשהי (לרוב דרשנו שהפונקציה תהיה קמורה וגזירה ברציפות). התכנות הקוני אומר שנשאיר את ניסוח הבעיה כמו שהוא, רק שנבצע הכללה עבור האילוף בלבד (שהוא מסוג Inequality):

בתכנות לינארי, באי-השוויון לעיל, הווקטור  $A\bar{x} - \bar{b}$  למעשה דורשים ממנו שיקיים:  $A\bar{x} - \bar{b} \in \mathbb{R}_n^+$  (כלומר כל רכיב של הווקטור הוא חיובי), כי מתקיים לפי ההגדרה:

$$\mathbb{R}_n^+ = \{\bar{u} : u_i \geq 0, \forall 1 \leq i \leq n\}$$



ולמעשה  $\mathbb{R}_n^+$  הוא קבוצה קונית אך אפשר גם להשתמש בקונוסים אחרים. נכליל את התכנות הלינארי, כלומר את הקבוצה הקונית  $\mathbb{R}_n^+$  ונגדיר:

K – Convex pointed cone



(הקבוצה הקונית נתחמת ע"י שני הווקטורים ה"חיצוניים" ומכילה את כל הווקטורים ביניהם)

הגדרה:

קבוצה קונית מקיימת את 3 התכונות הבאות. לכל שני איברים בקונוס,  $\bar{x}, \bar{y} \in K$ , מתקיימות:

- $\lambda \bar{x} \in K, \quad 0 \leq \lambda \in \mathbb{R}$
- $-\bar{x} \notin K, \quad \bar{x} \neq \bar{0}$
- $\bar{x} + \bar{y} \in K$  (חיבור ווקטורי כפי שלמדנו עד כה באלגברה)



נבחין כי בקבוצה הקונית  $\mathbb{R}_n^+$  מתקיימות הטענות הבאות:

$$\bar{a} \geq 0 \Leftrightarrow \bar{a} \in \mathbb{R}_n^+$$

$$\bar{a} \geq \bar{b} \Leftrightarrow \bar{a} - \bar{b} \in \mathbb{R}_n^+$$

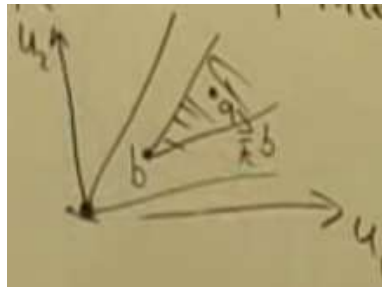
ונוכל להכליל טענות אלו בקבוצה הקונית  $K$  באופן הבא:

$$\bar{a} \geq_K 0 \Leftrightarrow \bar{a} \in K$$

$$\bar{a} \geq_K \bar{b} \Leftrightarrow \bar{a} - \bar{b} \in K$$

(כאשר האי-שוויון  $\geq_K, \leq_K$  הוא סימון לאי-שוויון בקבוצה הקונית  $K$ )

בצורה גרפית, מתקיים  $\bar{a} \geq_K \bar{b}$  אם "מ" הווקטור  $\bar{a}$  נמצא בקונוס ה-  $K$  שמוזן כך שקודקודו הוא בווקטור  $\bar{b}$  ולא בראשית הצירים):



דוגמאות לקבוצות קוניות:

דוגמה 1:

הקבוצה  $\mathbb{R}_n^+ = \{\bar{u} : u_i \geq 0, \forall 1 \leq i \leq n\}$ .

זוהי למעשה קבוצת הווקטורים האי-שלילים (ולכן כוללת גם את ווקטור האפס) ממימד  $n$ . ווקטור אי-שלילי מוגדר להיות ווקטור שכל רכיביו הם אי-שליליים.

דוגמה 2:

הקבוצה  $L^n$ : הקבוצה הקונית של לורנץ (Lorentz Cone) וקבוצה זו גם נקראת Ice cream cone כיוון שהיא נראית כמו גביע גלידה. למשל בתלת מימד:



בקבוצה  $L^n$  מתקיים:

$$u_n \geq \|(u_1, u_2, \dots, u_{n-1})\|_2$$

כלומר הרכיב האחרון בווקטור הוא בגודל שגדול מהנורמה של כל שאר הרכיבים בווקטור.

דוגמה 3:

הקבוצה  $S_+^m$  קבוצת המטריצות האי-שליליות (Positive semi-definite) הסימטריות מסדר  $[m \times m]$ . נוכיח את 3 התכונות הדרושות עבור קבוצה זו בכדי שתהיה קבוצה קונית:

- אם נכפיל מטריצה אי-שלילית סימטרית בסקלר חיובי התוצאה עדיין תהיה מטריצה אי-שלילית סימטרית ולכן בקבוצה.
- אם נכפיל מטריצה אי-שלילית סימטרית בסקלר שלילי נקבל כי המטריצה שלילית ולכן לא נמצאת בקבוצה.
- אם ניקח שתי מטריצות אי-שליליות סימטריות ונחבר אותן עדיין נקבל מטריצה אי-שלילית סימטרית ולכן בקבוצה.

הוכחנו את 3 התכונות ולכן קבוצה זו היא קבוצה קונית.

בעיות שנפתרות ע"י תכנות קוני:

בעיות שנפתרות ע"י תכנות קוני הן באופן כללי בעיות מהצורה הבאה:

$$LP: \min_{\bar{x} \in \mathbb{R}^m} \langle \bar{c}, \bar{x} \rangle = \min_{\bar{x} \in \mathbb{R}^m} \bar{c}^T \bar{x}, \quad \bar{c} \in \mathbb{R}^m$$

$$S.T (Subject To): \quad A\bar{x} - \bar{b} \geq_K \bar{0}, \quad A_{[n \times m]}, \quad \bar{b} \in \mathbb{R}^n, \quad A\bar{x} - \bar{b} \in K = \mathbb{R}_n^+$$

דוגמאות לבעיות שנפתרות ע"י תכנות קוני:

דוגמה 1:

בעיה סטנדרטית של תכנות קוני: (תכנות לינארי)

$$LP: \min_{\bar{x} \in \mathbb{R}^m} \langle \bar{c}, \bar{x} \rangle = \min_{\bar{x} \in \mathbb{R}^m} \bar{c}^T \bar{x}, \quad \bar{c} \in \mathbb{R}^m$$

$$S.T (Subject To): \quad A\bar{x} - \bar{b} \geq_K \bar{0}, \quad A_{[n \times m]}, \quad \bar{b} \in \mathbb{R}^n, \quad A\bar{x} - \bar{b} \in K = \mathbb{R}_n^+$$

דוגמה 2:

בעיה תכנות קוני ריבועית: (תכנות ריבועי)

$$LP: \min_{\bar{x} \in \mathbb{R}^m} \langle \bar{c}, \bar{x} \rangle = \min_{\bar{x} \in \mathbb{R}^m} \bar{c}^T \bar{x}, \quad \bar{c} \in \mathbb{R}^m$$

$$S.T (Subject To): \begin{cases} \begin{bmatrix} A\bar{x} - \bar{b} \\ \bar{d}^T \bar{x} - e \end{bmatrix} \succeq_{L^{n+1}} \bar{0} & , \bar{b} \in \mathbb{R}^n, A_{[n \times m]}, \bar{d} \in \mathbb{R}^m, e \in \mathbb{R}, \\ \|A\bar{x} - \bar{b}\|_2 \leq \bar{d}^T \bar{x} - e \end{cases}$$

דוגמה 3: (דוגמה חשובה!)

בעיות תכנות קוני, אשר הקבוצה הקונית היא  $S_+^m$ . תכנות זה נקרא: Semi Definite Programming או בקיצור SDP. נניח כי בידנו  $\bar{x} \in \mathbb{R}^n$  וקבוצת מטריצות אי-שליליות סימטריות:  $\{B, A_1, A_2, \dots, A_n\} \in S_+^m$ . נגדיר אופרטור שמופעל על ווקטור  $\bar{x} \in \mathbb{R}^n$  באופן הבא:

$$A(\bar{x}) = \left( \sum_{i=1}^n x_i A_i \right) - B \in S^m$$

(נבחין כי תמונת האופרטור זו מטריצה סימטרית כיוון שכל צירוף לינארי של מטריצות סימטריות הוא מטריצה סימטרית)

נתייחס ל- $A(\bar{x})$  כאל העתקה ממרחב  $\mathbb{R}^n$  למרחב המטריצות הסימטריות, כלומר  $A: \mathbb{R}^n \rightarrow S^m$ . בנוסף נסמן:

$$A(\bar{x}) \triangleq \sum_{i=1}^n x_i A_i, \quad \bar{x} \in \mathbb{R}^n$$

(הערה: נבחין כי האות  $A$  שונה מהאות  $A$  וכן  $A(\bar{x})$  זה גם אופרטור העתקה  $A: \mathbb{R}^n \rightarrow S^m$ )

ולכן קיבלנו:

$$A(\bar{x}) = A(\bar{x}) - B$$

וכעת הבעיה (תכנות SDP) היא:

$$LP: \min_{\bar{x} \in \mathbb{R}^n} \bar{c}^T \bar{x}, \quad \bar{c} \in \mathbb{R}^n$$

$$S.T: A(\bar{x}) - B = \left( \sum_{i=1}^n x_i A_i \right) - B \succeq_K 0, \quad B \in S_+^m, A(\bar{x}): \mathbb{R}^n \rightarrow S_+^m, A: \mathbb{R}^n \rightarrow S^m$$

$$, \quad K = S_+^m, \{A_i\}_{i=1}^n \in S_+^m$$

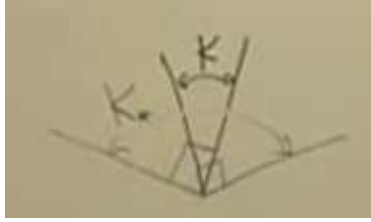
קבוצה קונית דואלית, קונוס דואלי (Dual Cone):

הגדרה:

קבוצה קונית דואלית  $K_*$  לקבוצה הקונית  $K$  היא:

$$K_* = \{ \bar{y} : \langle \bar{x}, \bar{y} \rangle \geq 0, \quad \forall \bar{x} \in K \}$$

ובצורה גרפית:



(אם הקונוס שלנו הוא הקונוס  $K = \mathbb{R}^2$ , ניתן לחשוב על הקבוצה הדואלית כעל קבוצה שמכילה את כל הווקטורים שהמכפלה הפנימית שלהם עם כל ווקטור בקבוצה הקונית היא גדולה או שווה לאפס, אך המכפלה הפנימית הסטנדרטית היא למעשה כפי שהכרנו אותה עוד בפיסיקה 1:  $\langle \bar{x}, \bar{y} \rangle = |\bar{x}| |\bar{y}| \cos(\theta)$  (כאשר  $\theta$  זו הזווית בין שני הווקטורים) ולכן הקבוצה הקונית הדואלית מכילה את כל הווקטורים שהם יוצרים זווית שגדולה מ- $90^\circ$  עם כל ווקטור בקבוצה הקונית)

הערה: הקבוצה הקונית הדואלית של הקבוצה הקונית  $K_*$  היא הקבוצה הקונית  $K$ , כלומר:

$$(K_*)_* = K$$

הגדרה:

קבוצה קונית  $K$  שדואלית לעצמה (Self-Dual), מקיימת:  $K = K_*$ . בצורה גרפית:



דוגמאות לקבוצות קוניות שדואליות לעצמן (תרגיל הוכחה לשיעורי הבית):

- $\mathbb{R}_m^+$  - קבוצת הווקטורים שכל רכיביהם חיוביים, והווקטורים כולם ממימד  $m$ .
- $L^m$  - הקבוצה הקונית של לורנץ (שבה הרכיב האחרון בווקטור גדול מהנורמה של הווקטור שמורכב משאר הרכיבים)
- $S_m^+$  - קבוצת המטריצות הסימטריות שהן גם אי-שליליות (Positive semi-definite) מסדר  $m \times m$ .

הוכחה כי הקבוצה הקונית  $S_m^+$  היא דואלית לעצמה:

נסמן  $K = S_m^+$ . אנחנו צריכים להוכיח את שתי הטענות הבאות:

- (1)  $\text{if } X \in K \Rightarrow X \in K_*$
- (2)  $\text{if } Y \notin K \Rightarrow \exists X \in K : \langle X, Y \rangle < 0$

(ניזכר כי הקבוצה הקונית הדואלית היא:  $K_* = \{Y : \langle X, Y \rangle \geq 0, \forall X \in K\}$ )

נוכיח את טענה (1). נניח כי בדינו  $X \in K = S_+^m$  ואנחנו נוכיח (לפי ההגדרה של הקבוצה הקונית) כי לכל  $Y \in S_m^+$  שנבחר נקבל  $\langle X, Y \rangle \geq 0$  ולכן למעשה לפי ההגדרה נוכיח כי  $X \in K_*$ . נניח כי אכן  $X \in K_*$ . עבור כל מטריצה  $Y \in S_m^+$  המכפלה הפנימית ביניהן מקיימת:

$$\langle X, Y \rangle = \text{Trace}(X^T Y) \underset{\substack{X \text{ symmetric} \\ \Rightarrow X = X^T}}{=} \text{Trace}(XY)$$

בנוסף, כל מטריצה סימטרית יכולה להיות מיוצגת ע"י מכפלה בין מטריצה משולשית למטריצה המשוכללת שלה:

$$X = U^T U, \quad Y = V^T V$$

ולכן:

$$\langle X, Y \rangle = \text{Trace}(XY) = \text{Trace}(U^T U V^T V) \underset{\substack{\text{Circular shift} \\ \text{under Trace}}}{=} \text{Trace}\left(\underset{=A^T}{V U^T} \underset{=A}{U V^T}\right) = \text{Trace}(A^T A) = \sum_{i=1}^n (a_{ii})^2 \geq 0$$

עבור הטענה (2), אם  $Y \notin K$  אז קיים ווקטור  $\bar{u} \in \mathbb{R}^m$  כך שמתקיים  $\bar{u}^T Y \bar{u} < 0$  (שזה למעשה אומר כי המטריצה היא שלילית ולכן  $Y \notin K = S_+^m$  ואז:

$$0 > \bar{u}^T Y \bar{u} \underset{\substack{\bar{u}^T Y \bar{u} \in \mathbb{R} \\ a \in \mathbb{R} \Rightarrow a = \text{Trace}(a)}}{=} \text{Trace}(\bar{u}^T Y \bar{u}) \underset{\substack{\text{Circular shift} \\ \text{under Trace}}}{=} \text{Trace}\left(\underset{\triangleq X^T \text{ (matrix)}}{\bar{u} \bar{u}^T} Y\right) = \langle X, Y \rangle$$

ואם נראה כי  $X \in S_m^+$  אז מתקיים למעשה  $\langle X, Y \rangle < 0$  ולכן הוכחנו את הטענה. נוכיח כי  $\bar{u} \bar{u}^T = X \in S_m^+$ . נבדוק את אי-שליליות המטריצה  $X = \bar{u} \bar{u}^T$  לפי ההגדרה:

$$\forall \bar{v} \in \mathbb{R}^n: \quad \bar{v}^T X \bar{v} = \bar{v}^T \bar{u} \bar{u}^T \bar{v} = a \cdot a^T = a^2 \geq 0$$

$$\underset{\substack{a \triangleq \bar{v}^T \bar{u} \in \mathbb{R} \\ \Rightarrow a^T = a}}{a \triangleq \bar{v}^T \bar{u} \in \mathbb{R}}$$

כעת נותר רק להוכיח כי המטריצה  $X$  היא סימטרית אך את זה קל לראות לפי ההגדרה של מכפלת ווקטור בווקטור המשוכלל (זה מקרה פרטי של Outer Product).

הוכחה כי הקבוצה הקונית  $\mathbb{R}_m^+$  היא דואלית לעצמה (קבוצת הווקטורים ממימד  $m$ ):

נרצה להוכיח שוב את שתי הטענות הבאות:

$$\begin{aligned} (1) \quad & \text{if } \bar{x} \in K \quad \Rightarrow \quad \bar{x} \in K_* \\ (2) \quad & \text{if } \bar{y} \notin K \quad \Rightarrow \quad \exists \bar{x} \in K: \langle \bar{x}, \bar{y} \rangle < 0 \end{aligned}$$

(נזכר כי הקבוצה הקונית, לפי ההגדרה זו הקבוצה:  $\{ \bar{y} : \langle \bar{x}, \bar{y} \rangle \geq 0, \forall \bar{x} \in K \}$   $K_*$ )

נתחיל מטענה (2), נוכיח כי אם קיים  $\bar{y} \notin K = \mathbb{R}_m^+$  אז קיים  $\bar{x} \in K = \mathbb{R}_m^+$  כך שמתקיים:

$$\langle \bar{x}, \bar{y} \rangle < 0 \quad \Leftrightarrow \quad \bar{x}^T \bar{y} < 0 \quad \Leftrightarrow \quad \sum_{i=1}^m x_i y_i < 0$$

אך אם  $\bar{y} \notin \mathbb{R}_m^+$ , זה אומר כי קיים (לפחות) אינדקס אחד  $1 \leq i \leq m$  כך שמתקיים  $y_i < 0$ . אנחנו נבחר  $\bar{x} \in \mathbb{R}_m^+$  באופן הבא:

$$\bar{x} = \{x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_m\} = \{0, \dots, 0, 1, 0, \dots, 0\}$$

ואז נקבל כי אכן מתקיים  $\langle \bar{x}, \bar{y} \rangle < 0$  כפי שרצינו להראות.

עבור טענה (1), נניח כי  $\bar{y} \in K = \mathbb{R}_m^+$  ולכן  $y_i \geq 0$  לכל  $1 \leq i \leq m$  וכן אם  $\bar{x} \in \mathbb{R}_m^+$  אז  $x_i \geq 0$  לכל  $1 \leq i \leq m$  אז מתקיים:

$$\langle \bar{x}, \bar{y} \rangle = \sum_{i=1}^m x_i y_i \geq 0$$

דואליות בתכנות קוני מסוג SDP:

כפי שראינו בתכנות לא לינארי (כאשר למדנו על הלגרנז'יאן התעסקנו עם פונקצית מטרה שהיא הייתה פונקציית עונשין מצרפית והיא הייתה לא לינארית כיוון שפונקציות העונשין היו לא לינאריות), לכל בעיה הצלחנו למצוא את הבעיה הדואלית שלה וגם במקרה זה נראה כיצד למצוא את הבעיה הדואלית של כל בעיה בתכנות קוני (הקונוס בבעיה הכללית שנציג הוא קונוס המטריצות הסימטריות האי-שליליות, כלומר  $S_+^m$  וזה כאמור קונוס שגם הוכחנו שהוא דואלי לעצמו וכן קראנו לתכנות זה תכנות SDP). הבעיה הכללית בתכנות קוני SDP, מנוסחת באופן הבא:

$$\begin{aligned} LP: \quad \min_{\bar{x} \in \mathbb{R}^n} \langle \bar{c}, \bar{x} \rangle &= \min_{\bar{x} \in \mathbb{R}^n} \bar{c}^T \bar{x}, & \bar{c} &\in \mathbb{R}^n \\ \\ ST (Subject To): \quad A(\bar{x}) - B &= \left( \sum_{i=1}^n x_i A_i \right) - B \succeq_K 0, & B &\in S^m \\ &A(\bar{x}) - B \in K = S_+^m, \\ &A(\bar{x}) \triangleq \left( \sum_{i=1}^n x_i A_i \right) \in S^m \\ &\forall 1 \leq i \leq n: A_i \in S^m \end{aligned}$$

(נבחין כי הקבוצה  $S^m$  זו קבוצת המטריצות הסימטריות מסדר  $[m \times m]$  ואילו  $S_+^m$  זו קבוצת המטריצות הסימטריות האי-שליליות מסדר  $[m \times m]$  ונבחין כי רק מטריצת האילון  $A(\bar{x})$  צריכה להיות סימטרית חיובית אך שאר המטריצות צריכות להיות רק סימטריות) תחילה, נבחין כי איבר בקבוצה הקונית הדואלית  $Y \in K_* = S_+^m$  (הוכחנו כי הקבוצה  $S_+^m$  דואלית לעצמה) מקיים את התכונות/עובדות הבאות:

- אם  $\bar{x} \in Feasible Set$  אז  $\langle A(\bar{x}) - B, Y \rangle \geq 0$ . הוכחה: נבחין כי אם  $\bar{x} \in Feasible Set$  כלומר,
- $\langle A(\bar{x}) - B, Y \rangle \geq 0$  אז לפי ההגדרה של קבוצה קונית דואלית, מתקיים  $A(\bar{x}) - B \succeq_K 0$ .
- $\langle A(\bar{x}), Y \rangle \geq_K \langle B, Y \rangle$ . הוכחה: ע"י העברת אגף נקבל:

$$0 \leq_K \langle A(\bar{x}), Y \rangle - \langle B, Y \rangle = A(\bar{x})^T Y - B^T Y = (A(\bar{x}) - B)^T Y = \langle A(\bar{x}) - B, Y \rangle$$

•  $\langle \bar{x}, A^*(Y) \rangle \geq \langle \bar{b}, A^*(Y) \rangle$ , כאשר  $A^*$  מוגדר להיות אופרטור שמעביר מטריצות לוקטורים (נגדיר בהמשך).

