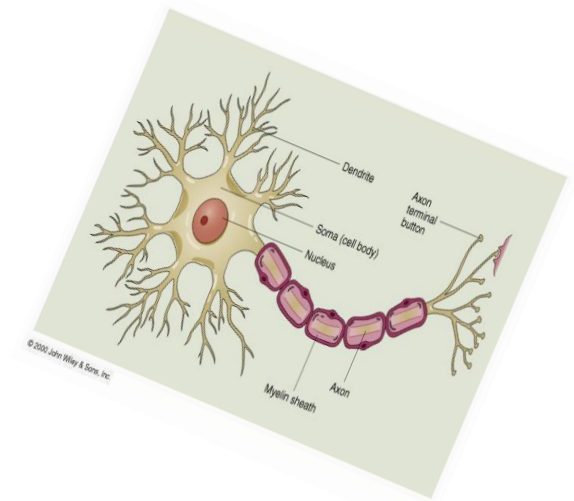
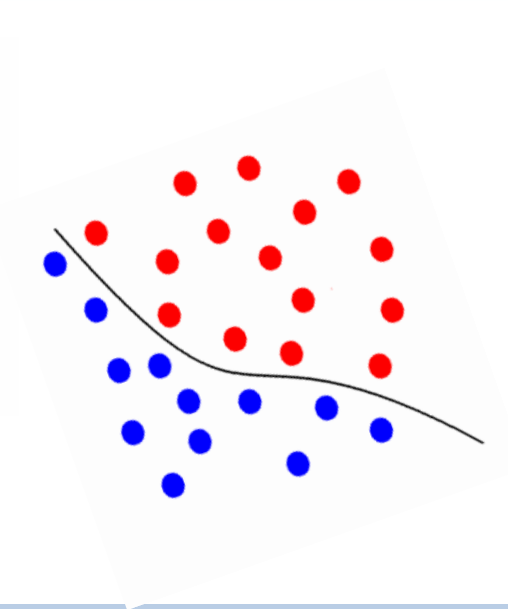
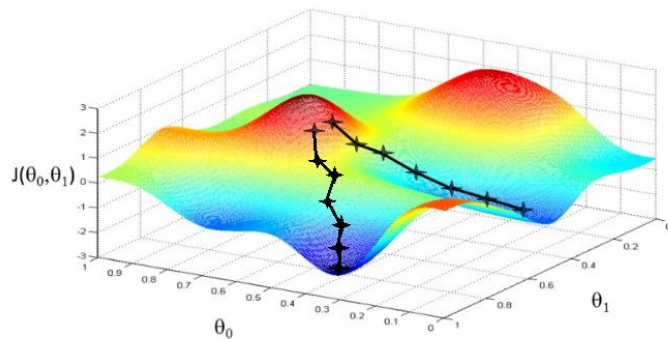




לא לשכוח להפעיל הקלטה!

Gradient descent



מה אנחנו צריכים כבר לדעת?

- ✓ Python – הגדרת משתנים, מטריצות, וקטורים, פונקציות מובנות, כתיבת פונקציות, גרפים
- ✓ מושגי יסוד בלמידת מכונה
- ✓ למידה מונחית Supervised learning
- ✓ למידה לא מונחית Unsupervised learning

Gradient descent

Gradient descent היא שיטה הנעזרת בגרדיאנט (מידע **מקומי**) על מנת למצוא את **המינימום** של פונקציה.

[קישור](https://towardsdatascience.com/a-visual-explanation-of-gradient-descent-methods-momentum-adagrad-rmsprop-adam-f898b102325c)

<https://towardsdatascience.com/a-visual-explanation-of-gradient-descent-methods-momentum-adagrad-rmsprop-adam-f898b102325c>



מעל"ה - מדע חישובי פיזיקה

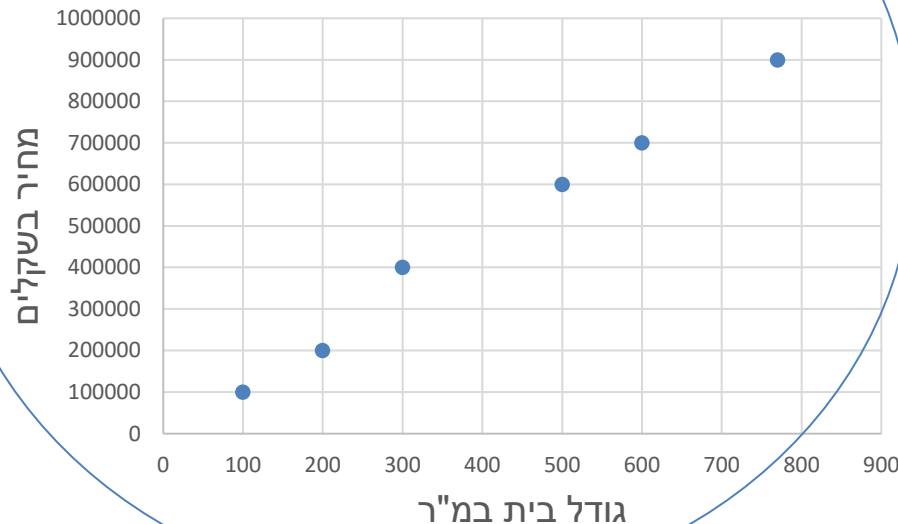
תזכורת - אלגוריתם בלמידת מכונה - למידה מונחית

