

20942) מבוא ללמידה חישובית | סיכום הרצאה 11

מנחה: ד"ר שי מימון

סמסטר: 2022'

נכתב על ידי: מתן כהן

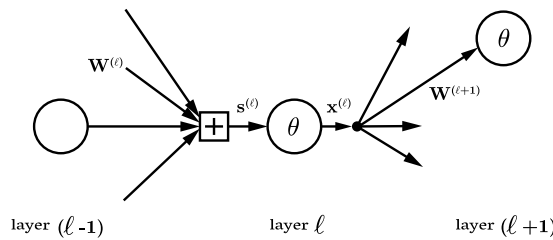
1 Backpropagation

מהות העניין הוא חישוב יעיל של Gradient Descent על פני כל שכבות הרשת. ניזכר כי על מנת לחשב את הגרדיאנט של E_{in} יש צורך בנגזרת שלו ביחס למטריצת המשקלים

$$\frac{\partial E_{in}}{\partial W^{(\ell)}} = \frac{1}{N} \sum_{n=1}^N \frac{\partial e_n}{\partial W^{(\ell)}}$$

נסמן:

- e - שגיאה על נקודה מסוימת (\mathbf{x}_n, y_n)
- $w_{ij}^{(\ell)}$ - משקל הקשת בין נירון i משכבה $\ell - 1$ לנירון j לשכבה ℓ



תחילה ניזכר כי:

$$(1.1) \quad \mathbf{s}_j^{(\ell)} = \sum_{\alpha=0}^{d^{(\ell-1)}} w_{\alpha j}^{(\ell)} \mathbf{x}_{\alpha}^{(\ell-1)} \Rightarrow \frac{\partial \mathbf{s}_j^{(\ell)}}{\partial w_{ij}^{(\ell)}} = \mathbf{x}_i^{(\ell-1)}$$

כעת נבחין כי המקום בו המשקולת $w_{ij}^{(\ell)}$ משפיעה על משתני הכניסה לשכבה ℓ הוא $\mathbf{s}_j^{(\ell)}$ ולכן:

$$\frac{\partial e}{\partial w_{ij}^{(\ell)}} \stackrel{\text{chain rule}}{=} \frac{\partial \mathbf{s}_j^{(\ell)}}{\partial w_{ij}^{(\ell)}} \cdot \frac{\partial e}{\partial \mathbf{s}_j^{(\ell)}} \stackrel{1.1}{=} \mathbf{x}_i^{(\ell-1)} \cdot \frac{\partial e}{\partial \mathbf{s}_j^{(\ell)}}$$

נסמן את ה-sensitivity:

$$\frac{\partial e}{\partial \mathbf{s}_j^{(\ell)}} = \delta_j^{(\ell)}$$

ונקבל:

$$\frac{\partial e}{\partial w_{ij}^{(\ell)}} = \mathbf{x}_i^{(\ell-1)} \cdot \delta_j^{(\ell)}$$

בעזרת מה שהראינו נוכל לכתוב:

$$\frac{\partial e}{\partial W^{(\ell)}} = \underbrace{\mathbf{x}^{(\ell-1)}}_{(d^{(\ell-1)}+1) \times 1} \cdot \left(\underbrace{\boldsymbol{\delta}^{(\ell)}}_{d^{(\ell)} \times 1} \right)^T$$

כאשר הממד של התוצאה הוא

$$(d^{(\ell-1)} + 1) \times d^{(\ell)}$$

כעת נתמודד עם $\delta_j^{(\ell)}$, לשם כך נבחין כי נוכל להשתמש ב $\mathbf{x}^{(\ell)}$ וכיוון שכל אלמנט $\mathbf{s}_j^{(\ell)}$ משפיע רק על אלמנט $\mathbf{x}_j^{(\ell)}$ נוכל להשתמש בזה על מנת לכתוב:

$$\delta_j^{(\ell)} = \frac{\partial e}{\partial \mathbf{s}_j^{(\ell)}} = \frac{\partial e}{\partial \mathbf{x}_j^{(\ell)}} \cdot \frac{\partial \mathbf{x}_j^{(\ell)}}{\partial \mathbf{s}_j^{(\ell)}} = \frac{\partial e}{\partial \mathbf{x}_j^{(\ell)}} \cdot \theta'(\mathbf{s}_j^{(\ell)})$$

