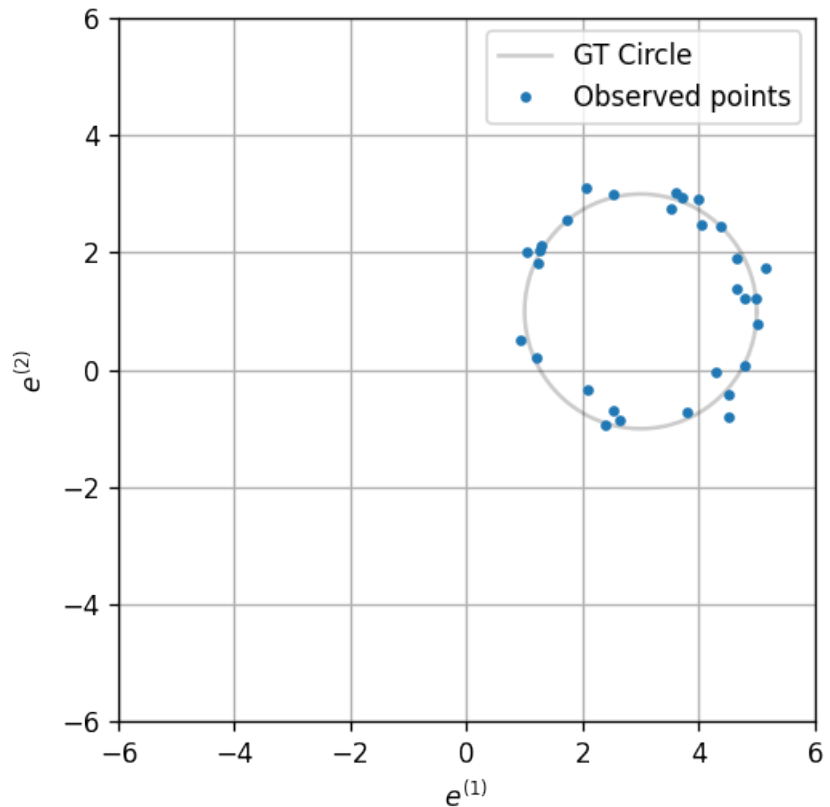


# תרגול 1 - בעיות אופטימיזציה וגזירה וקטורית

## הקדמה

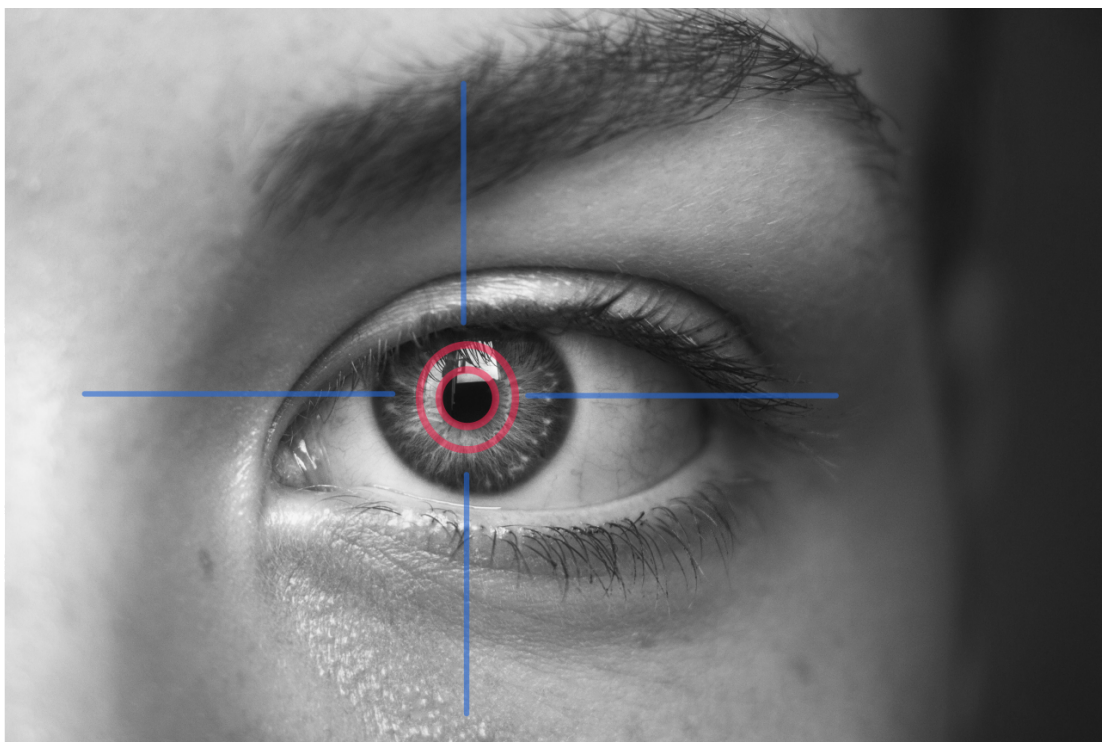
מרבית הבעיות והשיטות בתחום של מערכות לומדות עוסקות בבעיה של התאמת מודל מתמטי כך שיתאר בצורה מיטבית תופעה או תהליך מסויים על סמך אוסף נתון של תצפיות / מדידות. בהרצא הראשונה נציג את התחום באופן מפורט יותר, אך לעת אתה הגדרה זו תספק אותנו.