למידה מונחית - מאפיינים

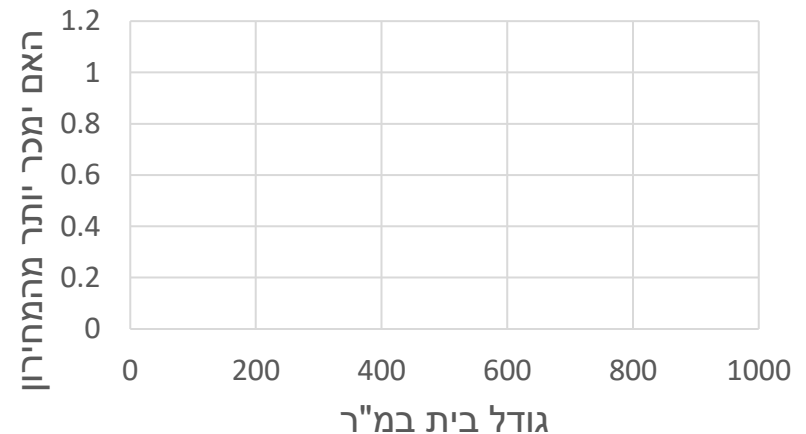
התשובה הנכונה ניתנת – Data Set

יש קשר בין משתני ה INPUT למשתנה ה OUTPUT

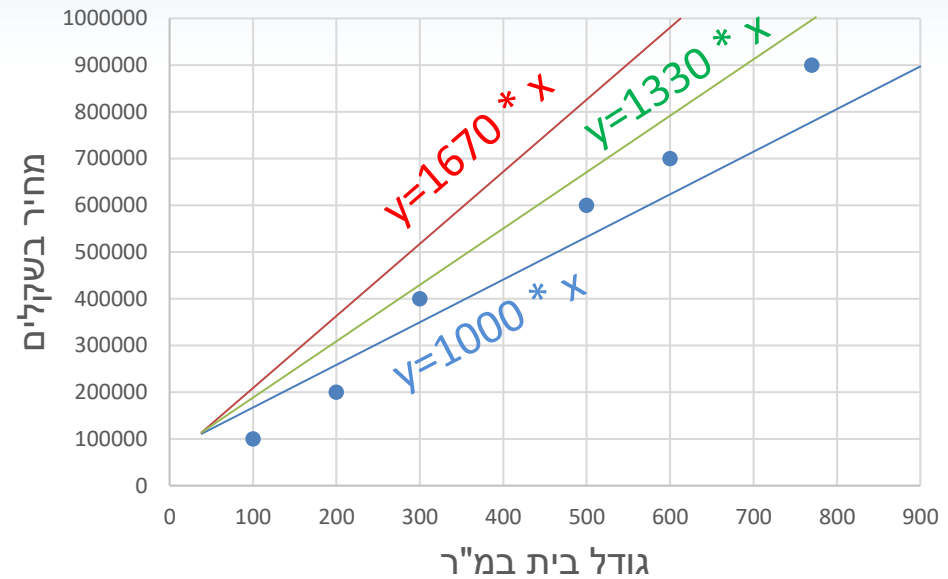
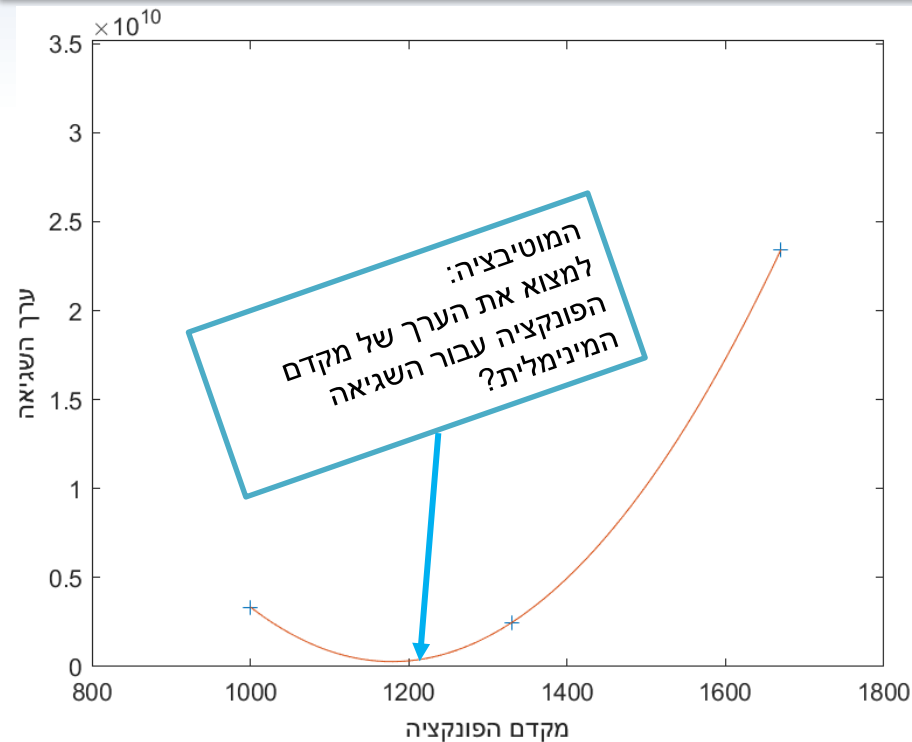
רגרסיה
Regression
פלט רציף



סיווג
classification
פלט - קטגוריות



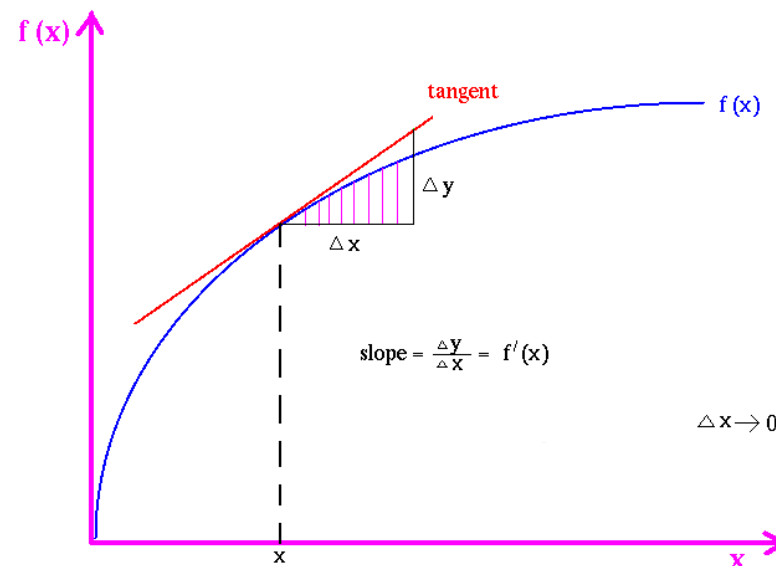
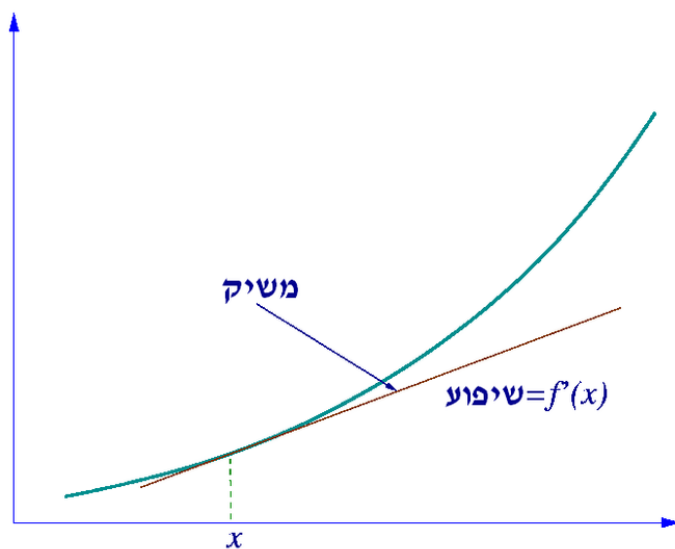
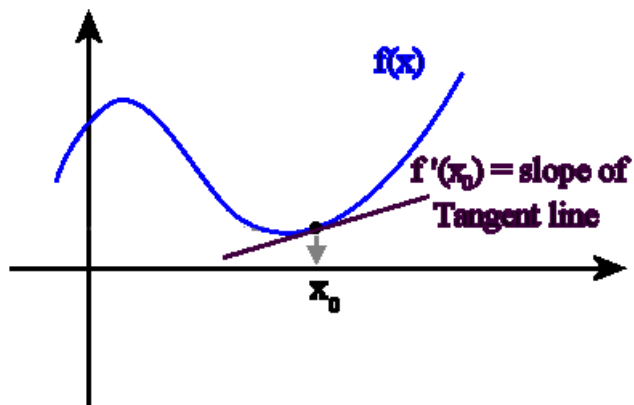
תזכורת – רגרסיה לינארית



