

**למידה עמוקה 046211**  
**גיליון יבש 1**

**315681593**  
**208731422**

**גל גרנות**  
**ניר טבת**

## שאלה 1

נתחיל עם הביטוי משמאל ונשתמש ביצוג של הפונקציה כאינטגרל:

$$\begin{aligned} & f(w_1) - f(w_2) - \nabla f(w_2)^T (w_1 - w_2) \\ &= \int_0^1 \nabla f(w_2 + t(w_1 - w_2))^T (w_1 - w_2) dt - \nabla f(w_2)^T (w_1 - w_2) \end{aligned}$$

הביטוי השני קבוע ביחס לאינטגרל, ולכן:

$$\begin{aligned} &= \int_0^1 \left( \nabla f(w_2 + t(w_1 - w_2))^T (w_1 - w_2) - \nabla f(w_2)^T (w_1 - w_2) \right) dt \\ &= \int_0^1 (\nabla f(w_2 + t(w_1 - w_2)) - \nabla f(w_2))^T (w_1 - w_2) dt \\ &\leq \left| \int_0^1 (\nabla f(w_2 + t(w_1 - w_2)) - \nabla f(w_2))^T (w_1 - w_2) dt \right| \\ &\leq \int_0^1 |(\nabla f(w_2 + t(w_1 - w_2)) - \nabla f(w_2))^T (w_1 - w_2)| dt \end{aligned}$$

כאשר המעבר האחרון מתבצע לפי אי שוויון המשולש האינטגרלי. נשתמש בחלקות  $\beta$  על הגורם הראשון במכפלה באינטגרנד:

$$\begin{aligned} &\leq \int_0^1 \beta |w_2 + t(w_1 - w_2) - w_2| |w_1 - w_2| dt = |w_1 - w_2| \int_0^1 \beta t |w_1 - w_2| dt \\ &= \beta |w_1 - w_2|^2 \int_0^1 t dt = \frac{\beta}{2} |w_1 - w_2|^2 \end{aligned}$$

כאשר השתמשנו במשפט קושי-שוורץ ביחס למכפלה פנימית אינטגרלית בין הפונקציות המתאימות.

## שאלה 2

### סעיף א':

הפונקציה  $f$  היא תבנית בילינארית עם מטריצה positive-definite, ולכן המינימום שלה ביחס ל- $w$  הוא וקטור האפס. נמצא את הפרמטרים  $\eta, A$ , עבורם האלגוריתם מתכנס בצעד אחד:

$$w(t+1) - w(t) = -\eta A \nabla f(w(t)) = -\eta A \nabla \frac{1}{2} w^T H w = -\eta A H w$$

עבור  $A = H^{-1}$  (הקיימת מכיוון ש- $H$  היא positive definite) נקבל:

$$w(t+1) = w(t) - \eta H^{-1} H w(t) = w(t) - \eta w(t) = (1 - \eta)w(t) = 0 \Rightarrow \eta = 1$$

כלומר, נבחר  $A = H^{-1}$  ו- $\eta = 1$  לקבלת התכנסות בצעד בודד. נשים לב שעבור  $H_{d \times d}$ , חישוב  $H^{-1}$  הוא אלגוריתם בסיבוכיות  $\theta(d^{2.373})$ , ולכן אינו פרקטי ב- $d$  גדול.

מה הקירוב המתאים??