לצורך המחשה נסתכל על הדוגמה הבאה. נניח ונתונים לנו מספר מדידות של נקודות שיושבות על מעגל בעל מרכז ורדיוס לא ידועים. בנוסף נניח ותהליך המדידה עצמו רועש ואנו מקבלים גרסאות מורעשות של הנקודות כפי שמודגם בשרטוט הבא:



מקובל להשתמש בשם **truth ground**, או בקיצור **GT**, כדי להתייחס למודל המקורי (הלא ידוע).

כעת נניח ואנו מעוניינים לשחזר את הפרמטרים של המעגל המקורי על פי הדגימות שבידינו. הבעיה של התאמת בעיה כדוגמאת זו, של מעגל לאוסף נקודות, מופיעה באלגוריתמי eye-tracking אשר מנסים לעקוב אחרי המיקום של האישון על מנת להבין מהו הכיוון שאליו אדם מביט.



נציין שהבעיה הזו לא מאד מייצגת ונחשבת ליחסית פשוטה בהשוואה לבעיות הטיפוסיות שאותם מנסים לפתור בתחום של מערכות לומדות, אך עם זאת היא תשמש כדוגמא טובה לעקרונות שבהם נעסוק בתרגול הנוכחי.

### פערים מתמטיים

לפני שנוכל לצלול לחומר העיקרי של הקורס עלינו להתחיל בהשלמה ורענון של הבסיס המתימטי אשר ישמש אותנו לתורך ניסוח הפתרון של הבעיות בהם נעסוק בקורס. ספציפית אנו נעשה שימוש ב:

- אלגברה לינארית: על מנת לתאר את המידע שהמודלים שאיתם נעבוד.
  - הסתברות: על מנת לתאר את האופן שבו נוצרים המדידות ובכדי לתאר את מידת הסבירות שבה מודל מסוים מתאים לתצפיות.
  - תורת האופטימיזציה: על מנת למצוא את המודלים אשר מתאימים באופן מיטבי לדגימות שבידינו.
- בתרגול הנוכחי נתעסק בבעיות אופטימיזציה סקלריות ווקטוריות ובתרגול הבא נוסיף את לכך את ההיבט ההיסטברותי.

### נוטציות

בקורס זה נצמד לנוטציות המתמטיות המופיעות בספר [Learning Deep](#) (מאת A. & Bengio Y. Goodfellow, I. Courville). ניתן למצוא את הפירוט המלא בקישור הבא. (למתקדמים: ניתן למצוא את פקודות ה-*LaTeX* הרלוונטיות כאן)

### אלגברה לינארית

בהתאם להגדרות אלו אנו נשתמש בסימונים הבאים בהקשר של אלגברה לינארית:

- $x$  - אותיות סטנדרטיות (lower (italic case) לועזיות או יווניות - סקלרים.

- $x$  - אותיות מודגשות - וקטורי עמודה
- $x^\top$  - וקטורי שורה
- $x_i$  - האיבר ה- $i$  בוקטור  $x$ .
- $\langle x, y \rangle (= x^\top y = \sum_i x_i y_i)$  - המכפלה הוקטורית הסטנדרטית בין  $x$  ל  $y$ .
- $\|x\|_2 (= \sqrt{\langle x, x \rangle})$  - הנורמה הסטנדרטית (נורמת  $l_2$ ) של הוקטור  $x$ .
- $\|x\|_l (= \sqrt[l]{\sum_i x_i^l})$  - נורמת  $l$  של  $x$
- $A$  - אותיות לועזיות גדולות מודגשות (bold) - מטריצה
- $A^\top$  - המטריצה  $A$  Transposed (המטריצה המשוכלפת).
- $A_{i,j}$  - האיבר ה- $j$  שורה ה- $i$  של  $A$ .
- $A_{i,:}$  - השורה ה- $i$  של  $A$ .
- $A_{:,i}$  - העמודה ה- $i$  של  $A$ .