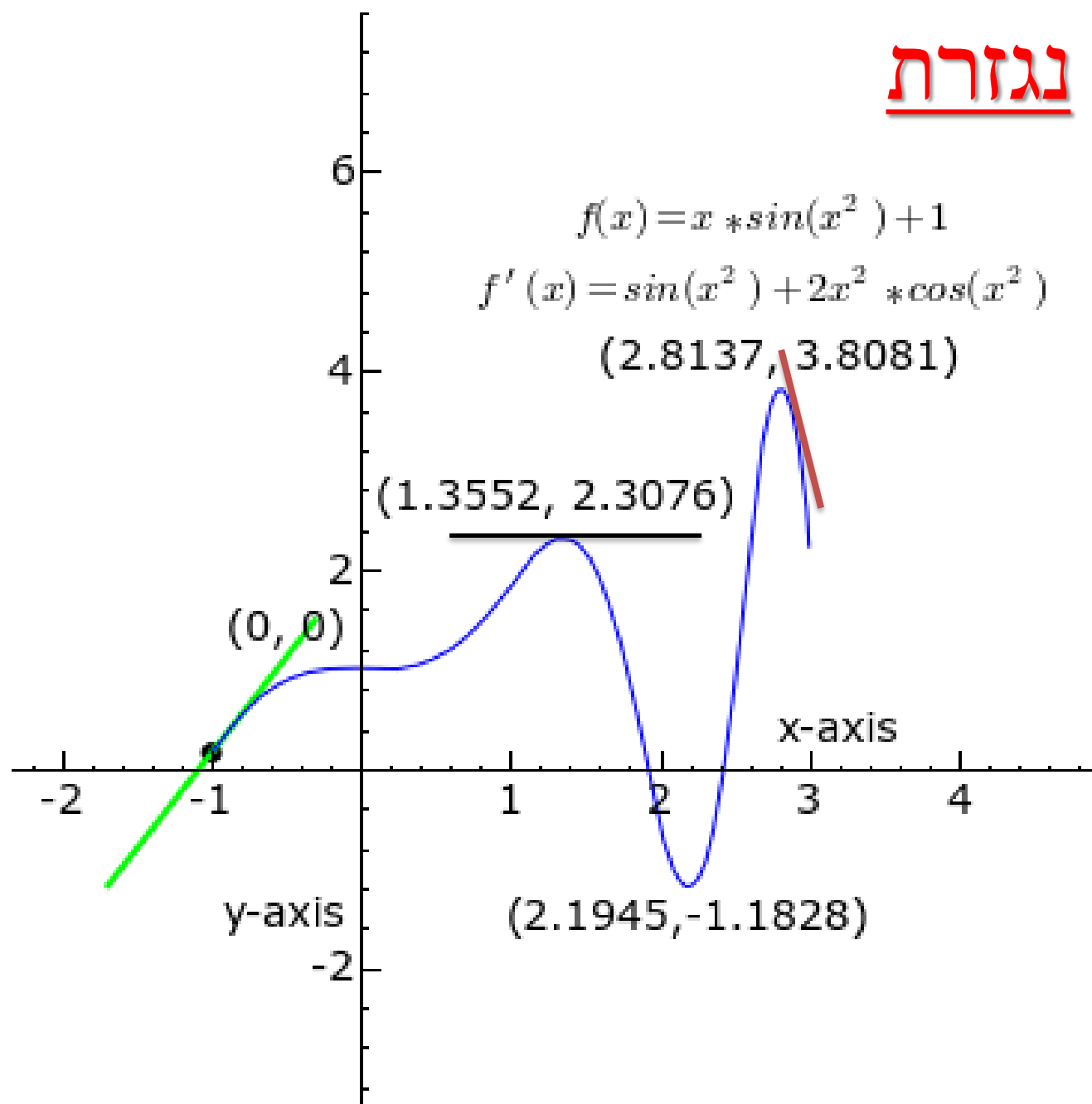
- ☐ מחשבים עבור כל קו את ערך הטעות
- ☐ ערך **הטעות** – סכימה של כל ריבועי ההפרשים בין הערך הרצוי למצוי לחלק ל 2
- ☐ משרטטים גרף של הטעות כפונקציה של המקדם: (1000,1330,1670)

Gradient descent – רקע מתמטי

• נגזרת



נגזרת

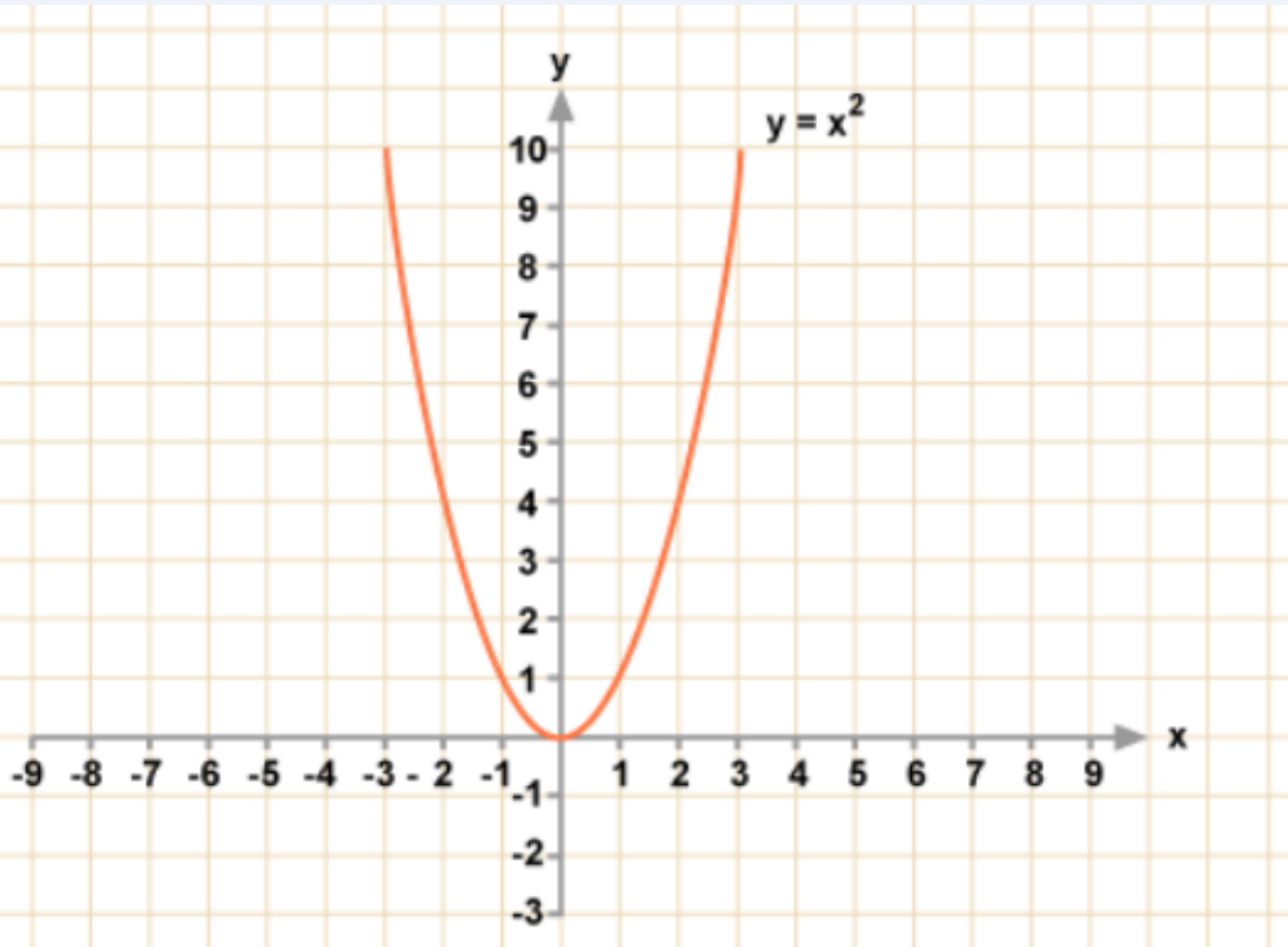


'נגזרת descent'

- על מנת להגיע למינימום של הפונקציה $f(x)$:
 - נתחיל מערך x רנדומלי.
 - עבור מספר צעדים מסוים:
- נחשב את הנגזרת של הפונקציה בנקודה זו.
- אם הנגזרת חיובית – נקטין את x (נוסיף ל- x את מינוס הנגזרת).
- אם הנגזרת שלילית – נגדיל את x (נוסיף ל- x את מינוס הנגזרת).
- כלומר – נלך בכיוון ההפוך לנגזרת.
- המטרה – למצוא את המשתנה x בנקודת המינימום של הפונקציה.