המטרה שלנו זה למצוא איבר  $Y$  (כאשר מתקיים  $A^*(Y) = \bar{y}$ ) כך שמתקיים  $A^*(Y) = \bar{c}$  ואז נוכל לומר כי לפי התכונה השלישית לעיל יתקיים:

$$\langle \bar{x}, \bar{c} \rangle \geq \langle \bar{b}, \bar{y} \rangle$$

ובכך למעשה נמצא ערך גבול תחתון עבור הבעיה שאנחנו מנסים לפתור שכן מתקיים:  $\langle \bar{x}, \bar{c} \rangle = \bar{c}^T \bar{x}$ . בנוסף, נרצה שערך תחתון זה יהיה כמו שיותר "צמוד" (מלמטה) לערך האמיתי של הפתרון, כלומר אנחנו רוצים למקסם את הביטוי  $\langle \bar{b}, \bar{y} \rangle$ , כלומר לפתור:

$$\begin{aligned} \max_{A(\bar{y}) \in K_*} \langle \bar{b}, \bar{y} \rangle &= \bar{b}^T \bar{y}, & A(\bar{b}) \in K_* = S_+^m \\ S.T \text{ (Subject To): } & A^*(Y) = \bar{c}, & A(\bar{c}) \in K_* = S_+^m \end{aligned}$$

וזו למעשה הבעיה הדואלית של תכנות קוני SDP. זו נקראת הבעיה הדואלית, כי במקום למצוא את המינימום של פונקצית המטרה שנדרשנו, אנחנו מנסים למצוא את המקסימום של פונקצית מטרה אחרת.

משפט הדואליות החלשה:

עבור הבעיה:

$$\begin{aligned} LP: \min_{\bar{x} \in \mathbb{R}^n} \langle \bar{c}, \bar{x} \rangle, & \bar{c} \in \mathbb{R}^n \\ S.T \text{ (Subject To): } & A(\bar{x}) = A(\bar{x}) - B \succeq_K 0, K = S_+^m \end{aligned}$$

ואם נניח כי  $Y \in K_*$  וכן  $A^*(Y) = \bar{c}$  אז מתקיים:

$$\langle \bar{x}, \bar{c} \rangle \geq \langle \bar{b}, \bar{y} \rangle$$

כאשר  $\bar{x} \in Feasible Set$  וכן  $\bar{y} \in Dual Feasible Set$ . (כאשר נזכור כי הוכחנו כי הקבוצה  $\mathbb{R}_n^+$  דואלית לעצמה)

משפט הדואליות החזקה:

כי שראינו, בעית התכנות הקוני והבעיה הדואלית לה הן:

$$\begin{aligned} \text{(Primary Problem = P):} & \begin{cases} LP: \min_{\bar{x} \in \mathbb{R}^n} \langle \bar{c}, \bar{x} \rangle, & \bar{c} \in \mathbb{R}^n \\ S.T \text{ (Subject To): } & A(\bar{x}) - B \succeq_K 0, K = S_+^m \end{cases} \\ \text{(Dual Problem = D):} & \begin{cases} \max_{A(\bar{y}) \in K_*} \langle \bar{b}, \bar{y} \rangle, & \bar{b} \in \mathbb{R}^n, A(\bar{y}) \in K_* = S_+^m \\ S.T \text{ (Subject To): } & A^*(Y) = \bar{c}, \bar{c} \in \mathbb{R}^n \end{cases} \end{aligned}$$

אם בבעיה  $(P)$ , מצאנו פתרון  $\bar{x}^* \in \mathbb{R}^n$  אשר מקיים  $A(\bar{x}^*) - B \succ_K 0$  כלומר נמצא בקבוצה הקונית ממש, כלומר  $\bar{x}^* \in \text{Strictly Feasible Set}$  אז לבעיה הדואלית  $(D)$  יש פיתרון והפתרון האופטימלי מקיים שהפער הדואלי (duality gap) הוא אפס:

$$\text{duality gap} = \langle \bar{c}, \bar{x}^* \rangle - \langle \bar{b}, \bar{y}^* \rangle = 0$$

או אם בבעיה  $(D)$ , מצאנו פתרון  $\bar{y}^* \in \mathbb{R}^n$  אשר  $A^*(Y)$  נמצא בקבוצה הקונית ממש, כלומר  $\bar{y}^* \in \text{Strictly Feasible Set}$  אז לבעיה  $(P)$  יש פיתרון והפתרון האופטימלי מקיים גם הוא את השוויון לעיל.

הערה:

במשפט הדואליות בתכנות לינארי, לא דרשנו את התנאי של Strict Feasibility. משפט הדואליות החזקה האחרון (לעיל) מזכיר במידה מסויימת את Slater Condition עבור דואליות חזקה בתכנות לא לינארי. התנאי בדבר Strict Feasibility הוא תנאי יחסית חזק שצריך לקיים בשביל הדואליות החזקה אז לרוב תנאי זה מתקיים.

הערה: Complimentary Slackness

כמו בתכנות הלינארי שדרשנו שכופלי לגרנז' יקיימו:  $(\bar{\lambda}^*)^T \bar{g}(\bar{x}^*) = 0$  אז במקרה של תכנות קוני מתקיים:

$$\langle Y^*, A(\bar{x}^*) - B \rangle = 0$$

או באופן שקול:

$$\langle Y^*, A(\bar{x}^*) - B \rangle - \langle Y^*, B \rangle = 0$$

אך זה ניתן גם לרישום אחר (באמצעות אופרטור האדג'וינט  $(A^*(\cdot))$ ):

$$\langle A^*(Y^*), \bar{x}^* \rangle - \langle \bar{y}^*, \bar{b} \rangle = 0$$

ולאחר שימוש באילוץ  $A^*(Y) = \bar{c}$  נקבל:

$$\langle \bar{c}, \bar{x}^* \rangle - \langle \bar{y}^*, \bar{b} \rangle = 0$$

שזהו בדיוק הפער הדואלי (duality gap).

דוגמה: Semi Definite Programming (SDP)

תזכורת: נניח  $\bar{x} \in \mathbb{R}^n$ , וכן קבוצת מטריצות:  $\{B, A_1, A_2, \dots, A_n\} \in S^m$ . האופרטור הוא:

$$A(\bar{x}) = \left( \sum_{i=1}^n x_i A_i \right) - B \triangleq A(\bar{x}) - B$$

ובעית האופטימיזציה היא:



$$(P): \begin{cases} LP: \min_{\bar{x} \in \mathbb{R}^n} \bar{c}^T \bar{x} & , \bar{c} \in \mathbb{R}^n \\ S.T (Subject To): A(\bar{x}) = A(\bar{x}) - B \succeq_K 0 & , K = S_+^m \end{cases}$$

והבעיה הדואלית:

$$(D): \begin{cases} \max_{Y \succeq_{K_*} 0} \langle B, Y \rangle \equiv Trace(B^T Y) & , B \in K_* = S_+^m \\ S.T (Subject To): A^*(Y) = \bar{c} & , \bar{c} \in \mathbb{R}^n \end{cases}$$

(כאשר נבחין כי  $S_+^m$  היא קבוצה דואלית לעצמה ולכן אם  $Y \in K$  אז גם  $Y \in K_*$ )

נמצא איך מחשבים את האילוץ  $A^*(Y)$  כלומר את הגדרת אופרטור האדג'וינט. נבחין כי מתקיים:

$$\langle \bar{x}, A^*(Y) \rangle = \langle A(\bar{x}), Y \rangle \equiv \left\langle \left( \sum_{i=1}^n x_i A_i \right), Y \right\rangle = \sum_{i=1}^n (x_i \langle A_i, Y \rangle) = \left\langle \bar{x}, \begin{pmatrix} \langle A_1, Y \rangle \\ \vdots \\ \langle A_n, Y \rangle \end{pmatrix} \right\rangle$$

כלומר מתקיים:

$$\langle \bar{x}, A^*(Y) \rangle = \left\langle \bar{x}, \begin{pmatrix} \langle A_1, Y \rangle \\ \vdots \\ \langle A_n, Y \rangle \end{pmatrix} \right\rangle$$

וזה אומר כי ניתן לחשב את הביטוי  $A^*(Y)$  באופן הבא:

$$A^*(Y) = \begin{pmatrix} \langle A_1, Y \rangle \\ \vdots \\ \langle A_n, Y \rangle \end{pmatrix} = \begin{pmatrix} Trace(A_1^T Y) \\ \vdots \\ Trace(A_n^T Y) \end{pmatrix}$$

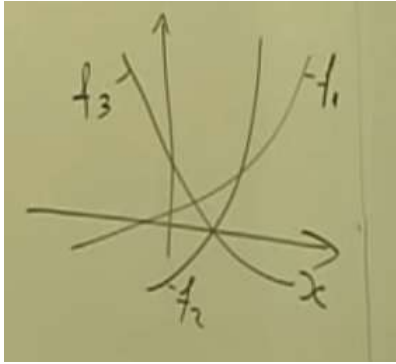
וזו הגדרת אופרטור האדג'וינט (אופרטור אדג'וינט) כאשר נזכור כי  $\{A_1, A_2, \dots, A_n\} \in S_+^m$

כיצד בעית Minmax קשורה לתכנות SDP (Semi Definite Programming):

נראה יישום של תכנות Semi Definite באמצעות דוגמה של בעית Minmax דיסקרטית. נפתור בעיות במימד אחד.