## Sets (קבוצות)

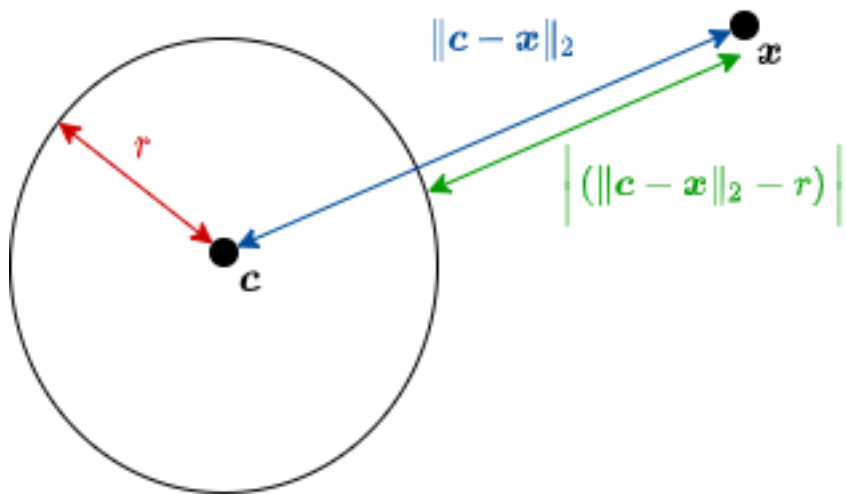
את סדרת התצפיות אנו נסמן כקבוצה (או סדרה) בעזרת הנוטציות הבאות:

- $\{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$  - סדרה של  $n$  וקטורים.

## בעיית האופטימיזציה

נחזור לדוגמא של התאמת המעגל. דרך אחת לגשת לבעיה מסוג זה הינה לחפש מבין כל המעגלים האפשריים את המעגל אשר המרחק הריבועי הממוצע (MSE - squared mean error) של הדגימות ממנו היא המינימאלית (בתרגול הבא אנו ניתן הצדקה מתימטית לשימוש בממדד שגיאה זה). נרשום זאת באופן מתמטי.

נרשום תחילה את המרחק של נקודה בודדת  $x$  ממעגל בעל מרכז ב  $c$  ורדיוס  $r$ . מרחק זה שווה להפרש בין המרחק בין  $x$  ל  $c$  והרדיוס, כפי שמופיע בשרטוט הבא:



נסמן את הריבוע של מרחק זה ב  $e$ :

$$e = (\|x - c\|_2 - r)^2$$

אם כן בעבור  $n$  נקודות,  $\{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$ , המרחק הריבועי הממוצע (MSE) של הנקודות מהמעגל הינו:

$$\frac{1}{n} \sum_{i=1}^n e_i = \frac{1}{n} \sum_{i=1}^n (\|x^{(i)} - c\|_2 - r)^2$$

מצאנו אם כן למעשה מדד טיב אשר מאפשר לנו לתת "ציון" לכל מעגל כתלות ב  $c$  ו  $r$ . נסמן את הפונקציה הזו כ:

$$f(c, r) = \frac{1}{n} \sum_{i=1}^n (\|x^{(i)} - c\|_2 - r)^2$$

כעת, אם ברצוננו למצוא עלינו למצוא את המעגל **האופטימאלי**, עלינו למצוא את הפרמטרים  $c$  ו  $r$  אשר מניבים את הערך הנמוך ביותר של  $f$ . מקובל לסמן את הפרמטרים האופטימאליים בעזרת \* באופן הבא:  $c^*$  ו  $r^*$ .

בעיות מסוג זה הם בדיוק סוגי הבעיות שתורת האופטימיזציה באה לפתור. לבנתיים נשים בצד את הפונקציה הזו ונחזור אליה לקראת סוף התרגול.

## המקרה הפשוט

בעיות אופטימיזציה עוסקות במציאת הארגומנט  $\theta$  שבעבורו פונקציה נתונה  $f(\theta)$  מחזירה את הערך המינימאלי או המקסימאלי שלה. לרוב מקובל לנסח בעיות אופטימיזציה כבעיות **minimization** (מזעור), כאשר ניתן כמובן לרשום כל בעיית **maximization** כבעיית minimization של  $\tilde{f}(\theta) = -f(\theta)$ . כמו כן, את פונקציה  $f(\theta)$  שאותה מנסים למזער (למקסם) נהוג לכנות ה**objective** או פונקציית המטרה. באופן פורמלי, בעיות אופטימיזציה נרשמות באופן הבא:

$$\theta^* = \arg \min_{\theta} f(\theta)$$

## פונקציות מרובות משתנים

בעיות אופטימיזציה כמובן לא מוגבלות רק לפונקציות של משתנה יחיד וכמו בדוגמה של המעגל, ניתן להסתכל גם על הבעיה של מציאת סט הערכים האופטימאליים  $\theta_1^*, \theta_2^*, \dots, \theta_d^*$  שמזערים פונקציה של מספר משתנים  $f(\theta_1, \theta_2, \dots, \theta_d)$ :

$$\theta_1^*, \theta_2^*, \dots, \theta_d^* = \arg \min_{\theta_1, \theta_2, \dots, \theta_d} f(\theta_1, \theta_2, \dots, \theta_d)$$

במקרים רבים נוח יותר לאגד את כל הארגומנטים של  $f$  לוקטור אחד  $\theta = (\theta_1, \theta_2, \dots, \theta_d)^\top$  ולרשום את בעיית האופטימיזציה כ:

$$\theta^* = \arg \min_{\theta} f(\theta)$$

## שיטות לפתרון בעיות אופטימיזציה

בעיות אופטימיזציה הם למעשה חלק מהתהליך של חקירת פונקציה אשר נלמד בחדו"א 1 ובתיכון, שם מצאנו נקודות הקיצון על ידי גזירה והשוואה ל-0. שיטה זו טובה למקרים פשוטים בהם המשוואה  $\nabla f(\theta) = 0$  ניתנת לפתרון באופן אנליטי, אך במקרים רבים זהו אינו המצב ויש צורך להשתמש בשיטות נומריות על מנת לפתור את הבעיה.