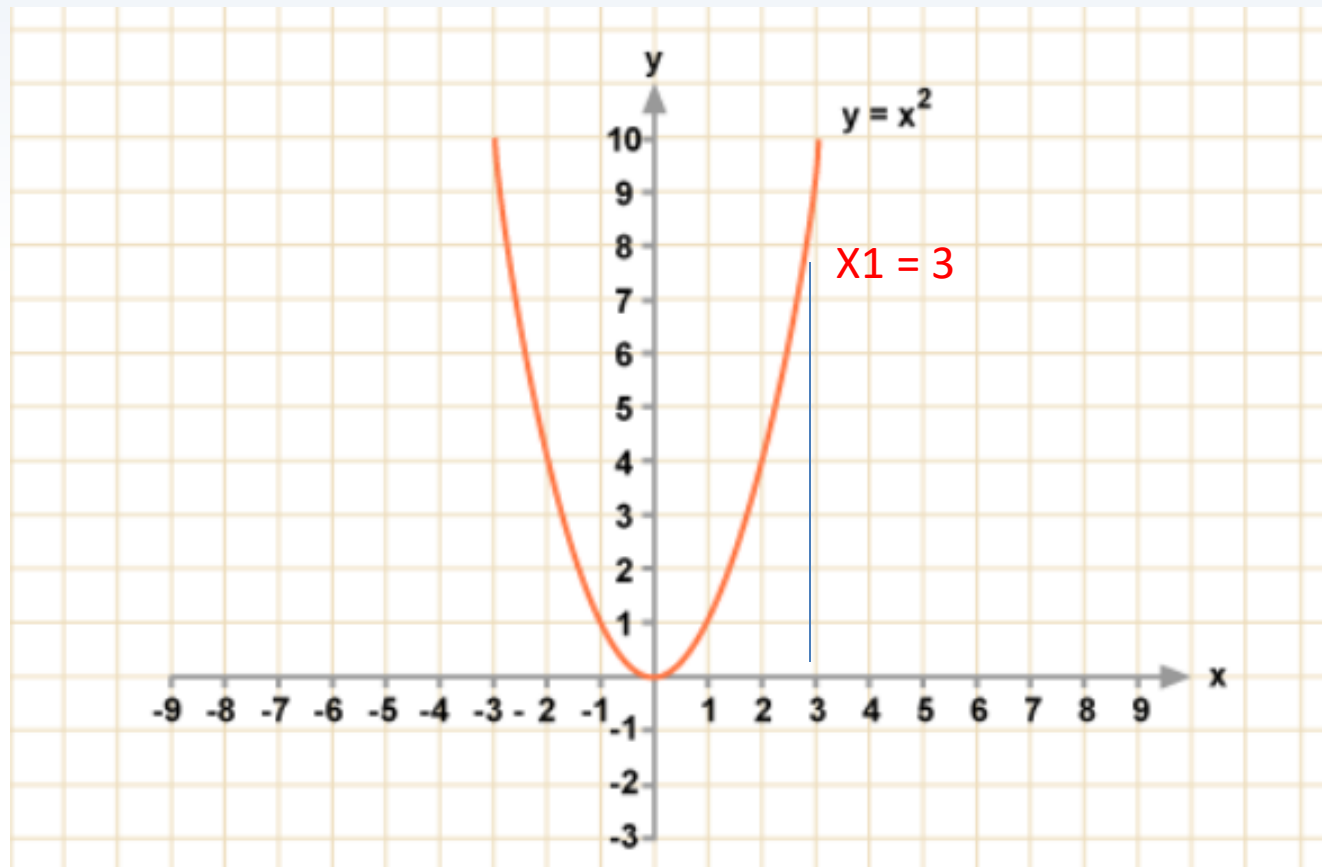
נתונה הפונקציה $y = x^2$

איך נמצא את המינימום בשיטת Gradient descent



- שלב 1 – איתחול משתנים:

נבחר נקודה רנדומלית, למשל $x_1 = 3$
נבחר את גודל הצעד (Learning Rate) $lr=0.01$



- שלב 2 - שינוי בערכי X עד שמגיעים לערך X בנקודת המינימום:

$$X(i+1) = X(i) - lr * dy/dx(i)$$

לפי החישוב:

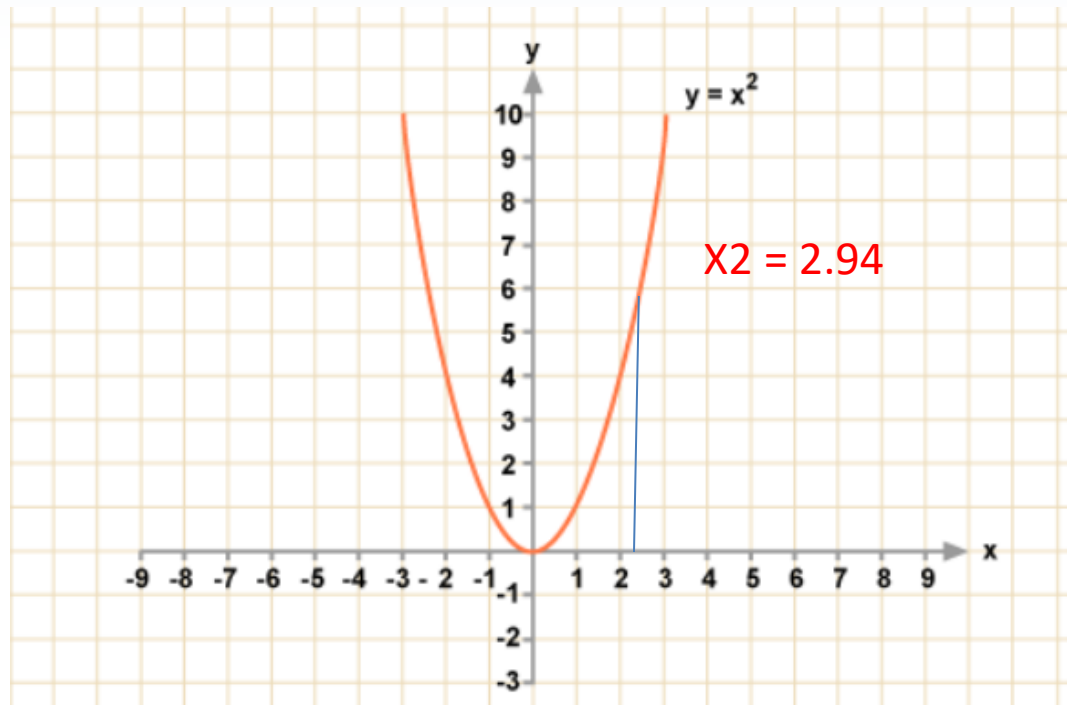
איטרציה #1

– שלב 2.1 - חישוב הנגזרת עבור X_1

$$dy/dx = 2 * X_1 = 2 * 3 = 6$$

– שלב 2.1 - חישוב הערך של X_2

$$X_2 = X_1 - 0.01 * 6 = 3 - 0.01 * 6 = 2.94$$



- שלב 2 - שינוי בערכי X עד שמגיעים לערך X בנקודת המינימום:

$$X(i+1) = X(i) - lr * dy/dx(i)$$

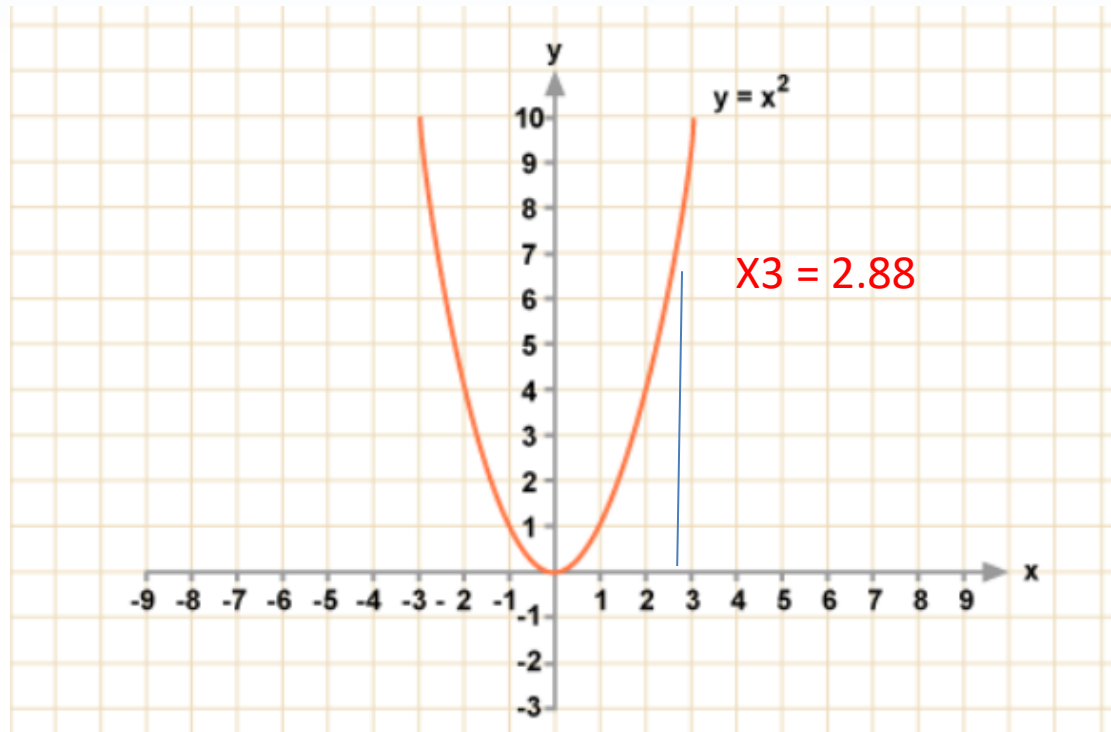
איטרציה #2

– שלב 2.1 - חישוב הנגזרת עבור X_2

$$dy/dx = 2 * X_2 = 2 * 2.94 = 5.88$$

– שלב 2.1 - חישוב הערך של X_2

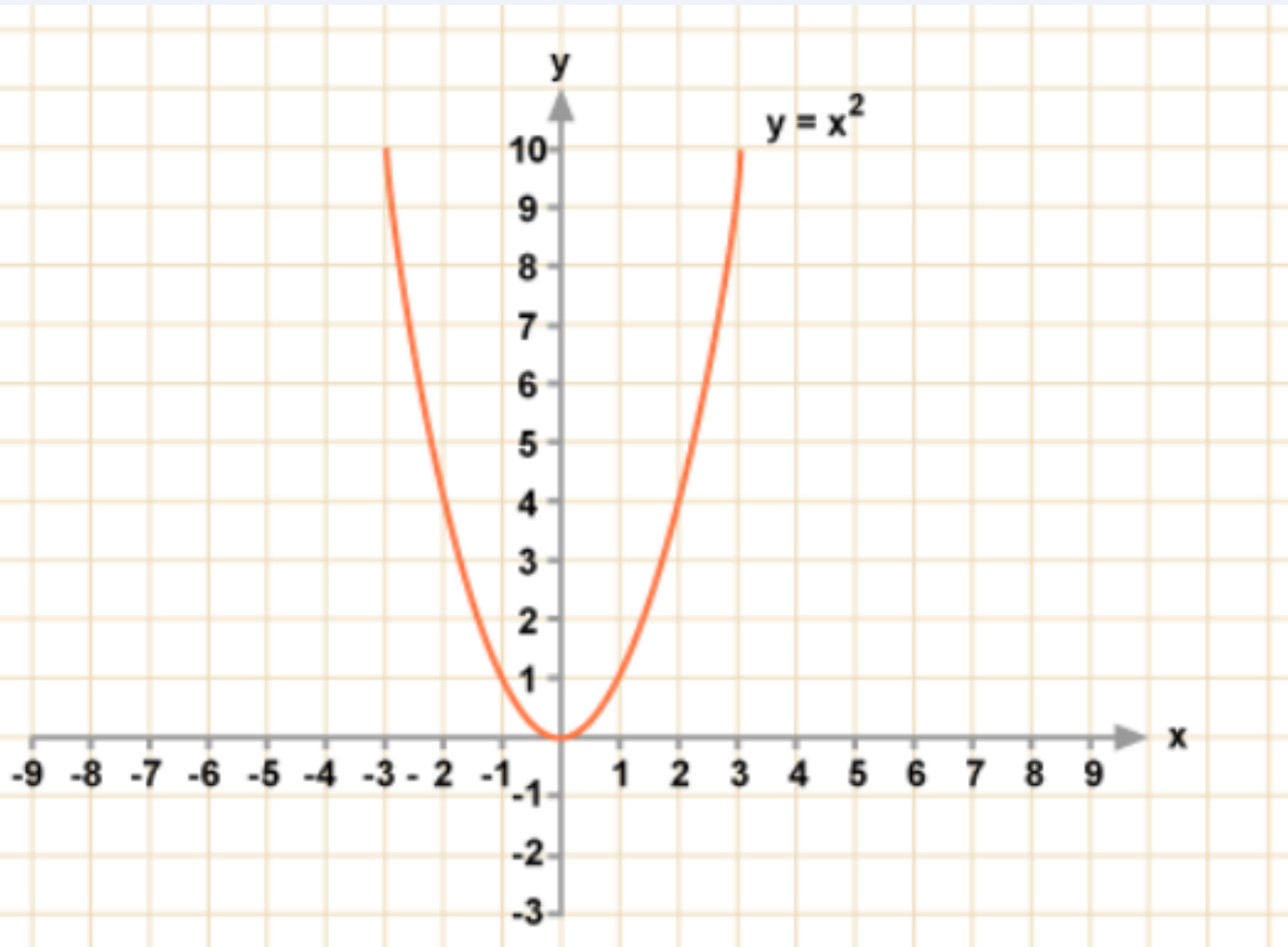
$$X_3 = X_2 - 0.01 * 5.88 = 2.94 - 0.01 * 5.88 = 2.88$$



נתונה הפונקציה $y = x^2$

עד מתי נבצע את האיטרציות?

על מה משפיע משתנה ה Learning Rate?



מקסימום האיטרציות 40

$lr=0.1$

נקודה התחלתית $x_1 = -5$

א. מצא את נקודת המינימום עבור

הפונקציה: $f(x) = x^2 - 10x + 5$

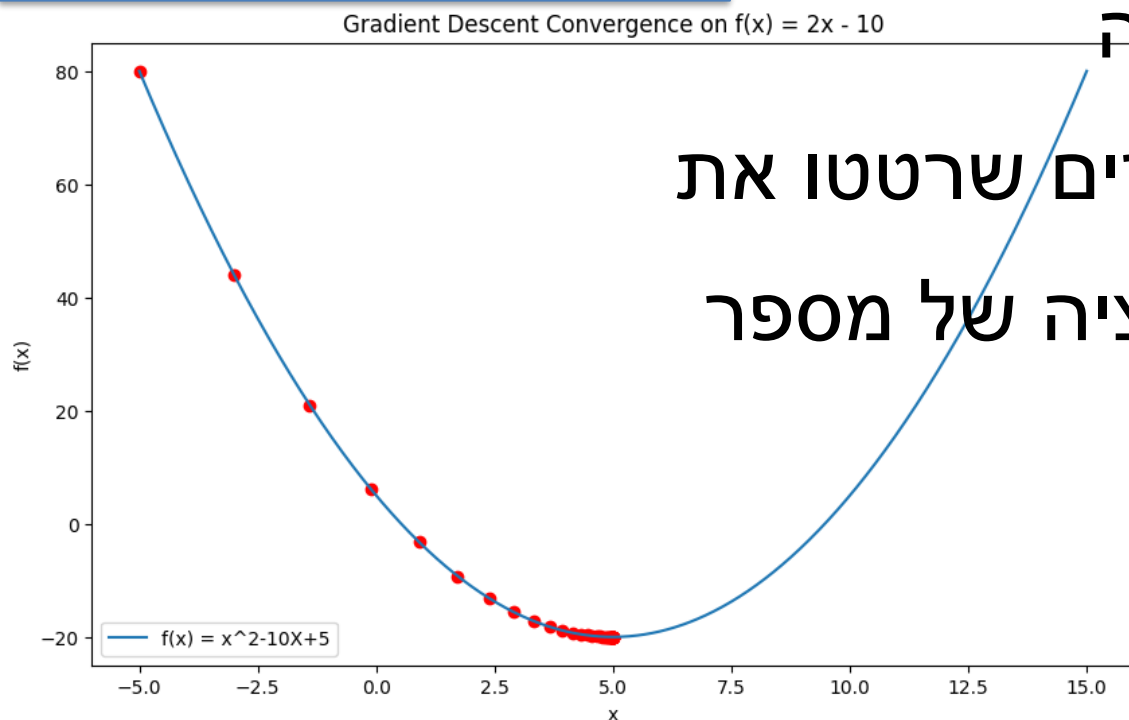
הנגזרת - $2x - 10$

ב. שרטטו את הפונקציה

ג. על אותה מערכת צירים שרטטו את

ההשתנות של x כפונקציה של מספר

האיטרציות



שאלות:

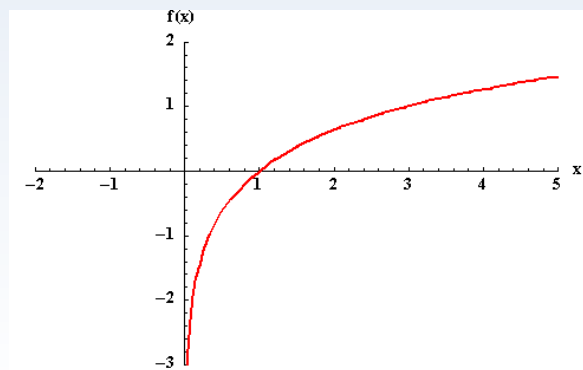
- ☐ מהי נקודת המינימום שהתקבלה? _____
- ☐ האם היא נקודת המינימום הצפויה? _____
- ☐ מה קורה לערכי המשתנה x מהרגע שהגעתם לנקודת המינימום? (קטנים/גדלים/לא משתנים) _____
- _____ מדוע?
- ☐ מה קרה כאשר קצב הלמידה היה קטן? האם הגעתם לנקודת המינימום של הפונקציה? מדוע? האם ניתן להגיע עם קצב למידה זה למינימום של הפונקציה? _____
- ☐ מה קרה כאשר קצב הלמידה היה גדול? האם הגעתם לנקודת המינימום של הפונקציה? מדוע? האם ניתן להגיע עם קצב למידה זה למינימום של הפונקציה? _____
- _____

שאלות:

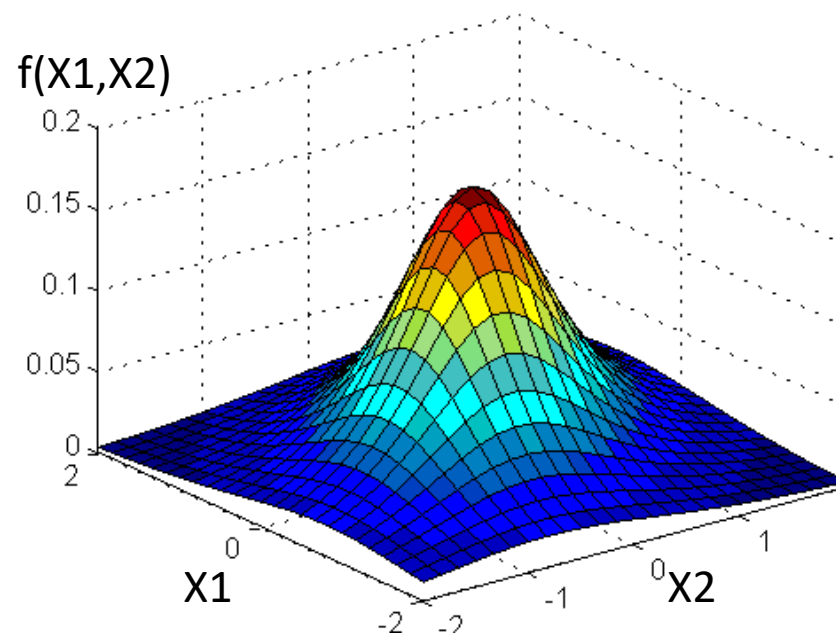
- מהי נקודת המינימום שהתקבלה? 5
- האם היא נקודת המינימום הצפויה? 5
- מה קורה לערכי המשתנה x מהרגע שהגעתם לנקודת המינימום? (קטנים/גדלים/ לא משתנים)
מדוע? הנגזרת בנקודה הזו שואפת לאפס
- שנו את קצב הלמידה: $\text{learning rate}=0.001$. האם הגעתם לנקודת המינימום של הפונקציה? מדוע? האם ניתן להגיע עם קצב למידה זה למינימום של הפונקציה? לא מגיעים
למינימום קצב הלמידה קטן מדי למספר האיטרציות
- שנו את קצב הלמידה: $\text{learning rate}=1$. האם הגעתם לנקודת המינימום של הפונקציה? מדוע? האם ניתן להגיע עם קצב למידה זה למינימום של הפונקציה? לא מגיעים
למינימום קצב הלמידה גדול, כל איטרציה מדלגת על מיקום המינימום

פונקציה מרובת-משתנים

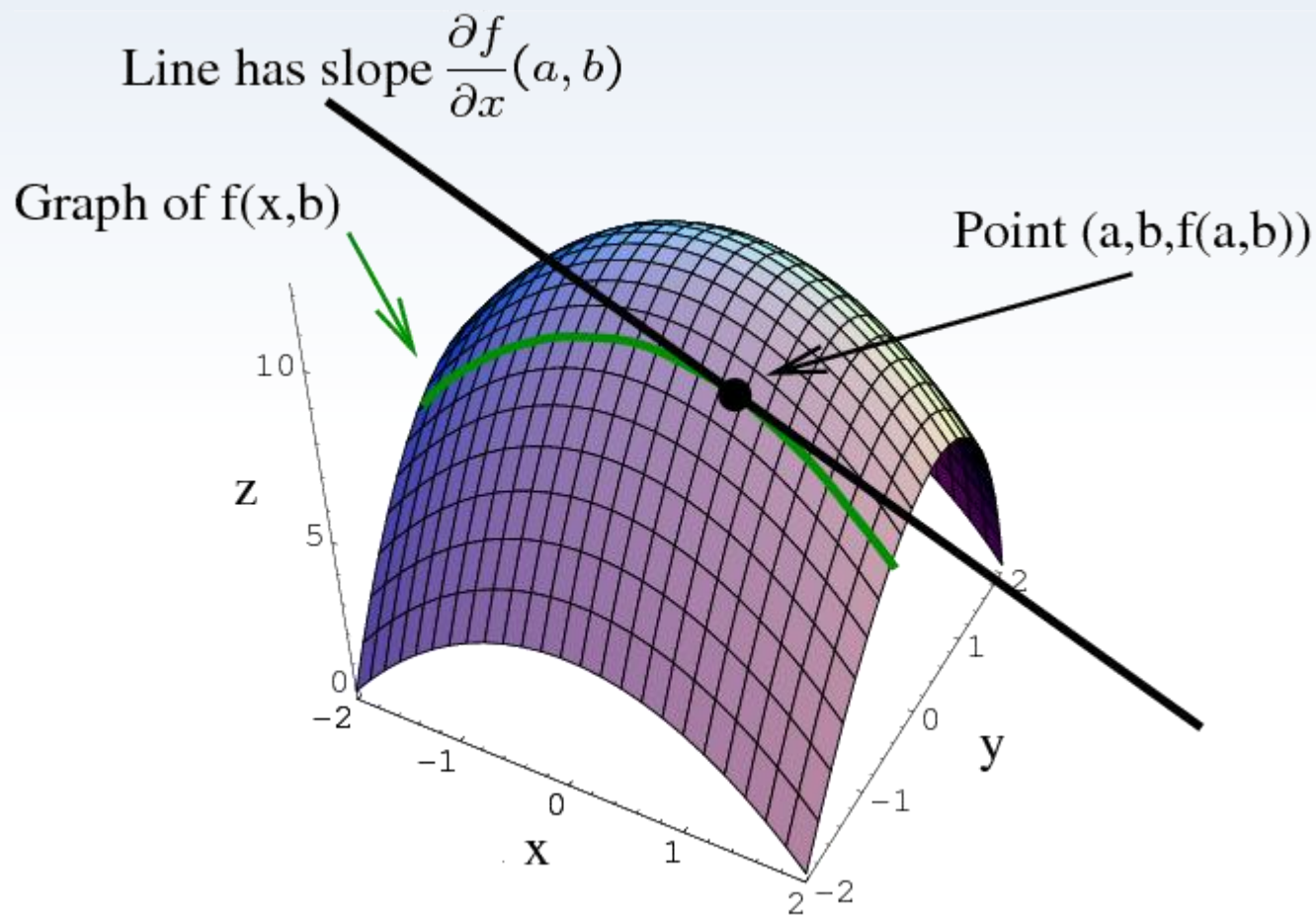
$f(x)$



$f(x_1, \dots, x_n)$



נגזרת חלקית



נגזרת חלקית

- נגזרת חלקית - גזירת פונקציה מרובת משתנים לפי אחד המשתנים.
- המשתנים האחרים נחשבים קבועים עבור הגזירה.