בעית Minmax ראשונה:

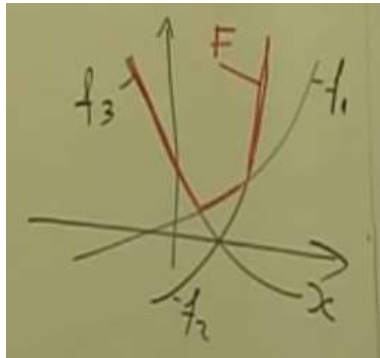
נניח את המצב הבא:



ונגדיר:

$$F(x) \triangleq \max_{i=\{1,\dots,m\}} f_i(x)$$

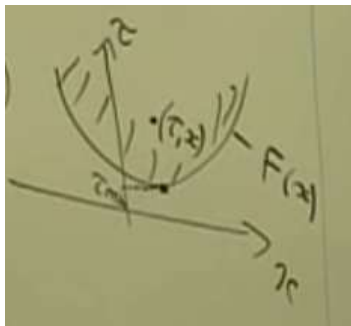
ואם נרצה לשרטט את  $F(x)$  נקבל (באדום):



ובעית האופטימיזציה היא:

$$\min_{x \in \mathbb{R}} F(x)$$

זו נראת כמו בעית אופטימיזציה ללא אילוצים אך נשים לב אליו הוא שהפונקציה  $F(x)$  היא לא גזירה בנקודות "חיבור", ויש לה "שפיצים" בהם הפונקציה לא "חלקה". נביט על הרישום הבא (כאשר למעשה קראנו לציר האנכי בשם  $\tau \equiv F(x)$ ):



נרצה לפתור את בעית האופטימיזציה עם אילוצים:

$$\min_{\tau, x} \tau$$

$$S.T \text{ (Subject To): } (x, \tau) \in \text{epigraph}(F)$$

ובאופן שקול ניתן גם לכתוב:

$$\min_{\tau, x} \tau$$

$$S.T \text{ (Subject To): } \tau \geq F(x)$$

או באופן שקול נוסף:

$$\min_{\tau, x} \tau$$

$$S.T \text{ (Subject To): } \tau \geq f_i(x), \quad i = 1, \dots, m$$

כלומר הפכנו בעית אופטימיזציה ללא אילוצים שפונקציית המטרה שלה לא "חלקה" לבעית אופטימיזציה עם אילוצים אשר כל אילוץ בה הוא "חלק" (אבל עדיין פונקציית המטרה לא "חלקה" – אז למה מה שעשינו הוא עדיף?)

בעית Minmax שניה:

נגדיר פעם נוספת (כמו בבעית Minmax הראשונה לעיל) את הפונקציה  $F(x)$  אך בצורה קצת שונה:

$$F(x) \triangleq \max_{i=\{1, \dots, m\}} |f_i(x)|$$

ונרצה לפתור את בעית האופטימיזציה:

$$\min_{x \in \mathbb{R}} F(x)$$

נרצה לבטא שוב את הבעיה לעיל, ע"י שני משתנים  $x, \tau$ . נוכל באופן דומה למה שראינו בבעיה הראשונה לרשום:

$$\min_{\tau, x} \tau$$

$$S.T \text{ (Subject To): } \tau \geq |f_i(x)|, \quad i = 1, \dots, m$$

אך זה ניסוח לא טוב כיוון שהאילוץ שלנו הוא לא "חלק" כיוון שפונקציית ערך מוחלט לא גזירה בראשית.

נוכל לנסח את הבעיה בצורה קצת שונה אם נבחין כי:

$$|f_i(x)| = \max \{-f_i(x), f_i(x)\}$$

ואז נקבל את בעית האופטימיזציה עם אילוצים הבאה:

$$\min_{\tau, x} \tau$$

$$S.T \text{ (Subject To): } \begin{aligned} \tau &\geq f_i(x) & , & \quad i = 1, \dots, m \\ \tau &\geq -f_i(x) & , & \quad i = 1, \dots, m \end{aligned}$$

או באופן קומפקטי:

$$\min_{\tau, x} \tau$$

$$S.T \text{ (Subject To): } -\tau \leq f_i(x) \leq \tau, \quad i = 1, \dots, m$$

קירוב צ'בישב (Chebyshev Approximation):

נניח ובידינו פונקציה של משתנה יחיד  $h(t)$  ונרצה לקרב את הפונקציה באמצעות קבוצה של פונקציות נתונות.

$$h(t) \approx \sum_i \alpha_i \varphi_i(t)$$

כאשר קבוצת הפונקציות  $\{\varphi_i(t)\}_i$  היא קבוצה כלשהי של פונקציות (ייתכן בסיס פוריה או כל קבוצה אחרת). אנחנו נתייחס כאילו קבוצה הפונקציות  $\{\varphi_i(t)\}_i$  היא סופית. כלומר:

$$h(t) \approx \sum_{i=1}^n \alpha_i \varphi_i(t)$$

ואנחנו נרצה למזער את הביטוי:

$$\left| h(t) - \sum_{i=1}^n \alpha_i \varphi_i(t) \right|$$

נגדיר:  $\psi(t) \triangleq \sum_{i=1}^n \alpha_i \varphi_i(t)$  ולכן אנחנו רוצים למזער את הביטוי:

$$|h(t) - \psi(t)|$$

אך באיזה מובן אנחנו רוצים למזער את הביטוי? אנחנו נרצה שבאינטרבל כלשהו  $T$  נקבל הפרש מקסימלי שהוא המינימלי שניתן לקבל בין  $h(t)$  לבין  $\psi(t)$  כאשר הקירוב שלנו  $\psi(t)$  הוא גם פונקציה של הפרמטרים  $\alpha_i$  שאנחנו בוחרים. באופן פורמלי אנחנו רוצים לפתור את הבעיה:

$$\min_{\alpha \in \mathbb{R}^n} \left( \max_{t \in T} \left| h(t) - \sum_{i=1}^n \alpha_i \varphi_i(t) \right| \right)$$

נבחין כי בבעיה זו, הפרמטרים שעליהם אנחנו עושים מינימיזציה הם רציפים, ואנחנו נעבור כעת למישור בדיד, כלומר המשתנה  $t$  יוכל לקבל ערכים בדידים בלבד:  $t = \{1, 2, \dots, T\}$ . ולכן בעית האופטימיזציה הופכת להיות:

$$\min_{\alpha \in \mathbb{R}^n} \max_{t \in \{1, \dots, T\}} \left| h_t - \sum_{i=1}^n \alpha_i \varphi_{t,i} \right|$$

כאשר  $\varphi$  היא למעשה מטריצה, שהשורה ה- $i$  מכילה את הערכים של הפונקציה  $\varphi_i$  בנקודות הבדידות  $t$ . כעת נבחין כי הבעיה הזו מתקרבת בניסוח שלה לבעיה שראינו בדוגמאות של Minmax. אם נסמן:  $f_t(\bar{\alpha}) \triangleq h_t - \sum_{i=1}^n \alpha_i \varphi_{t,i}$ , אז נבחין כי  $f_t(\bar{\alpha})$  היא פונקציה לינארית ביחס ל- $\bar{\alpha}$ . בדוגמאות של בעית ה-Minmax הגענו בסופו של דבר לניסוח הבעיה:

$$\min_{\tau, x} \tau, \tau \in \mathbb{R}$$

$$S.T \text{ (Subject To): } -\tau \leq f_i(x) \leq \tau, \quad i = 1, \dots, m$$

ולכן גם במקרה זה נוכל לנסח את הבעיה באופן דומה:

$$\min_{\tau, x} \tau, \tau \in \mathbb{R}$$

$$S.T \text{ (Subject To): } -\tau \leq f_t(\bar{\alpha}) \leq \tau, \quad t = 1, \dots, T$$

או באופן מפורש בעית האופטימיזציה שאנחנו צריכים לפתור היא:

$$\min_{\tau, x} \tau, \tau \in \mathbb{R}$$

$$S.T \text{ (Subject To): } -\tau \leq h_t - \sum_{i=1}^n \alpha_i \varphi_{t,i} \leq \tau, \quad t = 1, \dots, T$$

וזו בעית תכנות לינארי (שזה מקרה פרטי של תכנות קוני) וזה ניסוח בעית האופטימיזציה אשר הפתרון שלה הוא למעשה קירוב צ'בישב במישור הבדיד.

נוכל לרשום את כל מה שרשמנו על קירוב צ'בישב, בכתוב ווקטורי באופן הבא:

$$\bar{h} = \begin{pmatrix} h_1 \\ \vdots \\ h_T \end{pmatrix}, \quad \mathcal{G} = \begin{pmatrix} \varphi_1(1) & \cdots & \varphi_1(T) \\ \vdots & \ddots & \vdots \\ \varphi_n(1) & \cdots & \varphi_n(T) \end{pmatrix}, \quad \bar{\tau} = \begin{pmatrix} \tau \\ \vdots \\ \tau \end{pmatrix}$$

ונקבל את הבעיה:

$$\min_{\tau, x} \tau, \tau \in \mathbb{R}$$

$$S.T \text{ (Subject To): } -\bar{\tau} \leq \bar{h} - \mathcal{G}\bar{\alpha} \leq \bar{\tau}$$

אפשר לנסח את הבעיה לעיל בצורה כללית יותר, בצורה של תכנות לינארי כללי:

$$\min_{\bar{x}} \langle \bar{c}, \bar{x} \rangle = \min_{\bar{x}} \bar{c}^T \bar{x}$$

$$S.T \text{ (Subject To): } A\bar{x} \geq \bar{b}$$

כאשר אפשר למצוא את המטריצה  $A$  ואת הווקטורים  $\bar{c}, \bar{b}$  ולקבל בדיוק את בעית האופטימיזציה שניסחנו עבור קירוב צ'בישב.

קירוב צ'בישב מרוכב:

לעיתים מתעורר צורך לקרב פונקציה מרוכבת כלשהי (אשר הארגומנט שלה הוא ממשי) באמצעות קירוב צ'בישב (כלומר ע"י קבוצה סופית של פונקציות כלשהן אך הפעם הן מרוכבות). נתבונן בבעיה הבאה:

$$\min_{\bar{\alpha}} \max_{t=1, \dots, T} \left| h[t] - \sum_{i=1}^n \alpha_i \varphi_{ti} \right|, \quad t \in \mathbb{R}, \quad h[t], \varphi_{ti}, \alpha_i \in \mathbb{C}$$

ונסמן למען הנוחות את הפונקציה המרוכבת הבאה:  $z_t(\bar{\alpha}) \triangleq h[t] - \sum_{i=1}^n \alpha_i \varphi_i$  ולכן הבעיה היא:

$$\min_{\bar{\alpha}} \max_{t=1, \dots, T} |z_t(\bar{\alpha})|$$

כעת, בכדי להימנע מהתעסקות עם מספרים מרוכבים, נתאר כל ערך של הפונקציה המרוכבת  $z_t(\bar{\alpha})$  כווקטור בעל שתי קואורדינטות, רכיב ממשי ומדומה. נגדיר:

$$\bar{u}_t(\bar{\alpha}) = \begin{pmatrix} \text{Re}[z_t(\bar{\alpha})] \\ \text{Im}[z_t(\bar{\alpha})] \end{pmatrix}$$

ולכן הבעיה בניסוח הנוכחי היא:

$$\min_{\bar{\alpha}} \max_{t=1, \dots, T} |z_t(\bar{\alpha})| = \min_{\bar{\alpha}} \max_{t=1, \dots, T} \|\bar{u}_t(\bar{\alpha})\|_2$$

או באופן שקול:

$$\min_{\bar{\alpha}, \tau}$$

$$S.T \text{ (Subject To): } \tau \geq \|\bar{u}_t(\bar{\alpha})\|_2, \quad t = 1, \dots, T$$

וזו בעיה תכנות קוני ריבועית, כאשר האילוף הוא אילוף קוני מסדר שני כי אפשר לחשוב על אי-השוויון באילוף כעל הביטוי הבא:

$$\begin{pmatrix} \bar{u}_t(\bar{\alpha}) \\ \tau \end{pmatrix} = \begin{pmatrix} \text{Re}[z_t(\bar{\alpha})] \\ \text{Im}[z_t(\bar{\alpha})] \\ \tau \end{pmatrix} \in L^3$$

כאשר  $L^3$  זה הקונוס של לורנץ וזה נכון כיוון שלמדנו שבקונוס של לורנץ הרכיב האחרון (במקרה שלנו  $\tau$ ) צריך להיות גדול מהנורמה של כל הרכיבים האחרים (במקרה שלנו  $\bar{u}_t(\bar{\alpha})$ ), כאשר  $\bar{u}_t(\bar{\alpha})$  הוא למעשה שני רכיבים).

אנחנו נרצה לנסח את הבעיה האחרונה:

$$\min_{\bar{\alpha}, \tau}$$

$$S.T \text{ (Subject To): } \tau \geq \|\bar{u}_t(\bar{\alpha})\|_2, \quad t = 1, \dots, T$$

בצורה נוחה יותר, בכך שנצטרך למזער פונקציה לינארית שתלויה בווקטור שמורכב מהמשתנים  $\bar{\alpha}, \tau$  כך שהווקטור הנ"ל שייך לקונוס כלשהו שנגדיר. נגדיר:

$$\bar{v}_t \triangleq \begin{pmatrix} \bar{u}_t(\bar{\alpha}) \\ \tau \end{pmatrix} \triangleq \begin{pmatrix} \text{Re}[z_t(\bar{\alpha})] \\ \text{Im}[z_t(\bar{\alpha})] \\ \tau \end{pmatrix} \in L^3$$

ולכן נוכל לומר כי מתקיים:

$$\bar{v}_t \in \mathbb{R}^{3T}, \quad \bar{v}_t(\bar{\alpha}) = \begin{pmatrix} v_1(\bar{\alpha}) \\ v_2(\bar{\alpha}) \\ \vdots \\ v_T(\bar{\alpha}) \end{pmatrix} = \begin{pmatrix} \bar{u}_1(\bar{\alpha}) \\ \tau \\ \bar{u}_2(\bar{\alpha}) \\ \tau \\ \vdots \\ \bar{u}_T(\bar{\alpha}) \\ \tau \end{pmatrix}$$

אנחנו רוצים לקבל בעיה מהסוג הבא:

$$\begin{aligned} \min_{\bar{\alpha}, \tau} \tau \\ S.T \text{ (Subject To): } \bar{v}(\bar{\alpha}) \in K \end{aligned}$$

ולכן נגדיר את הקונוס הבא:

$$K = \underbrace{L^3 \times L^3 \times \cdots \times L^3}_{T \text{ times}}$$

ובכך למעשה הפכנו את הבעיה של תכנות קוני ריבועי לבעית תכנות קוני לינארית.

כל הניתוח והניסוח של הבעיה לעיל היה תחת ההנחה כי רכיבי הווקטור  $\bar{\alpha}$  הם מרוכבים, כלומר  $\alpha_i \in \mathbb{C}$ . אם ברצוננו לנסח בעית אופטימיזציה תוך שימוש במספרים ממשיים, אז ניתן להגדיר:

$$\bar{x} = \begin{pmatrix} \tau \\ \operatorname{Re}[u_1(\bar{\alpha})] \\ \operatorname{Im}[u_1(\bar{\alpha})] \\ \operatorname{Re}[u_2(\bar{\alpha})] \\ \operatorname{Im}[u_2(\bar{\alpha})] \\ \vdots \\ \operatorname{Re}[u_n(\bar{\alpha})] \\ \operatorname{Im}[u_n(\bar{\alpha})] \end{pmatrix}$$

וכעת אפשר לנסח את בעית האופטימיזציה בתור בעית תכנות קוני ריבועית (תוך שימוש בווקטור  $\bar{x}$  שהגדרנו זה עתה):

$$\begin{aligned} \min_{\bar{x}} \bar{c}^T \bar{x} \\ S.T \text{ (Subject To): } A\bar{x} - \bar{b} \geq_k 0 \end{aligned}$$

כאשר נשאר לקורא למצוא את המטריצה  $A$  והווקטורים  $\bar{c}, \bar{b}$  המתאימים לבעיה של קירוב צ'בישב.

## הרצאה 17 – Conversion of different problems to SDP (Semi Definite Programming)

הבסיס שעליו מתבססות על השיטות להפיכת בעיות אופטימיזציה שונות לבעית אופטימיזציה מסוג SDP (Semi Definite Programming) הוא אחד והוא למעשה הלמה של Schur.

למה: (Schur Complement)

נניח כי מטריצה  $A$  היא מטריצת בלוקים (כל איבר במטריצה הוא מטריצה בפני עצמו) סימטרית מהצורה:

$$A = \begin{pmatrix} P & Q^T \\ Q & R \end{pmatrix}, \quad R \succ 0$$

אז מתקיים:

- המטריצה  $A$  היא אי-שלילית  $\Leftrightarrow$  המטריצה  $P - Q^T R^{-1} Q$  אי-שלילית.
- המטריצה  $A$  היא חיובית  $\Leftrightarrow$  המטריצה  $P - Q^T R^{-1} Q$  חיובית.

הוכחה: (Schur Complement)

נבדוק את החיוביות של  $A$  באמצעות ההגדרה ונזכיר כי מטריצה  $W$  היא חיובית אם לכל  $\bar{x}$  מתקיים  $\bar{x}^T W \bar{x} > 0$ .

במקרה שלנו אנחנו צריכים לדרוש:

$$\forall \bar{u}, \bar{v}: \begin{pmatrix} \bar{u}^T & \bar{v}^T \end{pmatrix} \begin{pmatrix} P & Q^T \\ Q & R \end{pmatrix} \begin{pmatrix} \bar{u} \\ \bar{v} \end{pmatrix} \geq 0, \quad \bar{x} \triangleq \begin{pmatrix} \bar{u} \\ \bar{v} \end{pmatrix}$$

נבחין כי אם אנחנו דורשים כי אי-השוויון יתקיים לכל  $\bar{u}, \bar{v}$  אז בהכרח הוא יתקיים גם עבור  $\inf \bar{v}$  ולכן נוכל לומר כי אם אי-השוויון לעיל מתקיים אז גם מתקיים:

$$\forall \bar{u}: \inf_{\bar{v}} \begin{pmatrix} \bar{u}^T & \bar{v}^T \end{pmatrix} \begin{pmatrix} P & Q^T \\ Q & R \end{pmatrix} \begin{pmatrix} \bar{u} \\ \bar{v} \end{pmatrix} \geq 0$$

אם נפתח את הסוגריים ונבצע את המכפלה בצורה מפורשת נקבל פונקציה ריבועית ולכן אנחנו למעשה מתבקשים לפתור בעית אופטימיזציה לפי ווקטור  $\bar{v}$ , ואנחנו יודעים לנסח את התנאים ההכרחיים לפתרון בעיות מהסוג הזה, למשל שהגרדיאנט ביחס לווקטור  $\bar{v}$  צריך להתאפס בנקודה שהיא הפתרון וכו'. לאחר הניסוח אנחנו נראה כי הדרישה שאי-השוויון לעיל מתקיים הוא שקול (אמ"מ) לדרישה שיתקיים (כאשר נתון לנו כי  $R \succ 0$ ):

$$\bar{u}^T (P - Q^T R^{-1} Q) \bar{u} \geq 0$$

ולפי ההגדרה של מטריצה אי-שלילית, המטריצה  $P - Q^T R^{-1} Q$  היא אי-שלילית אמ"מ אי-השוויון לעיל מתקיים ולכן הוכחנו את הנדרש. נכתוב את הפתרון בתורה מפורשת (פתיחת הסוגריים וכו'):

$$\begin{pmatrix} \bar{u}^T & \bar{v}^T \end{pmatrix} \begin{pmatrix} P & Q^T \\ Q & R \end{pmatrix} \begin{pmatrix} \bar{u} \\ \bar{v} \end{pmatrix} = \bar{u}^T P \bar{u} + \bar{v}^T R \bar{v} + \bar{u}^T Q^T \bar{u} + \bar{v}^T Q \bar{u} = \bar{u}^T P \bar{u} + \bar{v}^T R \bar{v} + 2 \bar{u}^T Q^T \bar{u}$$



כעת נמצא את המינימום לפי  $\bar{v}$  ולשם כך נמצא את הגרדיאנט לפי  $\bar{v}$ :

$$\nabla_{\bar{v}} (\bar{u}^T P \bar{u} + \bar{v}^T R \bar{v} + \bar{u}^T Q^T \bar{u} + \bar{v}^T Q \bar{u}) = \bar{u}^T P \bar{u} + \bar{v}^T R \bar{v} + 2\bar{u}^T Q^T \bar{u} = 2Q\bar{u} + 2R\bar{v} = 0$$

$$\Rightarrow v = -R^{-1}Q$$

ולכן נציב את הנקודה  $v = -R^{-1}Q$  בחזרה לביטוי  $(\bar{u}^T \quad \bar{v}^T) \begin{pmatrix} P & Q^T \\ Q & R \end{pmatrix} \begin{pmatrix} \bar{u} \\ \bar{v} \end{pmatrix}$  ונקבל לאחר פישוט:

$$\bar{u}^T (P - Q^T R^{-1} Q) \bar{u} \geq 0$$

וזה נכון אם  $P - Q^T R^{-1} Q \succeq 0$ . מש"ל.

מזעור של הערך העצמי המקסימלי (Minimize Maximal Eigenvalue):

בעית התכונות הראשונה שנדגים כיצד מבצעים לה המרה לבעית תכונות SDP היא תהיה בעית המזעור של הערך העצמי המקסימלי. הבעיה מנוסחת כדלקמן:

נניח כי  $\bar{x} \in \mathbb{R}^n$  וכן נסמן:

$$A(\bar{x}) \triangleq A(\bar{x}) - \bar{b} \triangleq \sum_{i=1}^n x_i A_i - B, \quad A_1, A_2, \dots, A_n, B \in S^m$$

(כאשר  $S^m$  זו קבוצת המטריצות הסימטריות מסדר  $[m \times m]$ )

בנוסף נסמן את האופרטור  $\lambda_{\max}(\cdot)$  להיות אופרטור שמחזיר את הערך העצמי המקסימלי של הארגומנט של האופרטור (שהוא כמובן מטריצה). בעית האופטימיזציה היא:

$$\min_{\bar{x}} \lambda_{\max}(A(\bar{x}))$$

בהרבה שימושים הנדסיים, פעמים רבות המטרה היא להקטין את הערך העצמי של המטריצה ולכן זו בעיה מאד שימושית.

כיצד ניתן לפתור זאת? נניח כי הערכים העצמיים של מטריצה  $A$  הם  $\lambda_i$  ונמקם בתוך מטריצה אלכסונית שכל האיברים באלכסון שלה הם למעשה הע"ע. נוכל לרשום את הבעיה  $\min_{\bar{x}} \lambda_{\max}(A(\bar{x}))$  בתור בעית אופטימיזציה של סקלר כלשהו (שהוא למעשה הערך העצמי המקסימלי) ולכן הבעיה תהיה:

$$\min_{t, \bar{x}} t$$

$$S.T \text{ (Subject To): } t \geq \lambda_{\max}(A(\bar{x}))$$

את האילוץ לעיל ניתן לרשום גם בצורה אחרת. נוכל לרשום את זה בתור  $n$  אילוצים:

$$t - \lambda_i \geq 0, \quad \forall i = 1, 2, \dots, n$$

ואם נרשום בצורה מטריצית נקבל (דורש הוכחה):

$$tI - A(\bar{x}) \succeq 0$$

כלומר האילוץ הוא למעשה שהמטריצה  $tI - A(\bar{x})$  תהיה אי-שלילית. ואי-השוויון לעיל שקול (דורש הוכחה) לדרישה:

$$\begin{pmatrix} t - \lambda_1 & & 0 \\ & \ddots & \\ 0 & & t - \lambda_n \end{pmatrix} \succeq 0$$

ולכן הבעיה המקורית שלנו (מזעור הערך העצמי המקסימלי) בניסוח SDP היא:

$$\begin{aligned} \min_{t, \bar{x}} t \\ \text{S.T (Subject To): } tI - A(\bar{x}) \succeq 0 \end{aligned}$$

בקרב נלמד כי ישנן שיטות יעילות לחישוב ופתרון הבעיה הזו בניסוח SDP.

קירוב לינארי של מטריצה (Linear Matrix Approx.):

נושא שקשור לבעיית מזעור של הע"ע המקסימלי היא בעיית הקירוב הלינארי של מטריצה, ובעיית הקירוב הלינארי של מטריצה היא בעיה נפוצה מאד בקרב מהנדסים בתחומים שונים. בעיית הקירוב הלינארי של מטריצה נתונה היא:

ברצוננו לקרב את המטריצה  $A_0$  באמצעות קומבינציה לינארית של מטריצות בסיס (מטריצות כלשהן שנבחר), למשל מטריצות פשוטות כלשהן עם תכונות שאנחנו מעוניינים בהן  $\{A_1, \dots, A_n\}$ . כלומר נקבל:

$$A_0 \approx \sum_{i=1}^n x_i A_i, \quad \forall i = 1, \dots, n : x_i \in \mathbb{R}$$

ברצוננו לקרב את המטריצה  $A_0$  ע"י צירוף לינארי של מטריצות הבסיס שלנו  $\sum_{i=1}^n x_i A_i$  במובן של שגיאה מינימלית, כאשר השגיאה מוגדרת להיות:

$$\text{Error}(A(\bar{x})) \triangleq \left( \sum_{i=1}^n x_i A_i \right) - A_0$$

ואנחנו נרצה למזער את הנורמה של מטריצת השגיאה  $\text{Error}(A(\bar{x}))$ . ניתן להגדיר נורמות שונות עבור מטריצות ואנחנו נרצה למעשה למזער את האופרטור נורמה שפועל על המטריצה. אופרטור הנורמה של מטריצה הוא מוגדר ע"י:

$$\text{Norm}(A) \triangleq \|A\| \triangleq \max_{\|\bar{u}\|_2 \leq 1} \|A\bar{u}\|_2$$

כמו כן, ידוע מאלגברה לינארית כי הנורמה של מטריצה מקיימת:

$$\|A\| = \sigma_{\max}(A) = \sqrt{\lambda_{\max}(A^T A)}$$

כאשר  $\sigma_{\max}(A)$  זה הערך הסינגולרי המקסימלי של המטריצה  $A$ . בסך הכול, בעיית הקירוב הלינארי של מטריצה  $A_0$  כעת ניתנת לרישום ע"י הבעיה:

$$\min_{\bar{x}} \lambda_{\max} \left( \left( A(\bar{x}) \right)^T A(\bar{x}) \right)$$

אך את הבעיה  $\min_{\bar{x}} \lambda_{\max} \left( \left( A(\bar{x}) \right)^T A(\bar{x}) \right)$  אנחנו יודעים לפתור, כי זו בעיה ששקולה לבעיה:

$$\min_{t, \bar{x}} t, \quad t \in \mathbb{R}$$

$$S.T \text{ (Subject To): } \quad t^2 I - \left( A(\bar{x}) \right)^T A(\bar{x}) \succeq 0 \\ t \geq 0$$

כאשר האילוץ  $t^2$  הוא לשם נוחות. נבחין כי הבעיה שרשמנו לעיל היא אינה בעית SDP רגילה כי האילוץ הוא לא לינארי (כמו שהיינו רוצים) אך ניתן לרשום בעיה זו באופן שונה (תוך שימוש בלמה (Schur Complement):

$$\min_{t, \bar{x}} t, \quad t \in \mathbb{R}$$

$$S.T \text{ (Subject To): } \quad \begin{pmatrix} tI & A^T(\bar{x}) \\ A(\bar{x}) & tI \end{pmatrix} \succeq 0 \\ t > 0$$

$$\text{(כאשר } \begin{pmatrix} tI & A^T(\bar{x}) \\ A(\bar{x}) & tI \end{pmatrix} \text{ זו מטריצת בלוקים)}$$

לסיכום, מצאנו איך לבטא את בעית קירוב מטריצה כלשהי ע"י צירוף לינארי באמצעות בעית SDP.

תכנות לינארי באמצעות תכנות SDP (LP to SDP):

בחלק זה נלמד איך לבטא מספר בעיות אופטימיזציה לינאריות ידועות באמצעות בעיות SDP שקולות. נושא זה הוא בעיקר חשוב לידע התיאורטי שכן פתרון בעיות אופטימיזציה לינאריות הוא יעיל יותר מפתרון בעיות SDP שקולות אך חשוב לראות את הקשר בין שני סוגי הבעיות. למעשה בעית אופטימיזציה לינארית, היא מקרה פרטי של בעית SDP.