התחום של תורת האופטימיזציה הוא רחב והוא מציע מגוון רחב של שיטות כאלה המתאימות לפונקציות מטרה שונות. בפועל בקורס זה לא נצלול לעומקו של תחום זה, ולמעשה נפתור בעיות מסוג זה בעזרת אחת משלושת השיטות הבאות:

- **גזירה והשוואה ל-0**: מתאים למקרים בהם ניתן לפתור את  $\nabla f(\theta) = 0$  באופן אנליטי.
- **force Brute** (מעבר על כל האופציות): מתאים למקרים בהם הארגומנט של  $f$  לקבל רק אחד מסט ערכים סופי וקטן, לדוגמה  $\theta \in \{0, 1, \dots, 10\}$ .
- **descent Gradient** (שיטת הגרדיאנט): זוהי אחת משיטות האופטימיזציה הבסיסיות ביותר. בשיטה זו מחפשים מינימום לוקאלי על ידי התקדמות בכיוון ההפוך מהגרדיאנט. בתרגול זה אנחנו נדגים כיצד להשתמש בה.

## בעיית אופטימיזציה עם אילוצים

נשים לב שבדוגמא של התאמת המעגל, ישנו אילוץ מסויים על הפרמטרים של המעגל, רדיוס המעגל חייב להיות מספר חיובי. מגבלות מסוג זה מכונות אילוצים והם מופיעים בבעיות אופטימיזציה רבות. באופן כללי, בעיות אופטימיזציה יכולות להכיל אילוצים משני סוגים:

- אילוצי אי-שוויון מהצורה  $g(\theta) \geq 0$ . לדוגמא, האילוץ שכל הערכים של  $\theta$  יהיו קטנים מ-1, יופיע בתור:

$$g_i(\theta) = 1 - \theta_i \geq 0 \quad i = 1, \dots, d$$

- אילוצי שוויון מהצורה  $h(\theta) = 0$ . לדוגמא האילוץ שהנורמה של  $\theta$  תהיה שווה ל-1, יופיע בתור:

$$h(\theta) = \|\theta\|_2 - 1 = 0$$

אם כן הצורה הכללית ביותר של בעיית אופטימיזציה הינה:

$$\begin{aligned} \theta^* = \arg \min_{\theta} \quad & f(\theta) \\ \text{subject to} \quad & g_i(\theta) \leq 0, \quad i = 1, \dots, m \\ & h_j(\theta) = 0, \quad j = 1, \dots, p \end{aligned}$$

בתרגיל הבא נראה כיצד ניתן להתמודד עם אילוצים פשוטים ובהמשך הקורס נתייחס גם למקרים בהם האילוצים נעשים מסובכים יותר.

## תרגיל 1.1 - בעיית אופטימיזציה עם אילוצים

נרצה למצוא את הערך המקסימאלי של הפונקציה

$$f(\theta_1, \theta_2) = f(\theta) = e^{-(3\theta_1^2 + 3\theta_2^2 - 18\theta_1 - 24\theta_2 + 34)}$$

תחת האילוץ ש  $\theta$  נמצא בתוך או על השפה של מעגל היחידה (במישור של  $\theta_1 - \theta_2$ ).  
רשמו את בעיית האופטימיזציה ומצאו את פתרונה.

### פתרון 1.1

**רישום בעיית האופטימיזציה** הבעיה הנתונה הינה בעיית אופטימיזציה עם אילוץ אי-שוויון אחד. נרשום אותה באופן פורמלי:

$$\begin{aligned} \arg \min_{\theta = (\theta_1, \theta_2)^T} \quad & -e^{-(3\theta_1^2 + 3\theta_2^2 - 18\theta_1 - 24\theta_2 + 34)} \\ \text{subject to} \quad & 1 - (\theta_1^2 + \theta_2^2) \leq 0 \end{aligned}$$

**החלפת פונקציית המטרה** ראשית נשים לב כי ניתן לפשט את הבעיה על ידי השימוש בעובדה ש  $e^x$  היא פונקציה מונוטונית עולה ולכן נוכל לבצע את ההחלפה הבאה מבלי לשנות את תוצאת בעיית האופטימיזציה:

$$\begin{aligned} \arg \min_{\theta_1, \theta_2} \quad & -e^{-(3\theta_1^2 + 3\theta_2^2 - 18\theta_1 - 24\theta_2 + 34)} \\ = \arg \min_{\theta_1, \theta_2} \quad & 3\theta_1^2 + 3\theta_2^2 - 18\theta_1 - 24\theta_2 + 34 \end{aligned}$$

באופן דומה נוכל גם להיפתר מהתוספת של הקבוע  $+34$  (וגם לחלק את פונקציית המטרה, 3 ולקבל את בעיית האופטימיזציה הבאה:

$$= \arg \min_{\theta_1, \theta_2} \theta_1^2 + \theta_2^2 - 6\theta_1 - 8\theta_2$$

**טיפול באילוצי אי-השוויון** במקרים כגון זה, בהם מספר אילוצי אי-השוויון קטן (במקרה זה יש אילוץ בודד) נוכל לפרק כל אילוץ שוויון לשני מקרים פשוטים יותר.