$$z = f(x, y) = x^2 + xy + y^2$$

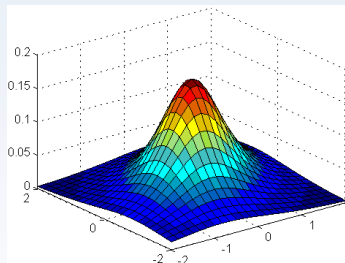
$$\frac{\partial z}{\partial x} = 2x + y$$

גרדיאנט

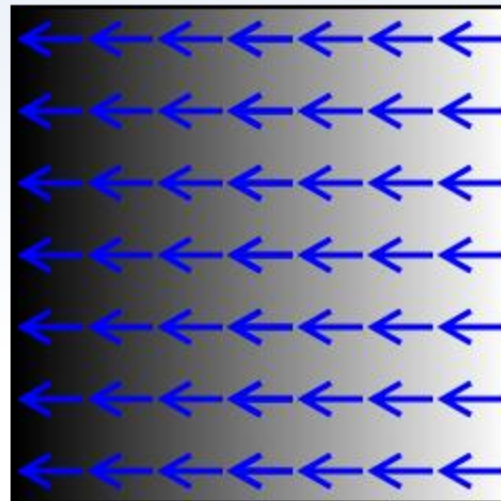
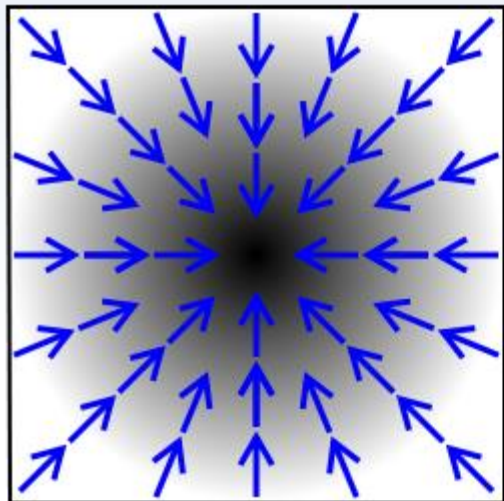
- אם f היא פונקציה גזירה של n משתנים, הגרדיאנט הוא וקטור באורך n של הנגזרות החלקיות של f .

$$\nabla f(x_1, \dots, x_n) = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)$$

גרדיאנט - משמעות גיאומטרית



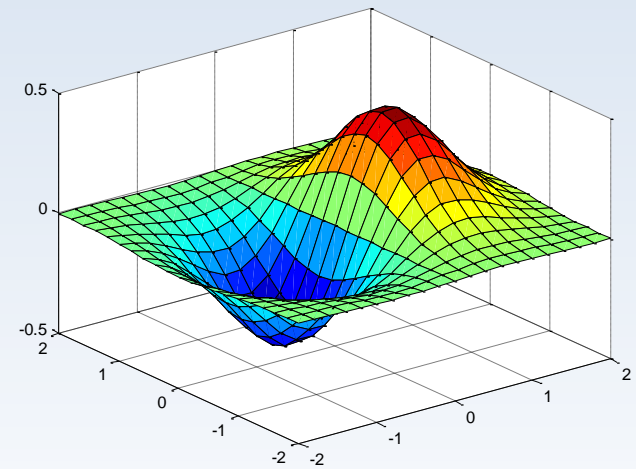
x2



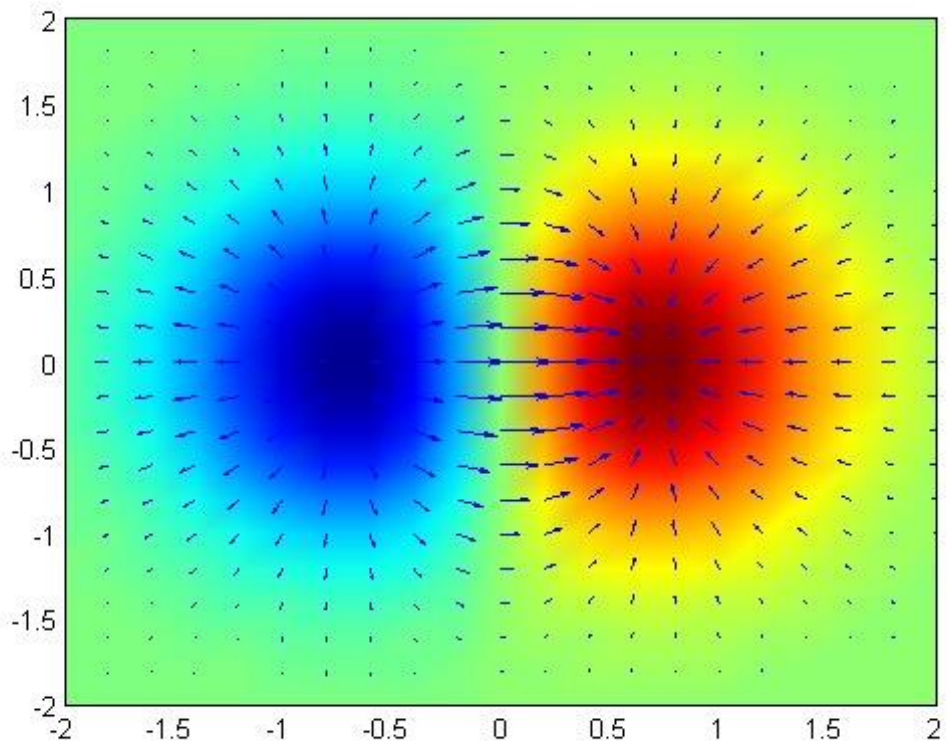
x1

ניתן לחשב את ערך הגרדיאנט של הפונקציה בכל נקודה ונקודה (עבור כל ערך של משתנים).