### סעיף ב':

נשתמש בהגדרות ה-GD וברמז:

$$w(t+1) = w(t) - \eta A \nabla f(w(t))$$

$$\begin{aligned} \frac{\partial f}{\partial t} &= \frac{\partial f}{\partial w} \frac{\partial w}{\partial t} = \nabla f(w(t))^T \dot{w}(t) = \nabla f(w(t))^T (-A \nabla f(w(t))) = \\ &= -\nabla f(w(t))^T A \nabla f(w(t)) \end{aligned}$$

נשתמש באי השוויון ברמז עם  $v = \nabla f(w(t)) \in \mathbb{R}^d$  ונקבל כי מתקיים:

$$0 \leq \lambda_{\min} \|\nabla f(w(t))\| \leq -\frac{\partial f}{\partial t} = \nabla f(w(t))^T A \nabla f(w(t)) \leq \lambda_{\max} \|\nabla f(w(t))\|$$

כלומר, הגרדיאנט של  $f$  אי-חיובי ולכן  $f$  לא עולה. כמו כן היא לא מתבדרת לפי הנתון, ולכן מתכנסת לנקודה קריטית.

## סעיף ג':

מאותם שיקולים נקבל:

$$\begin{aligned}
 f(w(t+1)) &\leq f(w(t)) + (w(t+1) - w(t))^T \nabla f(w(t)) + \frac{\beta}{2} \|w(t+1) - w(t)\|^2 \\
 &= -\eta \left( A \nabla f(w(t)) \right)^T \nabla f(w(t)) + \frac{\beta}{2} \left\| -\eta A \nabla f(w(t)) \right\|^2 \\
 &= f(w(t)) - \eta A \left\| \nabla f(w(t)) \right\|^2 + \frac{\beta \eta^2}{2} \left\| \nabla f(w(t)) \right\|^2 \\
 &= f(w(t)) - \eta \left( A - \frac{\beta \eta}{2} I \right) \left\| \nabla f(w(t)) \right\|^2
 \end{aligned}$$

נסמן  $A - \frac{\beta \eta}{2} I = C$  ונקבל:

$$f(w(t+1)) - f(w(t)) \leq -C \left\| \nabla f(w(t)) \right\|^2 \Rightarrow C \left\| \nabla f(w(t)) \right\|^2 \leq f(w(t)) - f(w(t+1))$$

התכונה מתקיימת לכל  $t$ . לכן, נבחר  $T \in \mathbb{N}$  כלשהו ונקבל:

$$\begin{aligned}
 C \sum_{t=0}^T \left\| \nabla f(w(t)) \right\|^2 &\leq \sum_{t=0}^T \left( f(w(t)) - f(w(t+1)) \right) = f(w(0)) - f(w(T+1)) \\
 &\leq f(w(0)) - \min_w f(w)
 \end{aligned}$$

כאשר המעבר לפני האחרון מכיוון שהטור המתקבל באגף ימין הוא טלסקופי. ניקח  $T \rightarrow \infty$  ונקבל:

$$0 \leq C \sum_{t=0}^{\infty} \left\| \nabla f(w(t)) \right\|^2 \leq f(w(0)) - \min_w f(w) < \infty$$

טור פונקציות מהצורה  $\sum_{t=0}^{\infty} \left\| \nabla f(w(t)) \right\|^2$  לא מתכנס אם האיבר הכללי שלו לא שואף ל-0, ולכן נסיק:

$$\lim_{t \rightarrow \infty} \left\| \nabla f(w(t)) \right\|^2 = 0 = \lim_{t \rightarrow \infty} C \nabla f(w(t)) = 0$$

כלומר, ב- $t \rightarrow \infty$  האלגוריתם מתכנס לנקודה סטציונרית כנדרש.

### שאלה 3

#### סעיף א':

ראשית מכלל השרשרת:

$$\frac{\partial L(w(t-1))}{\partial \eta_{t-1}} = \frac{\partial L}{\partial w(t-1)} \cdot \frac{\partial w(t-1)}{\partial \eta_{t-1}} = (**)$$

עבור האיבר השמאלי:

$$\frac{\partial L}{\partial w(t-1)} = \nabla L(w(t-1))$$

ועבור האיבר הימני:

$$w(t-1) = w(t-2) - \eta_{t-1} \nabla L(w(t-2))$$

$$\frac{\partial w(t-1)}{\partial \eta_{t-1}} = -\nabla L(w(t-2))$$

וסך הכל קיבלנו:

$$(**) = \nabla L(w(t-1)) \cdot (-\nabla L(w(t-2))) = -\nabla L(w(t-1))^T \nabla L(w(t-2))$$