1. המקרה בו המינימום נמצא על השפה של האילוץ (בתרגיל זה, זאת אומרת שהמינימום נמצאת ממש על מעגל היחידה). במקרה זה האילוץ הופך לאילוץ שוויון.  
2. המקרה בו המינימום לא נמצא על השפה של האילוץ. במקרה זה נחפש נקודות מינימום לוקאליות תוך התעלמות מהאילוץ ואחר כך נבדוק מי מהם מקיימת את האילוץ (אם בכלל יש נקודה כזו).  
הפתרון יהיה הפתרון הנמוך יותר מבין השניים. (כאשר יש יותר מאילוץ אי-שוויון יחיד צריך לפרק את כל אילוצי אי השוויון לשני מקרים. זאת אומרת, שכמות המקרים שיש לבדוק הינה 2 בחזקת מספר אילוצי אי-השוויון).

**החיפוש בתוך מעגל היחידה** נתחיל במציאת המינימום בחלקו הפנימי של העיגול. תחילה נתעלם מהאילוץ נחפש את כל נקודות המינימום (לוקאליות או גלובליות) של הבעיה. אחר כך נפסול את אלו שלא מקיימות את האילוץ. בעיית האופטימיזציה ללא האילוץ הינה:

$$\arg \min_{\theta_1, \theta_2} \theta_1^2 + \theta_2^2 - 6\theta_1 - 8\theta_2$$

בעיה זו ניתנת לפתרון בקלות על ידי גזירה והשוואה ל-0:

$$\begin{cases} \frac{d}{d\theta_1} \theta_1^2 + \theta_2^2 - 6\theta_1 - 8\theta_2 = 0 \\ \frac{d}{d\theta_2} \theta_1^2 + \theta_2^2 - 6\theta_1 - 8\theta_2 = 0 \end{cases} \\ \Leftrightarrow (\theta_1, \theta_2) = (3, 4)$$

נקודה זו אומנם חשודה כנקודת קיצון אך היא לא מקיימת את האילוץ ולכן היא אינה יכול להיות פתרון לבעיה. נוכל להסיק אם כן שאין נקודות מינימום בתוך המעגל ולכן נקודת המינימום תהיה חייבת להימצא על השפה.

**החיפוש על מעגל היחידה** על השפה אילוץ האי-השוויון הופך לשוויון:

$$\begin{aligned} \arg \min_{\theta=(\theta_1, \theta_2)^\top} \theta_1^2 + \theta_2^2 - 6\theta_1 - 8\theta_2 \\ \text{subject to } 1 - (\theta_1^2 + \theta_2^2) = 0 \end{aligned}$$

בתרגול הבא נלמד שיטה מסודרת לפתרון בעיות אופטימיזציה עם אילוצי שוויון כגון זה, אך לבנתיים נתאר כאן פתרון חליפי אשר משתמש בניסוח מחד של הבעיה (השיטה שמוצגת כאן לא רלוונטית להבנת החומר מופיעה פה רק לשם השלמות).

הדרך בה נימצא את המינימום על המעגל הינה על ידי החלפת הבעיה הנתונה בבעיית אופטימיזציה של משתנה יחיד ללא אילוצים. אנו נשתמש באילוץ (התנאי שהנקודה תמצא על מעגל היחידה) על מנת לבטא את  $\theta_2$  בעזרת  $\theta_1$  ורישום מחדש של פונקציית המטרה כפונקציה של  $\theta_1$  בלבד.

מתוך האילוץ נקבל ש:

$$\begin{aligned} 1 - (\theta_1^2 + \theta_2^2) &= 0 \\ \Leftrightarrow \theta_2 &= \pm \sqrt{1 - \theta_1^2} \end{aligned}$$

כאשר עלינו לבדוק את שני המקרים כאשר  $\theta_2$  חיובי (החלק העליון של מעגל היחידה) וכשאר הוא שלילי (החלק התחתון). לאחר ההחלפה נקבל את שני בעיות האופטימיזציה הבאות (המקרה החיובי והשלילי):

$$\begin{aligned} \arg \min_{\theta_1} \quad & \theta_1^2 + (1 - \theta_1^2) - 6\theta_1 \pm 8\sqrt{1 - \theta_1^2} \\ = \arg \min_{\theta_1} \quad & 1 - 6\theta_1 \pm 8\sqrt{1 - \theta_1^2} \\ = \arg \min_{\theta_1} \quad & -3\theta_1 \pm 4\sqrt{1 - \theta_1^2} \end{aligned}$$