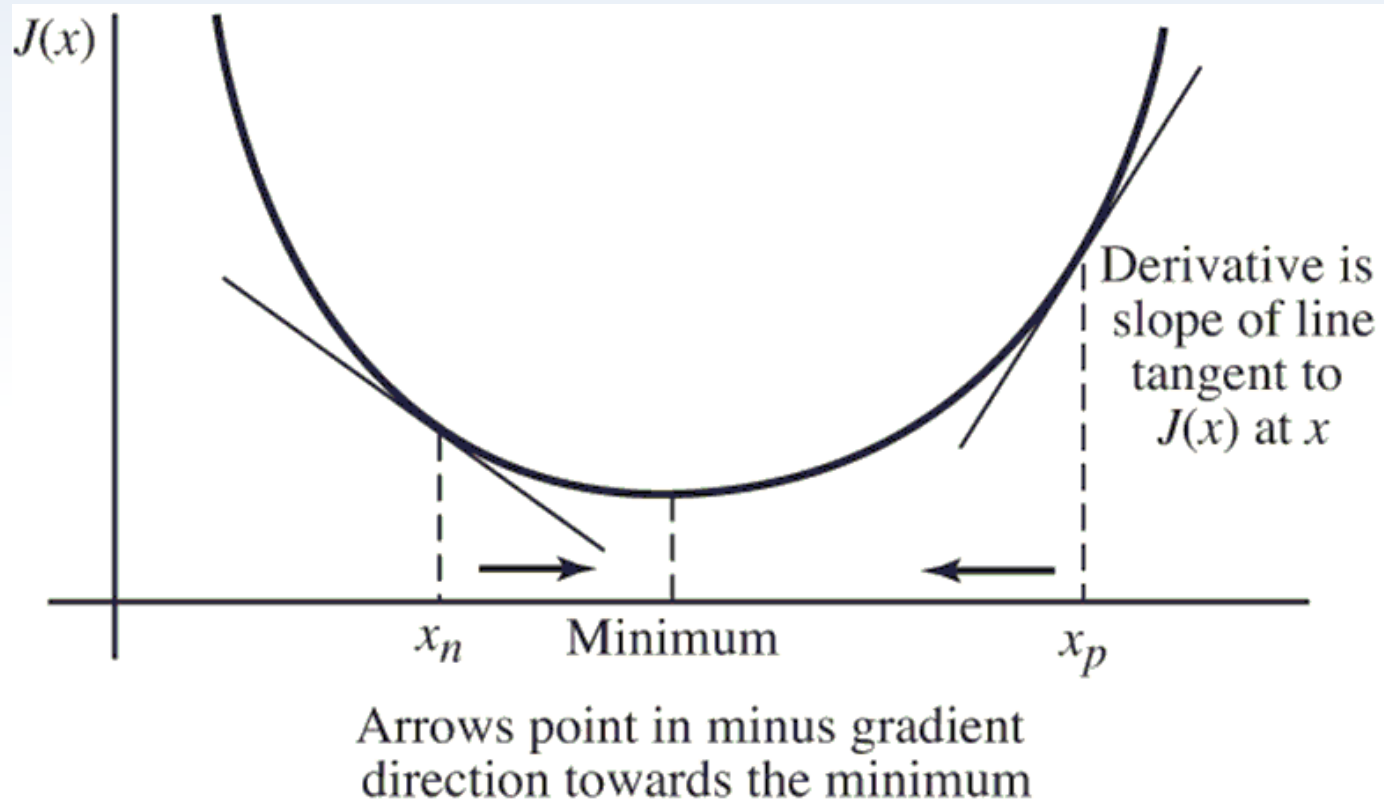
גרדיאנט



- הגרדיאנט מצביע לכיוון הגידול המקסימלי של הפונקציה.
- הליכה בכיוון הגרדיאנט תוביל להגדלת ערכי הפונקציה. הליכה נגד כיוון הגרדיאנט תוביל להקטנת ערכי הפונקציה.
- מעשית – עבור הפונקציה מרובת המשתנים, נשנה כל משתנה בהתאם לנגזרת החלקית שלו.



Gradient descent



כלל העדכון של המשקלות לפי gradient descent

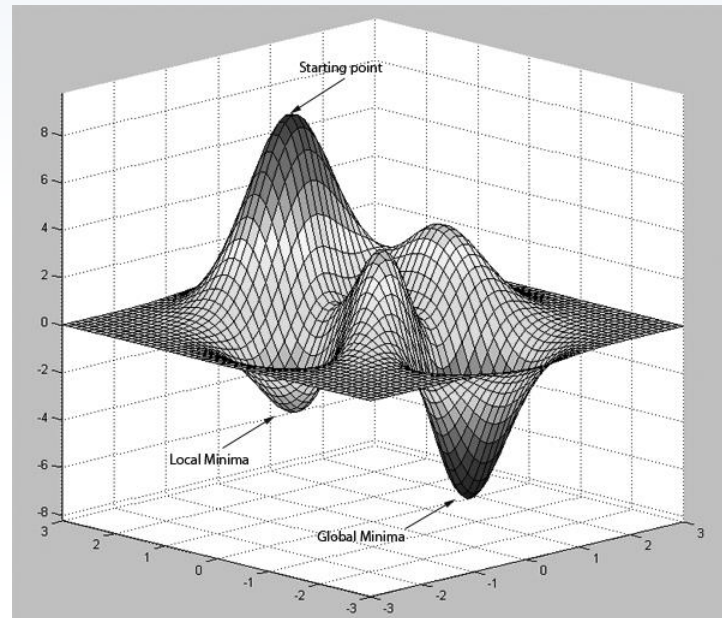
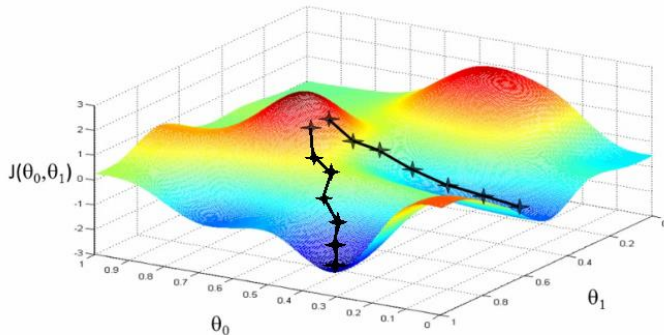
- על מנת לרדת בכיוון ההפוך לגרדיאנט, נעדכן כל משתנה בגודל השווה למינוס הנגזרת החלקית של פונקציית השגיאה ביחס למשקולת זו:

$$\Delta x_i = \varepsilon * -\frac{\partial f}{\partial x_i} \quad \Delta \vec{x} = \varepsilon * -gradient$$

- גודל הצעד בכיוון ההפוך לכיוון הגרדיאנט נקבע ע"י קצב הלמידה (**learning rate - ε**).
- העדכון יתבצע בכל המשתנים **בבת אחת** (לא נעדכן משתנה אחד ואז נחשב שוב את הגרדיאנט ונעדכן את השני וכו', אלא באותה נקודה נחשב את הגרדיאנט עבור כל המשתנים ונעדכן את כולם).
- **נגדיר מספר מקסימלי של צעדים.**
- נסו לכתוב למה שווה וקטור המשתנים החדש בכתיבה וקטורית.

Gradient descent

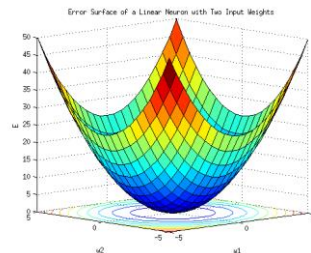
- קיימת סכנה להתכנסות למינימום לוקאלי.



- מהו הגרדיאנט בנקודת מינימום (לוקאלי או גלובאלי?).
כיצד יתקדם האלגוריתם?

חשיבות קצב הלמידה

- אפילו אם קיים רק מינימום יחיד במשטח השגיאה, אנחנו עלולים לפספס אותו אם נקבע קצב למידה גבוה מדי.
- לעומת זאת, אם נקבע קצב למידה נמוך מדי, מספר האיטרציות (הצעדים) הדרוש בשביל להגיע למינימום עלול להיות גבוה מאוד (למשל, גבוה יותר ממספר האיטרציות המקסימלי שהגדרנו).



Gradient descent – סיכום מעשי

- על מנת להגיע למינימום של הפונקציה $f(x(1), \dots, x(n))$:
 - נתחיל מנקודה מסוימת, שהיא וקטור המשתנים, המכיל את הערכים הרנדומלים של $x(1), \dots, x(n)$.
 - עבור מספר צעדים מסוים שהגדרנו:
 - נחשב את וקטור הגרדיאנט של הפונקציה בנקודה זו.
 - נוסיף לוקטור המשתנים את מינוס וקטור הגרדיאנט (נעדכן את כל המשתנים בבת אחת).
 - נפסיק כאשר אחד מהתנאים הבאים התממש:
 - הגענו למספר הצעדים המקסימלי.
 - הגענו לערך של הפונקציה f המקובל עלינו.

– מהי תשובת האלגוריתם?