נניח ונתונה לנו בעית האופטימיזציה הלינארית הבאה:

$$\min_{\bar{x} \in \mathbb{R}^n} \bar{c}^T \bar{x}$$

$$, \bar{c} \in \mathbb{R}^n$$

$$S.T \text{ (Subject To): } \quad \bar{a}_i^T \bar{x} - \bar{b}_i \geq 0, \quad i = 1, \dots, m, \quad \bar{a}_i, \bar{b}_i \in \mathbb{R}^n$$

ובעית ה-SDP השקולה שלה היא:

$$\min_{\bar{x} \in \mathbb{R}^n} \bar{c}^T \bar{x}$$

$$, \bar{c} \in \mathbb{R}^n$$

$$S.T \text{ (Subject To): } \quad \begin{pmatrix} \bar{a}_1^T \bar{x} - \bar{b}_1 & & 0 \\ & \ddots & \\ 0 & & \bar{a}_n^T \bar{x} - \bar{b}_n \end{pmatrix} \succeq 0$$

(קל לראות שזה שקול כיוון שהמטריצה

$$\begin{pmatrix} \bar{a}_1^T \bar{x} - \bar{b}_1 & & 0 \\ & \ddots & \\ 0 & & \bar{a}_n^T \bar{x} - \bar{b}_n \end{pmatrix}$$

תהיה אי-שלילית אם"מ כל הערכים העצמיים שלה

הם חיוביים, אך אם המטריצה היא אלכסונית (כמו במקרה זה) אז הערכים העצמיים נמצאים על האלכסון הראשי שלה ולכן למעשה אנחנו דורשים שכל אחד מהערכים העצמיים הוא אי-שלילי ונבחין כי כל אחד מהערכים העצמיים הוא למעשה האילוץ בבעיה המקורית ולכן זו אכן בעיה שקולה)

תכנות בעיות תכנות קוני ריבעיות (CQP) באמצעות תכנות SDP (CQP via SDP):

בעיות תכנות קוני ריבעיות, בדומה לבעיות תכנות לינארי, יכולות להיות מבוטאות ע"י בעיות אופטימיזציה מסוג SDP. נתבונן בבעית CQP כללית:

$$\begin{aligned} \min_{\bar{x} \in \mathbb{R}^n} \bar{c}^T \bar{x} \quad , \bar{c} \in \mathbb{R}^n \\ \text{S.T (Subject To):} \quad & \left\| \frac{A\bar{x} - \bar{b}}{\bar{u}(\bar{x})} \right\|_2 \leq \underbrace{\bar{e}^T \bar{x} - d}_{v(\bar{x}) \in \mathbb{R}} \quad , \begin{pmatrix} \bar{u}(\bar{x}) \\ v(\bar{x}) \end{pmatrix} \in L^n \quad , v \in \mathbb{R} \\ & , d \in \mathbb{R}, \bar{e} \in \mathbb{R}^n \end{aligned}$$

ניתנת לתיאור ע"י בעית SDP שקולה באופן הבא (תוך שימוש בלמה (Schur Complement):

$$\begin{aligned} \min_{\bar{x}} \bar{c}^T \bar{x} \quad , \bar{x} \in \mathbb{R}^n \\ \text{S.T (Subject To):} \quad & \begin{pmatrix} (v) & (u_1 \cdots u_n) \\ (u_1) & \begin{pmatrix} v & & 0 \\ & \ddots & \\ 0 & & v \end{pmatrix} \end{pmatrix} \succeq 0 \quad , \quad v > 0 \end{aligned}$$

נבחין כי תוך שימוש בלמה Schur Complement ניתן לכתוב את האילוץ לעיל בתור האילוץ:

$$\bar{v} - \bar{u} \begin{pmatrix} v & & 0 \\ & \ddots & \\ 0 & & v \end{pmatrix} \bar{u} \geq 0 \quad , \quad \bar{v} = \begin{pmatrix} v \\ \vdots \\ v \end{pmatrix} \quad , v \in \mathbb{R}$$

וזה שקול (לאחר מניפולציות) לאי-השוויון:

$$v - v^{-1} \bar{u}^T \bar{u} \geq 0$$

ואם נכפול בסקלר  $v$  נקבל:

$$v^2 - \bar{u}^T \bar{u} \geq 0$$

וזה אכן שקול לאילוץ המקורי.

שיטת המחסום עבור בעיות תכנות קוני (Barrier method for QP):

בחלק זה נלמד איך לפתור בעיות תכנות קוני כלליות באמצעות שיטת המחסום. כבר למדנו את שיטת המחסום בבעיות תכנות לא לינארי והרעיון הכללי הזה לרעיון שראינו בעבר. נניח ויש לנו בעית תכנות קוני מהצורה הבאה:

$$\min_{\bar{x} \in \mathbb{R}^n} \bar{c}^T \bar{x}, \quad \bar{c} \in \mathbb{R}^n$$

$$S.T \text{ (Subject To):} \quad A(\bar{x}) - B \geq_K 0$$

$$A(\bar{x}) \triangleq \sum_{i=1}^n A_i \bar{x}_i \quad \text{וכן } \{A_1, A_2, \dots, A_n\} \text{ מטריצות}$$

בנוסף נניח כי קיימת נקודה  $\bar{x}$  בתוך תחום האילוף (Feasible Area) ולא רק על גבולות התחום (כלומר ניתן למצוא סביבה שלמה של הנקודה שכל הסביבה נמצאת גם היא בתוך התחום) אשר מקיימת (אי-שיוויון חזק):

$$A(\bar{x}) - B >_K 0$$

ולכן נוכל להגדיר פונקציית מחסום אשר יהיו לה ערכים קטנים (מתונים) יחסית בתוך הקונוס ובאזור גבולות הקונוס ערכי הפונקציה ילכו וישאפו לאינסוף. נגדיר את פונקציית המחסום בצורה הבאה:

$$\text{Barrier function:} \quad \mathcal{G}: K \rightarrow \mathbb{R}$$

$$\lim_{\substack{u \rightarrow \partial K \\ u \in K}} \mathcal{G}(u) = \infty$$

כאשר  $\partial K$  זה השפה של הקונוס  $K$ .

פונקציית המחסום המצרפית (Barrier aggregate):

נגדיר את פונקציית המחסום המצרפית, בדומה לאיך שהגדרנו את פונקציית העונשין המצרפית (סכום פונקציית המטרה ועוד סכום פונקציית העונשין), בצורה הבאה:

$$F_\mu(\bar{x}) = \bar{c}^T \bar{x} + \frac{1}{\mu} \mathcal{G}(A(\bar{x}) - B), \quad \mu \in \mathbb{R}$$

נבחין שעבור פרמטר  $\mu$  גדול, נקבל כי עבור איבר  $A(\bar{x}) - B$  בתוך התחום ממש הביטוי  $\frac{1}{\mu} \mathcal{G}(A(\bar{x}) - B)$  יהיה מאד קטן, אך בגבולות הקונוס אנחנו עדיין נקבל כי הביטוי שואף לאינסוף ולכן ניתן לרשום:

$$F_\infty(\bar{x}) = \begin{cases} \bar{c}^T \bar{x} & , A(\bar{x}) - B >_K 0 \\ \infty & , A(\bar{x}) - B \leq_K 0 \end{cases}$$

ברור כי אם נמצא את המינימום של פונקציה המחסום המצרפית אז זה יהיה גם הפתרון של הבעיה:

$$\left. \begin{aligned} \min_{\bar{x} \in \mathbb{R}^n} F_\infty(\bar{x}) &= \min_{\bar{x} \in \mathbb{R}^n} \bar{c}^T \bar{x}, \quad \bar{c} \in \mathbb{R}^n \\ S.T \text{ (Subject To):} \quad &A(\bar{x}) - B \geq_K 0 \end{aligned} \right\} \Rightarrow \min_{\bar{x}} F_\infty(\bar{x}) = f^*$$

והאלגוריתם הוא להתחיל לפתור את בעית האופטימיזציה עם פרמטר  $\mu$  קטן יחסית, ובכל איטרציה להכפיל אותו פי כמה (למשל פי 2 עד 10) עד אשר הוא גדל להיות בסדר גודל של  $10^5$ .

דוגמאות של פונקציות מחסום מוכרות:

פונקצית מחסום $\mathcal{G}(\bar{u})$	קונוס
$-\sum_{i=1}^m \log(u_i)$	$\mathbb{R}_+^m$
$-\log\left(u_m^2 - \sum_{i=1}^{m-1} u_i^2\right)$	$L^m$ (קונוס לורנץ)
$-\log(\det(U))$	$S_+^m$ (קונוס המטריצות הסימטריות האי-שליליות)

פונקצית המחסום:  $-\log(\det(U))$

נניח ויש לנו מטריצה  $A$  וניתן לבטא אותה באמצעות מטריצה  $SAS^{-1}$  כאשר  $S$  מטריצה שעמודותיה הם הווקטורים העצמיים של מטריצה  $A$  וכן  $\Lambda$  זו מטריצה אלכסונית שעל האלכסון שלה נמצאים כל הערכים העצמיים של מטריצה  $A$  אז ניתן לכתוב:

$$\begin{aligned}\det(A) &= \det(S^{-1}\Lambda S) = \det(S^{-1})\det(\Lambda)\det(S) = \\ &= (\det(S))^{-1}\det(\Lambda)\det(S) = \det(\Lambda)(\det(S))^{-1}\det(S) = \\ &= \det(\Lambda) = \prod_i \lambda_i\end{aligned}$$

ולכן נקבל:

$$-\log(\det(U)) = -\log\left(\sum_{i=1}^m \lambda_i\right) = -\sum_{i=1}^m \log(\lambda_i)$$

פונקציה מטריצית (Matrix functions): (זו תזכורת, שכן כבר ראינו את הנושא בהרצאות הראשונות)

כאשר רוצים לפתור בעיה אופטימיזציה, למשל באמצעות שיטת המחסום, אנחנו למעשה צריכים לחשב את הגרדיאנט של פונקציה המחסום ולעיתים פונקצית המחסום מוגדרת ע"י מטריצה ולכן אנחנו צריכים להיזכר באנליזה מטריצית: כיצד גוזרים מטריצה וכיצד מפעילים אופרטורים שונים על מטריצה.

נניח ויש לנו פונקציה של משתנה יחיד (ממשי) בעל פיתוח פולינומי (למשל טיילור):

$$\varphi(t) = \sum_{i=0}^{\infty} c_i t^i$$

אז נגדיר פונקציה מטריצית באופן הבא (גם ע"י פיתוח פולינומי):

$$\varphi(A) = \sum_{i=0}^{\infty} c_i A^i$$

כאשר  $A^i$  אומר להכפיל את המטריצה  $A$  בעצמה,  $i$  פעמים.

נבחין כי אם ניתן לבטא את המטריצה  $A$  באמצעות פירוק למטריצת הערכים העצמיים,  $A = S^{-1} \Lambda S$  (שהיא מטריצה אלכסונית) אז החישוב של הפונקציה  $\varphi(A) = \sum_{i=0}^{\infty} c_i A^i$  שהגדרנו, נהיה פשוט הרבה יותר, כי נזכור כי להכפיל מטריצות אלכסוניות אחת בשניה זה לכפול את הערכים על האלכסון הראשי, אחד בשני.

פירוק לערכים עצמיים / מציאת הערכים העצמיים של מטריצה (Eigenvalue decomposition):

נזכיר תחילה כי ערך עצמי  $\lambda_i$  של מטריצה  $A$  הוא ערך אשר עבורו קיים וקטור  $\bar{s}_i \in \mathbb{R}^n$  כך שמתקיים:

$$A \bar{s}_i = \lambda_i \bar{s}_i$$

נניח כי למטריצה  $A$  הריבועית בגודל  $[m \times m]$  יש  $m$  ערכים עצמיים ווקטורים עצמיים (ולפי אלגברה לינארית אנחנו יודעים כי גם כל הווקטורים העצמיים הם בלתי תלויים לינארית אחד בשני). נגדיר את המטריצה  $S$  להיות מטריצת עמודות, שבה כל עמודה היא וקטור עצמי של מטריצה  $A$ . לכן נוכל לרשום (דורש הוכחה):

$$AS = S\Lambda$$

כיוון שכל הווקטורים העצמיים של  $A$  הם בלתי תלויים, כלומר כל העמודות של מטריצה  $S$  הן בלתי תלויות, אז זה אומר כי מטריצה  $S$  היא מדרגה מלאה ולכן קיימת לה הופכית. ולכן נכפול את השיוויון האחרון מימין במטריצה  $S^{-1}$  ונקבל:

$$A = SAS^{-1}$$

נזכר כי הגדרנו  $\varphi(A) = \sum_{i=0}^{\infty} c_i A^i$  ולכן אנחנו צריכים לחשב את הכפולות של מטריצה  $A$  בעצמה. נבחין כי מתקיים:

$$A^2 = (S\Lambda S^{-1})^2 = SAS^{-1}S\Lambda S^{-1} = S\Lambda(S^{-1}S)\Lambda S^{-1} = S\Lambda^2 S^{-1} = S \begin{pmatrix} \lambda_1^2 & & 0 \\ & \ddots & \\ 0 & & \lambda_n^2 \end{pmatrix} S^{-1}$$

וניתן לראות באופן דומה (ע"י אינדוקציה) כי מתקיים:

$$A^i = S\Lambda^i S^{-1} = S \begin{pmatrix} \lambda_1^i & & 0 \\ & \ddots & \\ 0 & & \lambda_n^i \end{pmatrix} S^{-1}$$

ולכן נקבל:

$$\varphi(A) = \sum_{i=0}^{\infty} c_i A^i = \sum_{i=0}^{\infty} c_i S\Lambda^i S^{-1} = S \left( \sum_{i=0}^{\infty} c_i \Lambda^i \right) S^{-1} = S\varphi(\Lambda)S^{-1} = S \begin{pmatrix} \varphi(\lambda_1) & & 0 \\ & \ddots & \\ 0 & & \varphi(\lambda_m) \end{pmatrix} S^{-1}$$

בנוסף, ניתן לראות כי אם מטריצה  $A$  היא סימטרית, מתקבל  $S^{-1} = S^T$  ואז נקבל:

$$\varphi(A) = S\varphi(\Lambda)S^T = S \begin{pmatrix} \varphi(\lambda_1) & & 0 \\ & \ddots & \\ 0 & & \varphi(\lambda_m) \end{pmatrix} S^T$$

תכונה נוספת של הפונקציה המטריצית שהגדרנו:  $\varphi(A) = \sum_{i=0}^{\infty} c_i A^i$ , מתקיים:

$$\text{Trace}(\varphi(A)) = \text{Trace}(S\varphi(\Lambda)S^{-1}) \underset{\substack{\text{cyclic shift} \\ \text{under Trace}}}{=} \text{Trace}(S^{-1}S\varphi(\Lambda)) = \text{Trace}(\varphi(\Lambda)) = \sum_{i=1}^m \varphi(\lambda_i)$$

גרדיאנט של עכבה של פונקציה מטריצית:

נרצה לדעת מהו הביטוי  $\nabla_A \text{Trace}(\varphi(A))$  באופן מפורש כאשר  $\varphi: [m \times m] \rightarrow [m \times m]$  (פונקציה מטריצית!). מתקיים:

$$\nabla_A \text{Trace}(\varphi(A)) = (\varphi'(A))^T$$

כאשר  $\varphi(t): [m \times m] \rightarrow \mathbb{R}$  היא פונקציה סקלרית.

לא נוכיח את הטענה האחרונה אך נדגים את השימושיות והנכונות של הביטוי על דוגמה כללית יחסית.

נניח פונקציה סקלרית  $\varphi(t) = t^3$  ולכן הנגזרת שלה היא  $\varphi'(t) = 3t^2$ . בנוסף, נסמן:  $f(A) \triangleq \text{Trace}(\varphi(A))$  (ונבחין כי  $f: [m \times m] \rightarrow \mathbb{R}$  כלומר היא פונקציה סקלרית!) ולכן:

$$f(A) = \text{Trace}(\varphi(A)) = \text{Trace}(A \cdot A \cdot A)$$

הדיפרנציאל של  $f(A)$  הוא:

$$df(A) = \text{Trace}(dA \cdot A \cdot A) + \text{Trace}(A \cdot dA \cdot A) + \text{Trace}(A \cdot A \cdot dA)$$

אך לפי תכונת הציקליות של העכבה מתקיים:

$$df(A) = \text{Trace}(dA \cdot A \cdot A) + \text{Trace}(A \cdot dA \cdot A) + \text{Trace}(A \cdot A \cdot dA) = 3 \cdot \text{Trace}(dA \cdot A^2)$$

בנוסף לפי ההגדרה של מכפלה פנימית בין מטריצות:  $\langle A, B \rangle \triangleq \text{Trace}(A^T B)$  נקבל:

$$df(A) = 3 \cdot \text{Trace}(dA \cdot A^2) = \left\langle dA, (3A^2)^T \right\rangle \Leftrightarrow df(A) = \langle dA, \nabla f \rangle$$

ולכן מצאנו כי מתקיים:

$$\nabla f = (3A^2)^T = 3(A^2)^T = 3(A^T)^T = (\varphi'(A))^T$$

ולכן "הוכחנו" את השוויון לעיל, לפחות שהוא נכון עבור הדוגמה  $\varphi(t) = t^3$ .



חישוב הגרדיאנט של פונקצית מחסום:

תוך שימוש בטענה שראינו:  $\nabla_A f(A) = \nabla_A \text{Trace}(\varphi(A)) = (\varphi'(A))^T$  נרצה לחשב את הגרדיאנט של הפונקציה:

$$f(A) = \log(\det(A)), \quad A \succ 0$$

כאשר  $f$  היא פונקציה סקלרית. מתקיים:

$$f(A) = \log(\det(A)) = \log\left(\prod_{i=1}^m \lambda_i\right) = \sum_{i=1}^m \log(\lambda_i) = \text{Trace}(\log(A))$$

ולכן הפונקציה הסקלרית שנבחר היא  $\varphi(t) = \log(t)$  והנגזרת שלה היא:  $\varphi'(t) = \frac{1}{t} = t^{-1}$  ולכן:

$$\nabla_A f(A) = (\varphi'(A))^T = (A^{-1})^T = (A^T)^{-1} \equiv A^{-T}$$

(ואם  $A = A^T$  סימטרית אז נזכור כי מתקיים:  $A = A^T$ )

לסיכום: אם פונקצית המחסום שלנו היא:  $\varphi(A) = -\log(\det(A))$  אז  $\nabla_A \varphi(A) = -A^{-1}$ .

בתכנות SDP, אנחנו לרוב נוהגים לכתוב את הארגומנט של פונקציה המחסום בצורה:  $A(\bar{x}) \triangleq A(\bar{x}) - B$  ולכן אם נרצה

למצוא את הגרדיאנט של  $\varphi(A(\bar{x}))$  ביחס לווקטור  $\bar{x}$  נרצה להשתמש בהגדרה הדיפרנציאלית:

$$d\varphi(A(\bar{x})) = \langle \nabla_A \varphi(A(\bar{x})), dA(\bar{x}) \rangle$$

$$d(A(\bar{x})) = \langle \nabla_A \varphi(A(\bar{x})), dA(\bar{x}) \rangle = \langle \nabla_A (A(\bar{x})), dA\bar{x} \rangle = \langle A^* (\nabla_A \varphi(A(\bar{x}))), d\bar{x} \rangle$$

ולכן מצאנו למעשה:

$$\nabla_{\bar{x}} \varphi(A\bar{x} - B) = A^* (\nabla_A \varphi(A(\bar{x})))$$

כעת לאחר שראינו כיצד מחשבים את הגרדיאנט של פונקצית המחסום, ניתן למעשה להתחיל לפתור באמצעות האלגוריתם שאנחנו מכירים לתכנות בעיות אופטימיזציה, למשל שיטת ניוטון, למרות שעבור שיטת ניוטון צריך לחשב את ההסיאן, וזה מעבר למסגרת הקורס הזה. ניתן להשתמש גם באלגוריתם BFGS שלמדנו שלא דורש את חישוב ההסיאן.