נפתור את את הבעיות האלה על ידי גזירה והשוואה ל-0:

$$\begin{aligned} \frac{d}{d\theta_1} - 3\theta_1 \pm 4\sqrt{1 - \theta_1^2} &= 0 \\ \Leftrightarrow -3 \pm 4 \frac{\theta_1}{\sqrt{1 - \theta_1^2}} &= 0 \\ \Leftrightarrow \pm \frac{4}{3}\theta_1 &= \sqrt{1 - \theta_1^2} \\ \Rightarrow \frac{16}{9}\theta_1^2 &= 1 - \theta_1^2 \\ \Leftrightarrow \theta_1^2 &= \frac{9}{25} \\ \Leftrightarrow \theta_1 &= \pm \frac{3}{5} \\ \Leftrightarrow (\theta_1, \theta_2) &= (\pm \frac{3}{5}, \pm \frac{4}{5}) \end{aligned}$$

על ידי בדיקה של הערכים שמניבים ארבעת הנקודות האלה מוצאים כי המינימום הגלובלי של פונקציית המטרה מתקבל במקרה של  $(\theta_1, \theta_2) = (\frac{3}{5}, \frac{4}{5})$ .

## סקלרים, וקטורים מטריצות ונגזרות

במהלך הקורס אנו נתקל פעמים רבות בצורך לחשב נגזרות המערבות וקטורים ומטריצות. נזכיר / נסביר בקצרה של כיצד נזגרת אלו מחושבות. נתחיל במקרה המוכר של הגרדיאנט, בו אנו מבצעים גזירה של פונקציה סקלרית לפי וקטור. לאחר מכאן נראה כיצד הגדרה זו מורחבת גם למקרים נוספים בהם הפונקציה לא בהכרח סקלרית והגזירה היא לא בהכרח לפי וקטור.

בעבור פונקציה  $f(x)$  אשר מקבלת וקטור  $x$  באורך  $d$  ומחזירה סקלר, פעולת הגרדיאנט מוגדרת באופן הבא:

$$\nabla_x f(x) = \frac{d}{dx} f(x) = \begin{bmatrix} \frac{d}{dx_1} f(x) \\ \frac{d}{dx_2} f(x) \\ \vdots \\ \frac{d}{dx_d} f(x) \end{bmatrix}$$

לדוגמא:

$$\frac{d}{dx} (a^\top x) = \frac{d}{dx} \left( \sum_{i=1}^d a_i x_i \right) = \begin{bmatrix} \frac{d}{dx_1} \left( \sum a_i x_i \right) \\ \frac{d}{dx_2} \left( \sum a_i x_i \right) \\ \vdots \\ \frac{d}{dx_d} \left( \sum a_i x_i \right) \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_d \end{bmatrix} = a$$

נסתכל כעת על מקרה מעט יותר מורכב בו אנו רוצים לגזור פונקציה וקטורית  $f(x)$  אשר מקבלת וקטור  $x$  באורך  $d$  ומחזירה וקטור באורך  $n$ :

$$\frac{d}{dx} f(x) = \begin{bmatrix} \left. \frac{d}{dx_1} f(x) \right| & \left. \frac{d}{dx_2} f(x) \right| & \dots & \left. \frac{d}{dx_d} f(x) \right| \end{bmatrix}$$

שימו לב שהתוצאה של הגזירה הינה מטריצה בגודל  $n \times d$ . באופן כללי, הגדול תוצאת פעולת הגזירה יהיה תמיד הגודל של האובייקט שאותו גוזרים בתוספת הגודל של האיבר שלפיו אנו מבצעים את הגזירה. דוגמאות:

- תוצאת הגזירה של מטריצה בגודל  $n \times m$  לפי וקטור באורך  $d$  תהיה טנזור בגודל  $n \times m \times d$ .
- תוצאת הגזירה של סקלר לפי מטריצה בגודל  $n \times m$  תהיה מטריצה בגודל  $n \times m$ .
- תוצאת הגזירה של מטריצה בגודל  $n \times m$  לפי מטריצה בגודל  $o \times p$  תהיה טנזור בגודל  $n \times m \times o \times p$ .

למורת חשוב להבין כיצד נגזרות אלו מודרות ומחושבות, בפועל אנחנו כמעט ולא ניתקל בצורך לחשב אותם על פי ההגדרה. במרבית המקרים, הנגזרות בהם ניתקל יבואו מתוך קבוצה מצומצמת של נגזרות מוכרות (או הרכבה שלהם עם פונקציות אחרות) ונוכל להשתמש בתוצאות מוכרות ולחסוך עבודה מיותרת. ניתן למצוא רשימה של נגזרות בדף הנוסחאות של הקורס.