כעת, אם נוכל לייצר מנגנון נוח לחישוב δ נוכל להתמודד איתה בכל שכבה בעת ה-Backpropagation בצורה פשוטה. לשם כך נתבונן בביטוי $\frac{\partial e}{\partial \mathbf{x}_j^{(\ell)}}$ בו השתמשנו על מנת לבטא את $\delta_j^{(\ell)}$ וזאת כאשר נזכור שכל $\mathbf{x}_j^{(\ell)}$ משפיע על כל $\mathbf{s}_k^{(\ell+1)}$ מהשכבה הבאה:

$$\frac{\partial e}{\partial \mathbf{x}_j^{(\ell)}} = \sum_{k=1}^{d^{(\ell+1)}} \frac{\partial \mathbf{s}_k^{(\ell+1)}}{\partial \mathbf{x}_j^{(\ell)}} \cdot \frac{\partial e}{\partial \mathbf{s}_k^{(\ell+1)}} = \sum_{k=1}^{d^{(\ell+1)}} w_{jk}^{(\ell+1)} \delta_k^{(\ell+1)}$$

קיבלנו בסוף שלכל $j = 1, \dots, d^{(\ell)}$:

$$\delta_j^{(\ell)} = \theta'(\mathbf{s}_j^{(\ell)}) \sum_{k=1}^{d^{(\ell+1)}} w_{jk}^{(\ell+1)} \delta_k^{(\ell+1)}$$

נגדיר שהתסכלות על איבר i עד איבר j בוקטור \mathbf{v} יכתב:

$$[\mathbf{v}]_i^j$$

מכפלת איבר באיבר של וקטורים \mathbf{u}, \mathbf{v} תכתב:

$$\mathbf{v} \otimes \mathbf{u}$$

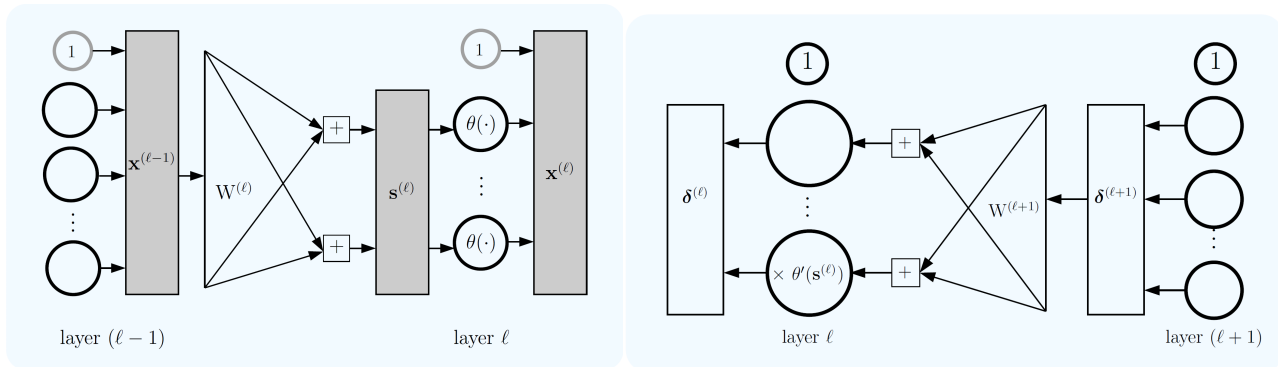
ועל פי הסימון נוכל לכתוב את $\delta^{(\ell)}$:

$$\delta^{(\ell)} = \theta'(\mathbf{s}^{(\ell)}) \otimes [W^{(\ell+1)} \delta^{(\ell+1)}]_1^{d^{(\ell)}}$$

כעת נוכל להסתכל על כלל ההתהליך.

1.1 תהליך ה- Backpropagation

ניזכר כי בהתחלה עשינו forward propagation על מנת לחשב את האיקסים $\mathbf{x}^{(\ell)}$ משכבת ה-input ועד שכבת ה-output כאשר חישבנו את $\mathbf{x}^{(\ell)}$ בעזרת $\mathbf{x}^{(\ell-1)}$.
ב-backpropagation נרצה לחשב את הדלתאות $\delta^{(\ell)}$ משכבת ה-output לאחור עד שכבת ה-input כאשר נחשב את $\delta^{(\ell)}$ בעזרת $\delta^{(\ell+1)}$.
ישנו דמיון בין שני התהליכים!



1.1.1 התחלת תהליך ה-backpropagation

- חשוב לזכור שבשכבה האחרונה אנחנו מדברים על סקלרים

נרצה להתחיל עם $\delta^{(L)}$ (הרי L זו השכבה האחרונה) ולכן על מנת לחשב את הדלתא:

$$\begin{aligned}\delta^{(L)} &= \frac{\partial e}{\partial \mathbf{s}^{(L)}} \stackrel{\text{last layer}}{=} \frac{\partial e}{\partial s^{(L)}} \\ &= \frac{\partial}{\partial s^{(L)}} \left(x^{(L)} - y \right)^2 \\ &= 2 \left(x^{(L)} - y \right) \cdot \frac{\partial x^{(L)}}{\partial s^{(L)}} \\ &= 2 \left(x^{(L)} - y \right) \cdot \theta' \left(s^{(L)} \right)\end{aligned}$$

כעת, אם θ היא \tanh הרי ש- $\theta' \left(s^{(L)} \right) = 1 - \left(x^{(L)} \right)^2$ ולכן:

$$\delta^{(L)} = 2 \left(x^{(L)} - y \right) \cdot \left(1 - \left(x^{(L)} \right)^2 \right)$$

אם מדובר ברגרסיה והשתמשנו לדוגמה בפונקציית הזהות הרי ש- $\theta' \left(s^{(L)} \right) = 1$:

$$\delta^{(L)} = 2 \left(x^{(L)} - y \right)$$

אלגוריתם ה-Backpropagation 1.1.2

אלגוריתם 1 Backpropagation to compute sensitivities $\delta^{(\ell)}$

Input: a data point (\mathbf{x}, y)

(1) Run forward propagation on \mathbf{x} to compute and save:

$$\mathbf{s}^{(\ell)} \text{ for } \ell = 1, \dots, L;$$

$$\mathbf{x}^{(\ell)} \text{ for } \ell = 0, \dots, L.$$

$$(2) \delta^{(L)} \leftarrow 2 (x^{(L)} - y) \theta' (s^{(L)})$$

[Initialization]

$$\theta' (s^{(L)}) = \begin{cases} 1 - (x^{(L)})^2 & \theta(s) = \tanh(s); \\ 1 & \theta(s) = s \end{cases}$$

(3) **for** $\ell = L - 1$ **to** 1 **do**

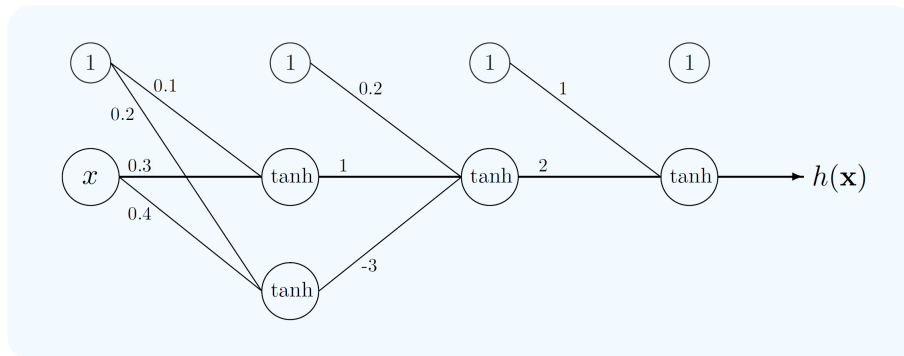
[Back-Propagation]

$$(a) \text{ Let } \theta' (\mathbf{s}^{(\ell)}) = [1 - \mathbf{x}^{(\ell)} \otimes \mathbf{x}^{(\ell)}]_1^{d^{(\ell)}} \quad \theta = \tanh$$

(b) Compute the sensitivity $\delta^{(\ell)}$ from $\delta^{(\ell+1)}$:

$$\delta^{(\ell)} \leftarrow \theta' (\mathbf{s}^{(\ell)}) \otimes [W^{(\ell+1)} \delta^{(\ell+1)}]_1^{d^{(\ell)}}$$

דוגמה. נתונה רשת עם 3 שכבות ($L = 3$) ונקודה $x = 2, y = 1$



נבחין כי:

$$d^{(0)} = 1, \quad d^{(1)} = 2, \quad d^{(2)} = 1, \quad d^{(3)} = 1$$

נגדיר בשלב הראשון את $W^{(\ell)}$ לכל $\ell = 1, \dots, 3$:

$$W^{(1)} = \begin{bmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \end{bmatrix}, \quad W^{(2)} = \begin{bmatrix} 0.2 \\ 1 \\ -3 \end{bmatrix}, \quad W^{(3)} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

טבלת ה-forward propagation:

$\mathbf{x}^{(0)}$	$\mathbf{s}^{(1)} = (W^{(1)})^T \mathbf{x}^{(0)}$	$\mathbf{x}^{(1)}$	$\mathbf{s}^{(2)} = (W^{(2)})^T \mathbf{x}^{(1)}$	$\mathbf{x}^{(2)}$	$\mathbf{s}^{(3)}$	$\mathbf{x}^{(3)}$
$\begin{bmatrix} 1 \\ 2 \end{bmatrix}$	$\begin{bmatrix} 0.1 & 0.3 \\ 0.2 & 0.4 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 0.7 \\ 1 \end{bmatrix}$	$\begin{bmatrix} 1 \\ \tanh(0.7) \\ \tanh(1) \end{bmatrix} = \begin{bmatrix} 1 \\ 0.6 \\ 0.76 \end{bmatrix}$	$\begin{bmatrix} 0.2 \\ 1 \\ -3 \end{bmatrix}^T \cdot \begin{bmatrix} 1 \\ 0.6 \\ 0.76 \end{bmatrix} = [-1.48]$	$\begin{bmatrix} 1 \\ -0.90 \end{bmatrix}$	$[-0.8]$	-0.66

נעת נעבור ל-backpropagation

נשתמש בהגדרת $\theta = \tanh$ ונחשב את $\delta^{(3)}$:

$$\delta^{(3)} = 2 \left(x^{(3)} - 1 \right) \left(1 - \left(x^{(3)} \right)^2 \right) = -1.855$$

נשתמש בלולאה:

for $\ell = L - 1$ to 1 **do**:

$$\text{Let } \theta'(\mathbf{s}^{(\ell)}) = [1 - \mathbf{x}^{(\ell)} \otimes \mathbf{x}^{(\ell)}]_1^{d^{(\ell)}}$$

והטבלה:

$\delta^{(3)}$	$\delta^{(2)}$	$\delta^{(1)}$
$[-1.855]$	$[(1 - 0.9^2) \cdot 2 \cdot -1.855] = [-0.69]$	$\begin{bmatrix} -0.44 \\ 0.88 \end{bmatrix}$

וכעת חישובי השגיאות:

$$\frac{\partial e}{\partial W^{(1)}} = \mathbf{x}^{(0)} \left(\delta^{(1)} \right)^T = \begin{bmatrix} -0.44 & 0.88 \\ -0.88 & 1.75 \end{bmatrix}$$

$$\frac{\partial e}{\partial W^{(2)}} = \begin{bmatrix} -0.69 \\ -0.42 \\ -0.53 \end{bmatrix}$$

$$\frac{\partial e}{\partial W^{(3)}} = \begin{bmatrix} -1.85 \\ 1.67 \end{bmatrix}$$

1.2 מרכיבים הכל ביחד

אלגוריתם 2 חישוב $E_{in}(\mathbf{w})$ והגרדיאנט $g = \nabla E_{in}(\mathbf{w})$

Input: $\mathbf{w} = \{W^{(1)}, \dots, W^{(L)}\}$; $\mathcal{D} = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$.

Output: error $E_{in}(\mathbf{w})$ and gradient $g = \{G^{(1)}, \dots, G^{(L)}\}$.

(1) Initialize: $E_{in} = 0$, and $G^{(\ell)} = 0 \cdot W^{(\ell)}$ for $\ell = 1, \dots, L$

(2) **for** each data point (\mathbf{x}_n, y_n) , $n = 1, \dots, N$, **do**

(a) Compute $\mathbf{x}^{(\ell)}$ for $\ell = 0, \dots, L$ [forward propagation]

(b) Compute $\delta^{(\ell)}$ for $\ell = L, \dots, 1$ [backpropagation]

(c) $E_{in} \leftarrow E_{in} + \frac{1}{N} \left(\mathbf{x}^{(L)} - y_n \right)^2$

(d) **for** $\ell = 1, \dots, L$ **do**

(i) $G^{(\ell)}(\mathbf{x}_n) = \left[\mathbf{x}^{(\ell-1)} \left(\delta^{(\ell)} \right)^T \right]$

(ii) $G^{(\ell)} \leftarrow G^{(\ell)} + \frac{1}{N} G^{(\ell)}(\mathbf{x}_n)$

כל התהליך הנ"ל נועד בשביל עדכון של איטרציה אחת בלבד - שיניב לנו סט משקלים חדש לכל שכבה באיטרציה הבאה:

$$W_{(t+1)}^{(\ell)} \leftarrow W_{(t)}^{(\ell)} - \eta G_{(t)}^{(\ell)}$$

1.3 בעיות נפוצות ושאלות פתוחות

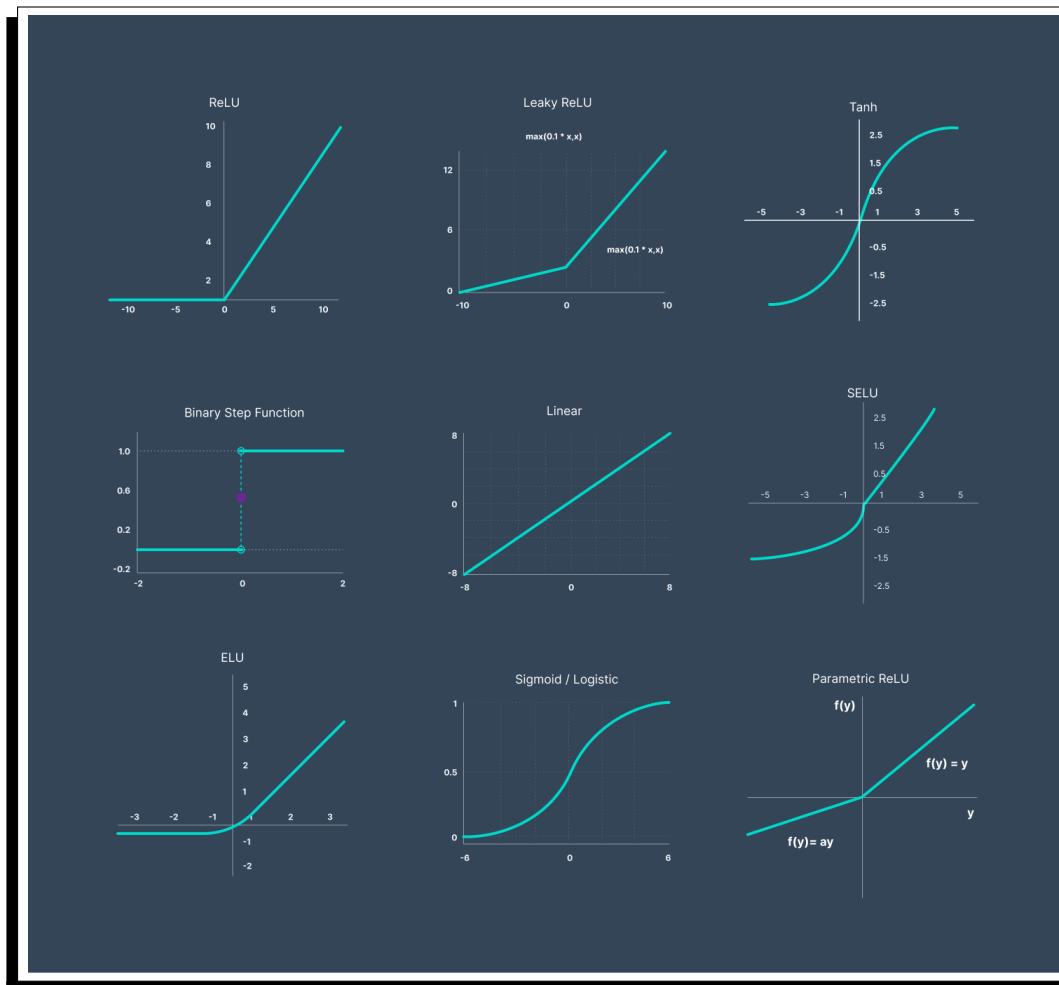
בעיה. מה עושים כאשר פונקציית המטרה שלנו לא קונבקסית? הרי יכול להיות ש GD יתכנס לנקודת מינימום מקומית

פתרון. במקום להתחיל מנקודה אחת את תהליך ה GD, נתחיל מכמה נקודות ולהשתמש בסט הולידציה על מנת להשוות בין החזאים השונים שקיבלנו.

הערה. אם נאתחל את המשקולות כאפסים ($W^{(\ell)} = 0$) אז לכל ℓ נקבל $\delta^{(\ell)} = 0$ וכתוצאה הנגזרות של השגיאה יהיו 0 וסיים לאחר איטרציה אחת ולא נתקדם.

בעיה. **Vanishing Gradient** הוא מצב בו ערכי $|w_{ij}|$ יהיו מאוד גבוהים ובקצוות פונקציית האקטיבציה שלנו (נניח ב-tanh) נגיע לערכי 1 ו-1- וכתוצאה מכך הנגזרת של θ תתקרב מאוד ל-0 וכתוצאה מכך גם ערכי $\delta^{(\ell)}$.

פתרון. שימוש בפונקציות אקטיבציה שונות (כמו $ReLU$) או החלפת פונקציית האקטיבציה בשכבת המוצא בפונקציית לינארית



1.3.1 אז כיצד נאתחל את המשקלים?

מומלץ לאתחל את המשקלים בצורה **בלתי תלויה סטטיסטית** עם משתנה אקראי **נורמלי**, בעל **תוחלת 0** ו**שונות קטנה**. בצורה זו המשקלים יהיו סביב נקודת האפס, אך לא אפס ממש.

1.3.2 סיום האלגוריתם

נהוג לשלב מספר קריטריוני עצירה

- מספר איטרציות
- שינוי קטן יחסית בשגיאה
- השגיאה האבסולוטית קטנה מערך מסוים

Stochastic (Online, Sequential) Gradient Descent 2

תהליך ה-GD שדיברנו עליו עד כה נקרא **Batch Gradient Descent** והוא מחושב עבור כל סט הדוגמאות לפני עדכון המשקלים.

גרסה אחרת של GD היא גרסה סטוכסטית שלו - **Stochastic Gradient Descent - SGD**. כמה נקודות על SGD:

- בכל פעם בוחרים באופן אקראי (אחיד) **דוגמה אחת** מסט הדוגמאות
- עבור הדוגמה הספציפית מחשבים את השגיאה ואת הגרדיאנט שמושפע מהדוגמה הספציפית הזו בלבד
- בעזרת החישובים שלנו על הנקודה \mathbf{x}_n הספציפית - נעדכן את סט הפרמטרים:

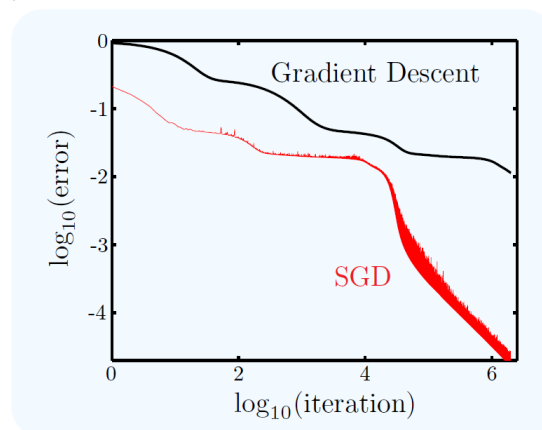
$$W^{(\ell)} = W^{(\ell)} - \eta G^{(\ell)}(\mathbf{x}_n)$$

- נקבל עדכונים מאוד **רועשים** אשר לאורך האימון יבטלו אחד את השני ויתכנסו לתוצאה טובה כמו Batch Gradient Descent

$$\eta \cdot \frac{1}{N} \sum_{n=1}^N \nabla e_n(\mathbf{w})$$

- האלגוריתם יעיל יותר מאלגוריתם GD המקורי כיוון שאין צורך למצוא N חישובי גרדיאנט, אלא רק אחד, משמע הוא **יעיל פי $\frac{1}{N}$**
- הרעשים שמתקבלים מאופן הפעולה תורמים לנו בהתחמקות ממינימום מקומי!
- אפשרי לשלב את 2 הטכניקות ולקחת תת-קבוצה של הדוגמאות ועליה לבצע GD

דוגמה. נתבונן בהבדלים בין Batch Gradient Descent ל-SGD. נתונות 500 דוגמאות, רשת בעלת 2 שכבות עם 5 שכבות חביות, קצב למידה $\eta = 0.01$



Variable Learning Rate Gradient Descent 3

היריסטיקה פשוטה שעוזרת בתהליך ה GD היא שינוי ערכי קצב הלמידה η
בכל צעד נתבונן בשגיאה

- אם השגיאה קטנה - נגדיל את ערך η
- אם השגיאה לא קטנה - נחזור צעד אחורה ונקטין את ערך η

אלגוריתם 3 Variable Learning Rate Gradient Descent

- (1) Initialize $\mathbf{w}(0)$, and η_0 at $t = 0$. Set $\alpha > 1$ and $\beta < 1$.
- (2) **while** stopping criterion has not been met **do**
 - (a) Let $\mathbf{g}(t) = \nabla E_{in}(\mathbf{w}(t))$, and set $\mathbf{v}(t) = -\mathbf{g}(t)$
 - (b) **if** $E_{in}(\mathbf{w}(t) + \eta_t \mathbf{v}(t)) < E_{in}(\mathbf{w}(t))$ **then**
 - (i) accept: $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta_t \mathbf{v}(t)$; $\eta_{t+1} = \alpha \eta_t$
 - (c) **else**
 - (i) reject: $\mathbf{w}(t+1) = \mathbf{w}(t)$; $\eta_{t+1} = \beta \eta_t$
- (d) Iterate to the next step, $t \leftarrow t + 1$

ונזכור:

$$\alpha \approx 1.05 - 1.1$$

$$\beta \approx 0.5 - 0.8$$

Steepest Descent 4

השאלה הנשאלת היא מדוע בכלל להסתבך עם קביעת η אם אפשרי לעשות חישוב אנליטי? הרי אנחנו רוצים לרדת בכיוון המנוגד לגרדיאנט.

לכן, במקום לבחור η כלשהו, נבחר η אופטימלי שיביא למינימום את E_{in}

אלגוריתם 4 Variable Learning Rate Gradient Descent

- (1) Initialize $\mathbf{w}(0)$, and η_0 at $t = 0$.
- (2) **while** stopping criterion has not been met **do**
 - (a) Let $\mathbf{g}(t) = \nabla E_{in}(\mathbf{w}(t))$, and set $\mathbf{v}(t) = -\mathbf{g}(t)$
 - (b) Let $\eta^* = \operatorname{argmin}_{\eta} E_{in}(\mathbf{w}(t) + \eta \mathbf{v}(t))$
 - (c) $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta^* \mathbf{v}(t)$
 - (d) Iterate to the next step, $t \leftarrow t + 1$

Regularization and Validation 5

- עד כה ראינו שרשתות נוירונים זהו כלי חזק מאוד שכולל שימוש ב GD ויכול לקרב אותנו לכל פונקציה שרק נרצה - משמע ישנה אופציה להתקל בהתאמת יתר (overfitting)
- גודל שכבת הכניסה אמנם אינו נתון לשינוי (גודל הקלט/פיצ'רים) וכך גם גודל שכבת היציאה (תלוי בבעיה אותה פותרים)
- מספר וגודל השכבות החבויות הוא פרמטר הניתן לשינוי
 - המחשבה הראשונית היא להקטין את מספר דרגות החופש שלנו ברשת (מספר השכבות החבויות והנוירונים)
 - מאידך מה שנהוג לעשות זה **בדיוק ההיפך** - לתת לרשת להיות כמה שיותר גדולה ולשלוט על התאמת היתר בעזרת רגולריזציה וולידציה

Weight Based Complexity Penalties 5.1

אחת שיטות הרגולריזציה המוכרות היא **Weight decay** - מה שמוכר לנו בתור Ridge Regression בבעית הרגרסיה ברשת הנוירונים נגדיר את השיטה בצורה דומה אך שונה

הגדרה. Weight squared decay - הוספת סכום ריבועי המשקולות בריבוע

$$E_{aug}(\mathbf{w}) = E_{in}(\mathbf{w}) + \frac{\lambda}{N} \sum_{\ell, i, j} \left(w_{ij}^{(\ell)} \right)^2$$

הרגולריזציה הנ"ל נקראת כך כיוון שבעת הגזירה נקבל איבר התלוי ב $W^{(\ell)}$:

$$\frac{\partial E_{aug}(\mathbf{w})}{\partial W^{(\ell)}} = \frac{\partial E_{in}(\mathbf{w})}{\partial W^{(\ell)}} + \frac{2\lambda}{N} W^{(\ell)}$$

ובעדכון המקדמים נכפיל בקצב הלמידה η ובמינוס (כדי ללכת נגד כיוון הגרדיאנט) - מה שיגרור דעיכה (decay) לכיוון האפס.

Weight Elimination 5.2

תפקידה "להעלים" משקלים קטנים

הגדרה. Weight Elimination

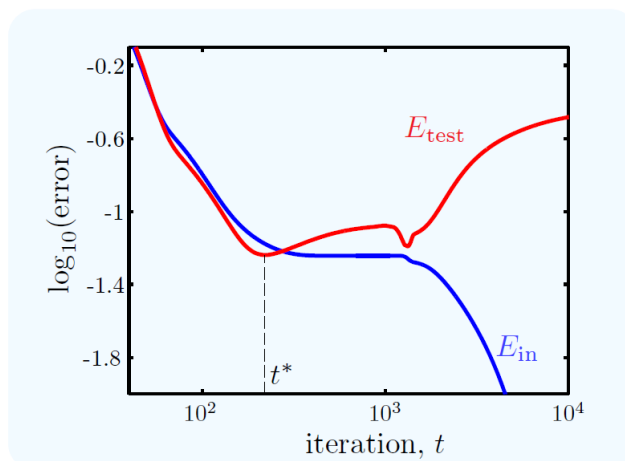
$$E_{aug}(\mathbf{w}, \lambda) = E_{in}(\mathbf{w}) + \frac{\lambda}{N} \sum_{\ell, i, j} \frac{\left(w_{ij}^{(\ell)} \right)^2}{1 + \left(w_{ij}^{(\ell)} \right)^2}$$

הנגזרת:

$$\frac{\partial E_{aug}}{\partial w_{ij}^{(\ell)}} = \frac{\partial E_{in}}{\partial w_{ij}^{(\ell)}} + \frac{2\lambda}{N} \cdot \frac{\left(w_{ij}^{(\ell)} \right)^2}{\left(1 + \left(w_{ij}^{(\ell)} \right)^2 \right)^2}$$

Early Stopping 5.3

שיטת נוספת לרגולריזציה - כוללת הפעלת GD אך במקום לבצע המון איטרציות, עוצרים בנקודה בה השגיאה על סף הולידציה עולה.

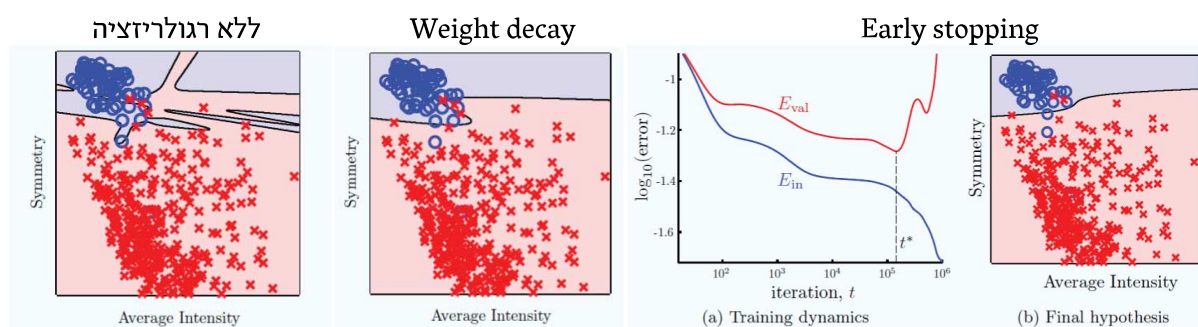


מהות שיטה זו לרגולריזציה מושתתת על העובדה שמספר איטרציות מוגבל גורר מספר היפתוזות מוגבל - משמע אנחנו לא נבדוק את כלל ההיפותוזות שלנו, אלא נתמקד בתת-קבוצה של היפותוזות!

הערה. שיטה זו קשורה מאוד לשיטת Weight decay כיוון שעצם העובדה שאנחנו מגבילים את האיטרציות גורמת לכך שהמשקולות תשארנה נמוכות וקרובות לאפס.

דוגמה לשימוש ברגולריזציה

500 דוגמאות, רשת נוירונים עם שכבה חבויה אחת בעלת 10 נוירונים, 2 פיצ'רים



	E_{train}	E_{val}	E_{in}	E_{out}
No regularization	-	-	0.2%	3.1%
Weight decay	-	-	1.0%	2.1%
Early stopping	1.1%	2.0%	1.2%	2.0%

Approximation vs Generalization 6

אפשר להראות שאם נבחר מספיק נגזרים עבור רשת בעלת 2 שכבות נוכל לקרב כל פונקציה שרצוה. באופן כללי ניתן לרשום את החזאי של רשת כזו:

$$h(\mathbf{x}) = \theta \left(w_{01}^{(2)} + \sum_{j=1}^m w_{j1}^{(\ell)} \theta \left(\sum_{i=0}^d w_{ij}^{(1)} x_i \right) \right)$$

בצורה יותר נוחה:

$$h(\mathbf{x}) = \theta \left(w_0 + \sum_{j=1}^m w_j \theta (\mathbf{v}_j^T \mathbf{x}) \right)$$

מאידך, נוכל להשתמש בפרספטון עם טרנס' לא לינארית ולבנות פרספטון שתלוי לא לינארית בפיצ'רים שלנו ופונקצית θ לא לינארית ולקבל:

$$\mathbf{x} = \Phi(\mathbf{x}) \rightarrow \mathbf{z} = [1, \phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_m(\mathbf{x})]^T$$

$$h(\mathbf{x}) = \theta \left(w_0 + \sum_{j=1}^M w_j \phi_j(\mathbf{x}) \right)$$

אך הדבר **השונה** הוא עצם העובדה שהטרנספורמציה מופעלת **אך ורק** על הפיצ'רים שלנו: $\phi_j(\mathbf{x})$ לעומת רשת הנגזרים **שפועלת** גם על המשקולות $\theta(\mathbf{v}_j^T \mathbf{x})$ - משמע בעזרת רשת הנגזרים אנחנו מתאימים את עצמנו **יותר טוב** למידע שלנו כי לומדים את ה- w_{ij} בהתאם לדוגמאות לעומת הפרספטון בו קבענו את ϕ_i עוד לפני האימון!