#### סעיף ב':

באופן דומה לסעיף א', נשתמש בכלל השרשרת:

$$\frac{\partial L(w(t-1))}{\partial \alpha_{t-1}} = \frac{\partial L(w(t-1))}{\partial \eta_{t-1}} \cdot \frac{\partial \eta_{t-1}}{\partial \alpha_{t-1}} = (**)$$

האיבר השמאלי הוא כמובן התוצאה מסעיף א'. נחשב את האיבר הימני:

$$\eta_{t-1} = \eta_{t-2} - \alpha_{t-1} \cdot \frac{\partial L(w(t-2))}{\partial \eta_{t-2}}$$
$$\frac{\partial \eta_{t-1}}{\partial \alpha_{t-1}} = -\frac{\partial L(w(t-2))}{\partial \eta_{t-2}} = \nabla L(w(t-2))^T \nabla L(w(t-3))$$

כאשר המעבר האחרון הוא הזזה של תוצאת סעיף א' ב-1 והכפלה במינוס.

סך הכל קיבלנו:

$$(**) = (-\nabla L(w(t-1))^T \nabla L(w(t-2))) \cdot (\nabla L(w(t-2))^T \nabla L(w(t-3)))$$

## סעיף ג':

נבצע כעת את כלל השרשרת הבא:

$$\frac{\partial L(w(t-1))}{\partial \eta_{t-2}} = \frac{\partial L}{\partial w(t-1)} \cdot \frac{\partial w(t-1)}{\partial \eta_{t-1}} = \frac{\partial L}{\partial w(t-1)} \cdot \frac{\partial w(t-1)}{\partial w(t-2)} \cdot \frac{\partial w(t-2)}{\partial \eta_{t-2}}$$

כעת האיבר הראשון הוא:

$$\frac{\partial L}{\partial w(t-1)} = \nabla L(w(t-1))$$

האיבר הימני, זהה לאיבר הימני מסעיף א' בהזזה של 1:

$$\frac{\partial w(t-1)}{\partial \eta_{t-1}} = -\nabla L(w(t-2)) \rightarrow \frac{\partial w(t-2)}{\partial \eta_{t-2}} = -\nabla L(w(t-3))$$

נותר למצוא רק את האיבר האמצעי, מתוך הקשר:

$$w(t-1) = w(t-2) - \eta_{t-1} \nabla L(w(t-2))$$

$$\frac{\partial w(t-1)}{\partial w(t-2)} = I - \eta_{t-1} \nabla^2 L(w(t-2))$$

סך הכל קיבלנו:

$$(**) = \left( \nabla L(w(t-1)) \right) \cdot \left( I - \eta_{t-1} \nabla^2 L(w(t-2)) \right) \cdot \left( -\nabla L(w(t-3)) \right)$$

## סעיף ד':

סעיף זה הוא הכללה של סעיף ג' למספר צעדים גדול, נחשב באותה דרך:

$$\begin{aligned}\frac{\partial L(w(t-1))}{\partial \eta_{t-\tau}} &= \frac{\partial L}{\partial w(t-1)} \cdot \frac{\partial w(t-1)}{\partial \eta_{t-\tau}} \\ &= \frac{\partial L}{\partial w(t-1)} \cdot \frac{\partial w(t-1)}{\partial w(t-2)} \cdot \frac{\partial w(t-2)}{\partial w(t-3)} \cdot \frac{\partial w(t-3)}{\partial w(t-4)} \dots \frac{\partial w(t-\tau)}{\partial \eta_{t-\tau}}\end{aligned}$$

ניעזר בסעיף ג' ונאמר שעבור איברים כלליים:

$$\begin{aligned}\frac{\partial w(t-k)}{\partial w(t-k-1)} &= I - \eta_{t-k} \nabla^2 L(w(t-k-1)) \\ \frac{\partial w(t-\tau)}{\partial \eta_{t-\tau}} &= -\nabla L(w(t-\tau-1))\end{aligned}$$

נכפול בכל האיברים ונקבל:

$$(**) = \left( \nabla L(w(t-1)) \right) \cdot \left( \prod_{k=1}^{\tau-1} \left( I - \eta_{t-k} \nabla^2 L(w(t-k-1)) \right) \right) \cdot \left( -\nabla L(w(t-\tau-1)) \right)$$

## סעיף ה':

עבור הגישה הראשונה שבה מחשבים כל מחזור פרמטר אחד:

- לכל עדכון יש סיבוכיות חישוב נמוכה שכן אנחנו מעדכנים רק פרמטר אחד
- קל יותר למימוש
- מתאים יותר לפונקציות שכדי להתכנס בהן גודל הצעד צריך להשתנות משמעותית בין סיבוב לסיבוב.

עבור הגישה השנייה שבה מחשבים כל מחזור  $T$  פרמטרים:

- לכל עדכון יש סיבוכיות חישוב גדולה שכן אנחנו מעדכנים  $T$  פרמטרים בו זמנית
- מסובך יותר למימוש
- מתאים יותר לפונקציות בהן גודל הצעד צריך להתחשב בשינויים רחבים יותר בפונקציה ולא רק בצעד האחרון. עבור  $T$  גדול מדי, ייתכן שיהיה קשה להתכנס כי גודל הצעד לא ישתנה במשך  $T$  מחזורים.

#### שאלה 4

נתונה הפונקציה:

$$f(x) = e^{e^x + e^{2x}} + \sin(e^x + e^{2x})$$

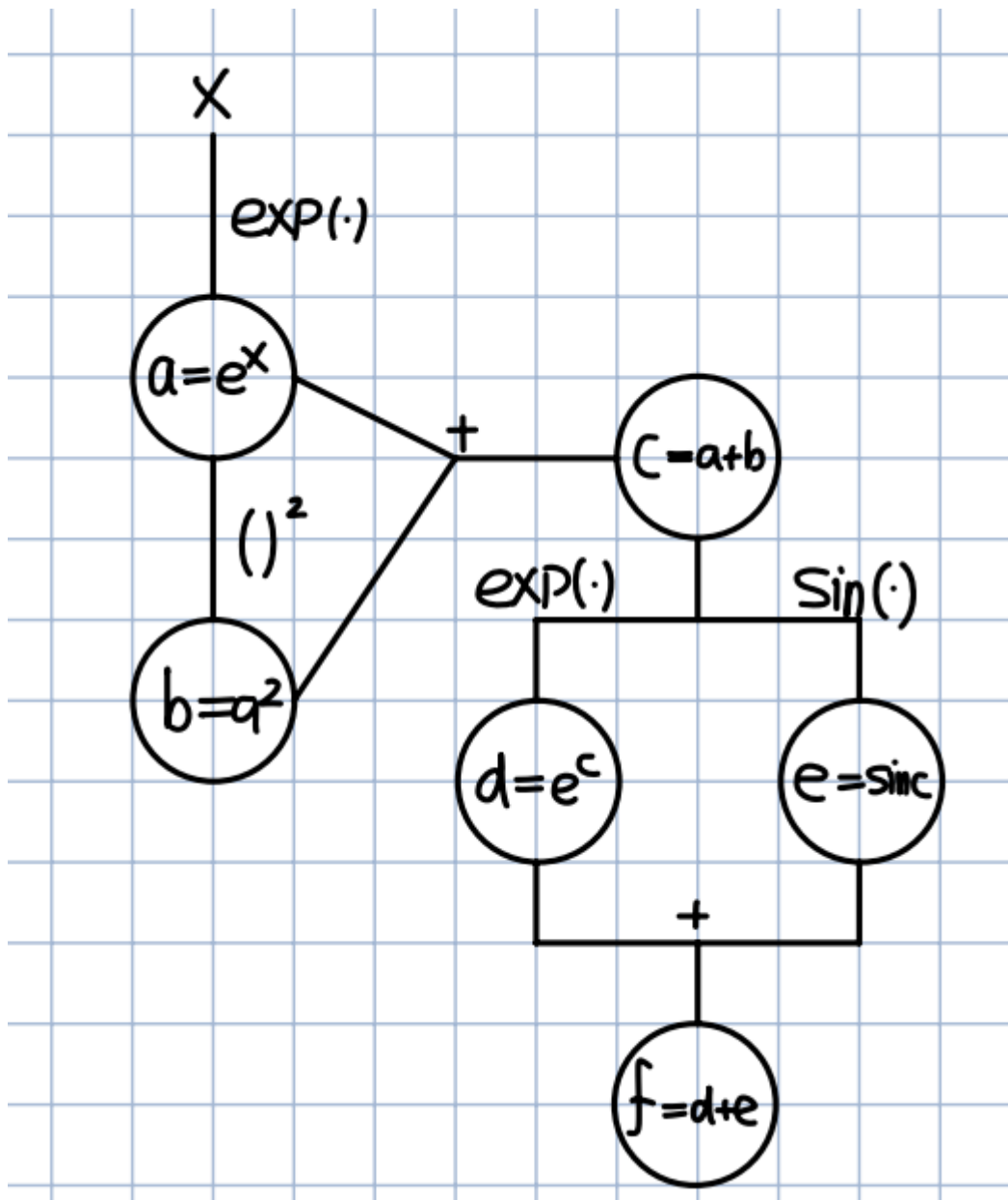
סעיף א':

נחשב ישירות את הנגזרת:

$$\frac{df}{dx} = (e^x + 2e^{2x})e^{e^x + e^{2x}} + (e^x + 2e^{2x})\cos(e^x + e^{2x})$$

$$= (e^x + e^{2x})(e^{e^x + e^{2x}} + \cos(e^x + e^{2x}))$$

סעיף ב':





## סעיף ג':

מכלל השרשרת נקבל:

$$f = f(d, e) = d + e$$

$$e = e(c) = \sin c$$

$$d = d(c) = e^c$$

$$c = c(a, b) = a + b$$

$$b = b(a) = a^2$$

$$a = a(x) = e^x$$

ולכן:

$$\frac{df}{dx} = \frac{df}{dd} \frac{dd}{dc} \left( \frac{dc}{da} \frac{da}{dx} + \frac{dc}{db} \frac{db}{da} \frac{da}{dx} \right) + \frac{df}{de} \frac{de}{dc} \left( \frac{dc}{da} \frac{da}{dx} + \frac{dc}{db} \frac{db}{da} \frac{da}{dx} \right)$$

$$= \left( \frac{df}{dd} \frac{dd}{dc} + \frac{df}{de} \frac{de}{dc} \right) \left( \frac{dc}{da} \frac{da}{dx} + \frac{dc}{db} \frac{db}{da} \frac{da}{dx} \right)$$

$$= (1 \cdot e^c + 1 \cdot \cos c)(1 \cdot e^x + 1 \cdot 2a \cdot e^x) = (e^{a+b} + \cos(a+b))(e^x + 2e^x e^x)$$

$$(e^{e^x+a^2} + \cos(e^x + a^2))(e^x + 2e^{2x}) = (e^{e^x+e^{2x}})(e^x + 2e^{2x})$$

כלומר, התוצאה המתקבלת זהה לתוצאה המתקבלת לפי החישוב הישיר. פרקטית, גם את המשתנה  $a + b$  חישבנו יותר מפעם אחת ויכולנו לשמור בצד לחישוב יעיל יותר.