## תרגיל 1.2 - תרגיל בנגזרות

חשבו את הנגזרות הבאות:

1.  $\frac{d}{dx} \|x\|_2^2$
2.  $\frac{d}{dx} \|x\|_2$  הנחיה: השתמש בכלל השרשרת
3.  $\frac{d}{dx} (x^\top A x)$
4.  $\frac{d}{dA} (x^\top A x)$

## פתרון 1.2

### סעיף 1

$$\frac{d}{dx} \|x\|_2^2 = \frac{d}{dx} (x^\top x) = \frac{d}{dx} \left( \sum_{i=1}^d x_i^2 \right) = \begin{bmatrix} \frac{d}{dx_1} \left( \sum x_i^2 \right) \\ \frac{d}{dx_2} \left( \sum x_i^2 \right) \\ \vdots \\ \frac{d}{dx_d} \left( \sum x_i^2 \right) \end{bmatrix} = \begin{bmatrix} 2x_1 \\ 2x_2 \\ \vdots \\ 2x_d \end{bmatrix} = 2x$$

**סעיף 2** נסמן  $h(x) = \|x\|_2^2$  ונשים לב ש  $\|x\|_2 = \sqrt{h(x)}$ . כמו כן נשתמש בעובדה שאת הנגזרת של  $h(x)$  כבר חישבנו בסעיף הקודם:

$$\frac{d}{dx} \|x\|_2 = \frac{d}{dx} \sqrt{h(x)} = \frac{1}{2\sqrt{h(x)}} \cdot \frac{d}{dx} h(x) = \frac{x}{\sqrt{h(x)}}$$

### סעיף 3 נגזור על פי הגדרה:



$$\begin{aligned}
\frac{d}{dx}(x^\top Ax) &= \frac{d}{dx} \left( \sum_{i,j} a_{i,j} x_i x_j \right) = \begin{bmatrix} \frac{d}{dx_1} \left( \sum_{i,j} a_{i,j} x_i x_j \right) \\ \frac{d}{dx_2} \left( \sum_{i,j} a_{i,j} x_i x_j \right) \\ \vdots \\ \frac{d}{dx_d} \left( \sum_{i,j} a_{i,j} x_i x_j \right) \end{bmatrix} = \begin{bmatrix} \sum_i a_{i,1} x_i + \sum_j a_{1,j} x_j \\ \sum_i a_{i,2} x_i + \sum_j a_{2,j} x_j \\ \vdots \\ \sum_i a_{i,d} x_i + \sum_j a_{d,j} x_j \end{bmatrix} \\
&= \begin{bmatrix} A_{:,1}^\top x + A_{1,:} x \\ A_{:,2}^\top x + A_{2,:} x \\ \vdots \\ A_{:,d}^\top x + A_{d,:} x \end{bmatrix} = A^\top x + Ax = (A^\top + A)x
\end{aligned}$$

**סעיף 4** נגזור על פי הגדרה:

$$\begin{aligned}
\frac{d}{dA}(x^\top Ax) &= \frac{d}{dA} \left( \sum_{i,j} a_{i,j} x_i x_j \right) \\
&= \begin{bmatrix} \frac{d}{da_{1,1}} \left( \sum_{i,j} a_{i,j} x_i x_j \right) & \frac{d}{da_{1,2}} \left( \sum_{i,j} a_{i,j} x_i x_j \right) & \cdots & \frac{d}{da_{1,m}} \left( \sum_{i,j} a_{i,j} x_i x_j \right) \\ \frac{d}{da_{2,1}} \left( \sum_{i,j} a_{i,j} x_i x_j \right) & \frac{d}{da_{2,2}} \left( \sum_{i,j} a_{i,j} x_i x_j \right) & \cdots & \frac{d}{da_{2,m}} \left( \sum_{i,j} a_{i,j} x_i x_j \right) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{d}{da_{n,1}} \left( \sum_{i,j} a_{i,j} x_i x_j \right) & \frac{d}{da_{n,2}} \left( \sum_{i,j} a_{i,j} x_i x_j \right) & \cdots & \frac{d}{da_{n,m}} \left( \sum_{i,j} a_{i,j} x_i x_j \right) \end{bmatrix} \\
&= \begin{bmatrix} x_1 x_1 & x_1 x_2 & \cdots & x_1 x_m \\ x_2 x_1 & x_2 x_2 & \cdots & x_2 x_m \\ \vdots & \vdots & \ddots & \vdots \\ x_n x_1 & x_n x_2 & \cdots & x_n x_m \end{bmatrix} = xx^\top
\end{aligned}$$

### תרגיל 1.3 - בחזרה לבעיית המעגל

נחזור לפונקציית ה"ציון" של מידת ההתאמה של מעגל לנקודות מהדוגמא בתחילת התרגול.

$$f(c, r) = \frac{1}{n} \sum_{i=1}^n (\|x^{(i)} - c\|_2 - r)^2$$

חשבו את הנגזרות של פונקציה זו:  $\nabla_c f$  ו  $\frac{d}{dr} f$ . (לצורך החישוב של  $\nabla_c f$  השתמשו בשיטה דומה לזו שהופיע בסעיף 2 של השאלה הקודמת)

### פתרון 1.3

נתחיל בגזירה לפי  $r$ :

$$\frac{d}{dr} f(c, r) = \frac{2}{n} \sum_{i=1}^n (r - \|x^{(i)} - c\|_2) r$$

על פי ההנחיה, לצורך הגזירה לפי  $c$  נסמן את פונקציית העזר  $h(x, c) = \|x - c\|_2^2$  בעזרתה נוכל לרשום את  $f$  כ:

$$f(c, r) = \frac{1}{n} \sum_{i=1}^n (\sqrt{h(x^{(i)}, c)} - r)^2$$

נחשב תחילה את הנגזרת של  $h(x, c)$ :

$$\begin{aligned}\frac{d}{dc}h(x, c) &= \frac{d}{dc}\|x - c\|_2^2 \\ &= \frac{d}{dc}(x - c)^\top (x - c) \\ &= \frac{d}{dc}(\|x\|_2^2 - 2c^\top x + \|c\|_2^2) \\ &= 2(c - x)\end{aligned}$$

נשתמש כעת בתוצאה זו על מנת לחשב את הנגזרת הכוללת:

$$\begin{aligned}\frac{d}{dc}f(c, r) &= \frac{2}{n} \sum_{i=1}^n (\sqrt{h(x^{(i)}, c)} - r) \cdot \left( \frac{d}{dh(x^{(i)})} \sqrt{h(x^{(i)}, c)} \right) \cdot \frac{d}{dx}h(x^{(i)}) \\ &= \frac{2}{n} \sum_{i=1}^n (\sqrt{h(x^{(i)}, c)} - r) \frac{1}{\sqrt{h(x^{(i)}, c)}} (c - x^{(i)}) \\ &= \frac{2}{n} \sum_{i=1}^n (r - \|x^{(i)} - c\|_2) \frac{x^{(i)} - c}{\|x^{(i)} - c\|_2}\end{aligned}$$

## descent Gradient (שיטת הגרדיאנט)

כפי שצינו קודם, במקרים רבים לא נוכל פתור את בעיות האופטימיזציה על ידי גזירה והשוואה ל 0 ונאלץ לעשות שימוש בשיטות נומריות. נתאר כאן בקצרה את שיטת ה **descent gradient** (שיטת הגרדיאנט) שהיא אחת השיטות הבסיסיות ביותר אשר מנסה לפתור את בעיות מסוג זה. אנו עוד נדון בהמשך הקורס בהרחבה בתכונות ובבעיות הקיימות בשיטה זו, אך בשלב זה נסתפק בלתאר את אופן פעולתה.

הרעיון מאחורי שיטה זו הינו להתחיל בנקודה אקראית כל שהיא במרחב ולהתחיל לזוז בצעדים קטנים לכיוון שבו פונקציית המטרה קטנה באופן המהיר ביותר. הכיוון הזה הוא כמובן הכיוון ההפוך לגרדיאנט של הפונקציה. זהו אלגוריתם חמדן (greedy) אשר מנסה בכל צעד לשפר במעט את מצבו ביחס לשלב הקודם. אלגוריתמים מסוג זה מתכנסים לרוב למינימום לוקאלי ולא למינימום הגלובלי של הפונקציה. האלגוריתם זה רחוק מלתת מענה מושלם לבעיה, אך במקרים רבים הוא מצליח לספק פתרון סביר.

## האלגוריתם

- מאתחלים את  $\theta^{(0)}$  בנקודה אקראית כל שהיא
- חוזרים על צעד העדכון הבא עד להתכנסות:

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_{\theta} f(\theta^{(t)})$$

הפרמטר  $\eta$  אשר קובע את גודל הצעדים אשר נעשה בתהליך ההתכנסות. (את הדיון על קריטריון ההתכנסות ועל הבחירת של  $\eta$  נשאיר לשלב מאוחר יותר)

## תרגיל 1.4 - fitting circle for descent Gradient

רשמו את צעד העדכון של אלגוריתם descent gradient בעבור המקרה של התאמת העיגול.

## פתרון 1.4

על פי הנגזרת שחיבנו בתרגיל 1.3 נסיק כי צעד העדכון הוא:

$$r^{(t+1)} = r^{(t)} - \frac{2\eta}{n} \sum_{i=1}^n (r^{(t)} - \|x^{(i)} - c^{(t)}\|_2) r^{(t)}$$
$$c^{(t+1)} = c^{(t)} - \frac{2\eta}{n} \sum_{i=1}^n (r^{(t)} - \|x^{(i)} - c^{(t)}\|_2) \frac{x^{(i)} - c^{(t)}}{\|x^{(i)} - c^{(t)}\|_2}$$

## הרצה של האלגוריתם

למטה מוצג תהליך ההתכנסות של אלגוריתם הגרדיאנט בעבור  $\eta = 0.01$  ו 500 צעדים כאשר מתחילים את התהליך ממעגל היחידה:

