

מבוא ללמידה חישובית | סיכום הרצאה 10 (20942)

מנחה: ד"ר שי מימון

סמסטר: 2022'

נכתב על ידי: מתן כהן

1 ולידציה - Validation

- ולידציה היא עוד דרך התמודדות עם התאמת יתר (overfitting)
- גם רגולריזציה וגם ולידציה הן שיטות אשר מנסות להביא למינימום את ה- E_{out} בנוסף להקטנת E_{in}

הגדרה. קבוצת המבחן - Test set

קבוצת המבחן (Test set) היא תת-קבוצה של \mathcal{D} (סט הדוגמאות הכולל שלנו) **שלא נכללת בתהליך הלמידה** ומשמשת על אותנו בהערכת ההיפותזה הסופית שלנו $g \in \mathcal{H}$. במילים אחרות - קבוצת המבחן משמשת אותנו להערכת ביצועי המודל.

הערה. תוצאות השגיאה E_{test} על קבוצת המבחן הן שיערוך **בלתי מוטה** (unbiased) ל- E_{out}

הגדרה. קבוצת ולידציה - Validation set

- קבוצת הולידציה היא תת-קבוצה של דוגמאות **שלא נכללת בתהליך הלמידה** ומשמשת לאמידת ה- E_{out}
- אף על פי שקבוצת הולידציה לא נכללת בצורה ישירה בתהליך אימון המודל, היא כן עוזרת **בקבלת החלטות** בזמן תהליך הלמידה
- קבוצת הולידציה משמשת אותנו **בעת תהליך הלמידה** והיא **קטנה עד כדי** שהשערוך שלה ל- E_{out} לא אופטימיסטי.

הערה. **ברגע** שקבוצת דוגמאות **משפיעה** בצורה מסוימת על תהליך הלמידה - היא לא חלק מ**קבוצת המבחן** ולכן קבוצת הולידציה היא לא חלק מקבוצת המבחן.

1.1 K מתוך N

בהינתן סט דוגמאות $\mathcal{D} = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ נחלק את הסט ל-2 קבוצות, \mathcal{D}_{train} ו- \mathcal{D}_{val} בצורה **שלא מבוססת על ערכי הדוגמאות** כך ש:

- בקבוצה \mathcal{D}_{val} יש K דוגמאות
- בקבוצה \mathcal{D}_{train} יש $N - K$ דוגמאות

כעת, החזאי האופטימלי שיתבסס על דוגמאות מתוך \mathcal{D}_{train} נסמנו ב- g^- (כיוון שהוא לא כולל את כל הדוגמאות). כמו-כן נסמן את השגיאה בין ערך החזאי האופטימלי של g^- על נקודה \mathbf{x}_n לבין הערך האמיתי של הנקודה y_n בתור $e(g^-(\mathbf{x}_n), y_n)$ ונבחין כי **תוחלת השגיאה היא ה- $E_{out}(g^-)$** :

$$\mathbb{E}_{\mathbf{x}_n} [e(g^-(\mathbf{x}_n), y_n)] = E_{out}(g^-) \quad (1.1)$$

על מנת לחשב את ה- E_{val} על g^- נחשב:

$$E_{val}(g^-) = \frac{1}{K} \sum_{\mathbf{x}_n \in \mathcal{D}_{val}} e(g^-(\mathbf{x}_n), y_n)$$

כעת, נרצה לחשב את התוחלת על השגיאה E_{val} על פי g^- כיוון שהאיבר $e(g^-(\mathbf{x}_n), y_n)$ הוא בעצם משתנה אקראי. כמו-כן מההנחה שלנו שכלל הדוגמאות בלתי תלויות סטטיסטית נוכל להשתמש בתכונות התוחלת ובמה שהראינו ב- 1.1 ולקבל:

$$\begin{aligned}\mathbb{E}_{\mathcal{D}_{val}} [E_{val}(g^-)] &= \frac{1}{K} \sum_{\mathbf{x}_n \in \mathcal{D}_{val}} \mathbb{E}_{\mathcal{D}_{val}} [e(g^-(\mathbf{x}_n), y_n)] \\ &\stackrel{1.1}{=} \frac{1}{K} \sum_{\mathbf{x}_n \in \mathcal{D}_{val}} E_{out}(g^-)\end{aligned}$$

ולכן כיוון שסוכמים K קבועים (כי בקבוצה \mathcal{D}_{val} ישנן K דוגמאות) נקבל:

$$\mathbb{E}_{\mathcal{D}_{val}} [E_{val}(g^-)] = E_{out}(g^-)$$

בצורה דומה נגדיר את שונות השגיאה:

$$\text{Var}_{\mathbf{x}} [e(g^-(\mathbf{x}_n), y_n)] = \sigma^2(g^-)$$

ומכך שכלל הדוגמאות בלתי תלויות סטטיסטית נוכל להשתמש בהגדרת השונות ובעובדה שלכל משתנה אקראי X וקבוע α :

$$\text{Var}(\alpha \cdot X) = \alpha^2 \text{Var}(X)$$

ולכל שני משתנים אקראיים בת"ס (או חסרי קורלציה) X_1, X_2 :

$$\text{Var}(X_1 + X_2) = \text{Var}(X_1) + \text{Var}(X_2)$$

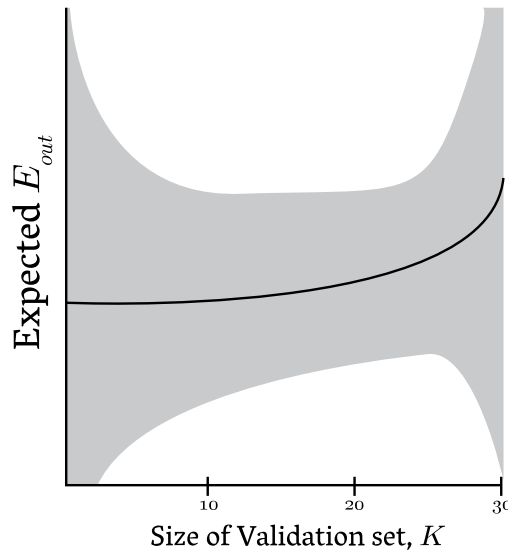
ובאופן דומה לחלוטין למה שנעשה מעלה:

$$\text{Var}_{\mathcal{D}_{val}} [E_{val}(g^-)] = \frac{1}{K^2} \sum_{\mathbf{x}_n \in \mathcal{D}_{val}} \text{Var} [e(g^-(\mathbf{x}_n), y_n)] = \frac{1}{K} \sigma^2$$

מכך נסיק:

במידה ונרצה שהשונות תהיה כמה שיותר קטנה (ואנחנו רוצים) הרי שהיא תלויה ב- K ולכן נסיק שככל שתהיינה לנו יותר דגימות בולידציה כך הדיוק בשיעורן שלנו E_{out} ישתפר!

דוגמה. בדומה לדוגמה מהשיעור הקודם, גם כאן נתבונן בהיפותזה מתוך \mathcal{H}_2 (פולינומים מדרגה 2), עבור $Q_f = 10$ (דרגת פולינום על פי נוצרו הדוגמאות), עבור $N = 40$ דוגמאות ודרגת רעש 0.4.



\mathcal{H}_2 , $Q_f = 10$, $N = 40$, $\text{noise} = 0.4$

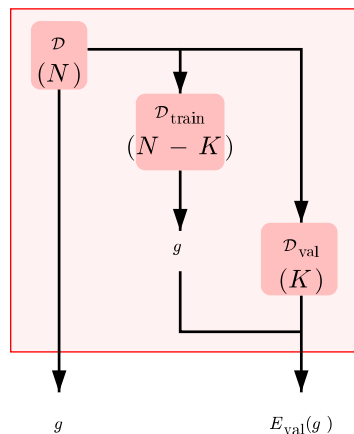
- על ציר ה- y נמצאת תוחלת $E_{out}(g^-)$
- על ציר ה- x נמצא את מספר הדוגמאות בסט הולידציה
- הקו השחור הוא אותו יחס בין מספר הדוגמאות בסט הולידציה לבין התוחלת
- השטח האפור מציג את השונות/סטיית התקן סביב הממוצע - $\sqrt{\frac{\sigma^2}{K}}$

נקודות חשובות על סמך הדוגמה:

- בנקודה מסוימת נראה שתוחלת השגיאות עולה - זו כיוון שבנקודה זו K גדול מידי וסט האימון שכולל $N - K$ דוגמאות דליל מאוד - דבר שפוגע בתהליך האימון
- השונות (השטח האפור) הולכת וקטנה עד נקודה מסוימת כאשר מגדילים את K
 - הנקודה בה השונות חוזרת לגדול היא אותה נקודה בה השגיאה עולה - משמע הגענו למצב בו K גדול מידי ואין לנו מספיק דוגמאות להתאמן עליהן
- כלל אצבע הוא הקצאה של $K = \frac{N}{5}$ עבור סט הולידציה

מסקנה. K צריך להיות מצד אחד גדול כדי שהדיווח שלנו על טיב המודל יהיה אמין אך מצד שני הוא צריך להיות מספיק קטן על מנת שנוכל לקבל חזאי טוב.

בסופו של דבר נוכל לאמן את החזאי שלנו על סט מצומצם עם $N - K$ דוגמאות ולקבל חזאי g^- ובנוסף לאמן מודל נוסף על כל N הדוגמאות ולקבל חזאי g שבו נשתמש.



1.2 בחירת מודל - Model Selection

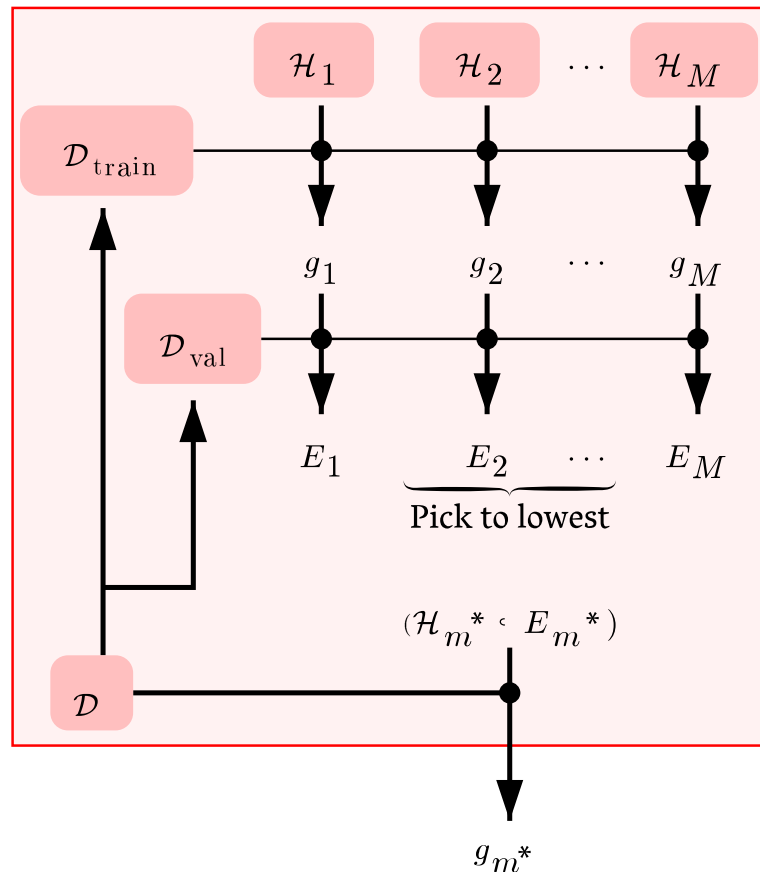
אחד השימושים הכי משמעותיים של הולידציה זה להדריך אותנו בעת הלמידה שלנו - וזהו הדבר המבדיל בין Test set ל- Validation set.

בפרט, הולידציה עוזרת לנו לבחור מודל מבין כלל המודלים שהתקבלו מתהליך הלמידה.

- בחירה בין מודל לינארי ללא לינארי
- בחירת דרגת פולינום במודל
- בחירת ערכים לפרמטרי רגולריזציה
- או כל פרמטר אחר שמשפיע על תהליך הלמידה שלנו

מסקנה. כל היפר-פרמטר שקיים בתהליך הלמידה ניתן לקביעה על פי הולידציה ובעזרת עובדה זו נוכל לאמן מודלים שונים עם היפר-פרמטרים שונים על מנת למצוא מודל אופטימלי.

נעשה זאת על ידי קביעת סט דיסקרטי של פרמטרים מראש (ניתן לבחור פילוג בצורה אחידה, לוגריתמית וכו') ונאמן חזאים שונים על פי אותם פרמטרים. לאחר מכן נוכל לבחור את המודל שנתן את הביצועים הכי טובים ביחד עם הפרמטרים שלו.



הסבר:

- (1) מכלל המודלים וההיפותזות $\mathcal{H}_1, \dots, \mathcal{H}_M$ נקבל חזאי אופטימלי מתאים g_1^-, \dots, g_M^- על פי סט האימון D_{train}
- (2) נבחן את ה- E_{out} על סט הולידציה D_{val}
- (3) נבחר את המינימלי מבין E_1, \dots, E_M בתקווה שייצג את הביצועים הכי טובים על Out-of-sample data
- (4) על פי אותו E_m שבחרנו נקבל היפר-פרמטרים אופטימליים ונפעיל את אותו חזאי עם הפרמטרים על D (כל הדוגמאות)
- (5) נקבל g_m^* שהוא החזאי הרצוי

Cross Validation 1.3

כפי שאמרנו, ככל ש- K קטן, השגיאה של g מתקרבת אל השגיאה של g^- כיוון שאנחנו מאמנים על "כמעט" אותו מספר של נקודות.

כמו-כן, ככל ש- K גדל, השגיאה של g^- מתקרבת לשגיאה על סט הולידציה של g^- ובאופן יותר ברור:

$$E_{out}(g) \underset{\text{small } K}{\approx} E_{out}(g^-) \underset{\text{large } K}{\approx} E_{val}(g^-)$$

אז מצד אחד אנחנו רוצים K קטן אך מצד שני אנחנו רוצים K גדול, מה עושים!?

1.3.1 Leave one out cross validation

בשיטה זו נפתור את הדילמה לעיל.

אלגוריתם Leave one out

(1) לכל n מ-1 ועד N :

(א) הגדר את סט האימון:

$$\mathcal{D}_n = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{n-1}, y_{n-1}), (\mathbf{x}_{n+1}, y_{n+1}), \dots, (\mathbf{x}_N, y_N)$$

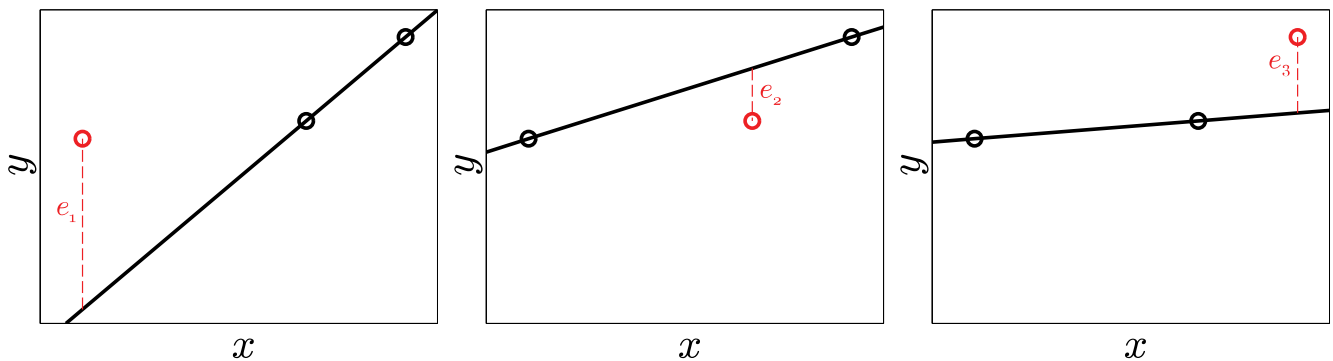
(ב) סמן את ההיפותזה הסופית שנלמדה מ- \mathcal{D}_n ב- g_n^-

(ג) סמן את השגיאה של ההיפותזה שהתקבלה על ידי $e_n = E_{val}(g_n^-) = e(g_n^-(\mathbf{x}_n), y_n)$

(2) הגדר את שגיאת ה- Cross validation error:

$$E_{cv} = \frac{1}{N} \sum_{n=1}^N e_n$$

דוגמה. ל-Leave one out Cross Validation



$$E_{cv} = \frac{1}{3} (e_1 + e_2 + e_3)$$

הסבר:

ישנן 3 נקודות, לכן ביצענו 3 איטרציות כאשר בכל אחת מהאיטרציות השארנו נקודה אחת בחוץ לבדיקת השגיאה. לבסוף מצאנו את E_{cv} בתור ממוצע השגיאות.

K-fold cross validation 1.3.2

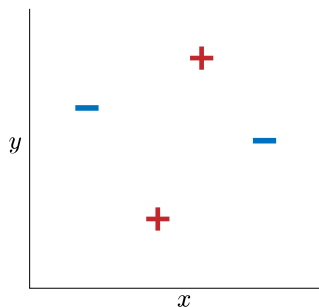
זוהי דרך נוספת להתמודדות עם הבעיה, היא זהה ל-Leave one out מבחינת המתודולוגיה אך כאן נרצה להשאיר קבוצה יותר גדולה של דוגמאות בחוץ בגודל K , ישנם $\frac{N}{K}$ מודלים שנאמן בצורה זו.

2 רשתות נוירונים - Neural Networks

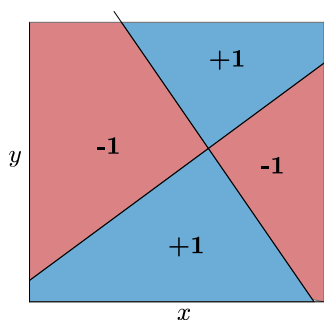
- הכללה של הפרספטרון
- יכולת למדל עם פונקציות מטרה מורכבות ()

2.1 Multi-layer Perceptron

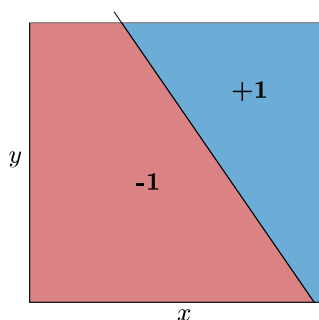
ניזכר כי מודל הפרספטרון לא הצליח להתמודד עם בעית ה-XOR שכן אף קו לינארי לא יכול להפריד את הדוגמה הבאה:



אך אם נוכל להשתמש בשני חזאים? האם נוכל לפתור את הבעיה?
התשובה היא - כן, אם נגדיר שאם שני החזאים מסכימים על הסימן שלהם, נחזיר -1 ואם לא מסכימים נחזיר $+1$.



$$h_1(\mathbf{x}) = \text{sign}(\mathbf{w}_1^T \mathbf{x})$$



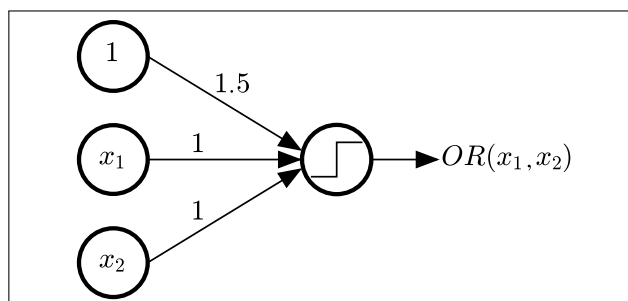
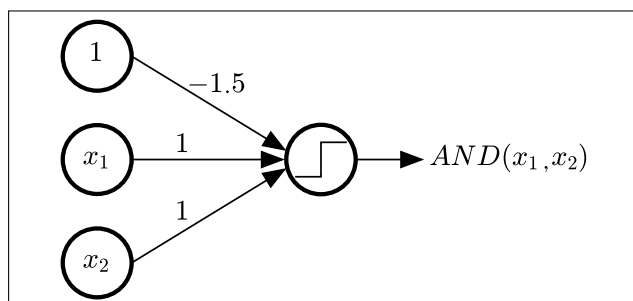
$$h_2(\mathbf{x}) = \text{sign}(\mathbf{w}_2^T \mathbf{x})$$

נבחין כי זו פעולת XOR ולכן נרצה להביע את הפעולה בעזרת:

$$h = \bar{h}_1 \cdot h_2 + h_1 \bar{h}_2$$

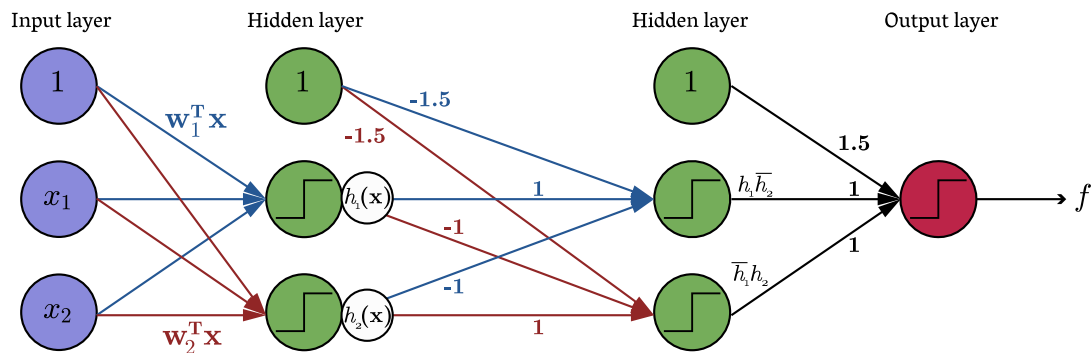
ולכן נרצה לממש שערי AND, OR ו-NOT.

2.1.1 מימוש שערי AND, OR ו-NOT:



על מנת לממש שער NOT כל שצריך הוא להפוך את סימני המשקולות של הקלט.

כעת נוכל לייצר את שער ה-XOR בעזרת שערי OR, AND והפיכת סימני המשקולות (NOT):



Sigmoidal Neural Networks 2.2

ראינו כי השימוש בפונקציית $sign$ הוא בעייתי כיוון שזוהי אינה פונקציה "חלקה" מה שגורם לבעיה הקומבינטורית להיות אפילו קשה יותר.

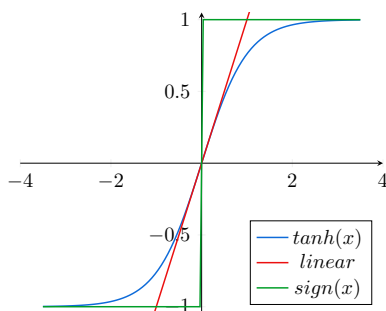
היינו רוצים פונקציה חלקה וגזירה שדומה לפונקציית ה- $sign$ שתאפשר לנו להשתמש בכלים אנליטיים כמו גרדיאנט על מנת למצוא משקלים אופטימליים.

דוגמה. נוכל להשתמש ב \tanh :

$$\tanh(x) = \frac{\sinh(x)}{\cosh(x)} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\lim_{x \rightarrow \infty} \tanh = 1, \quad \lim_{x \rightarrow -\infty} \tanh = -1$$

ואז בעזרת המקדם של x נוכל לשלוט על מהירות ההתכנסות של \tanh ל $+1$ ו- -1 .
באופן כללי נוכל לראות את השוני בין הפונקציות:



בעזרת שימוש בפונקציות הללו נוכל לבנות רשתות נוירונים.

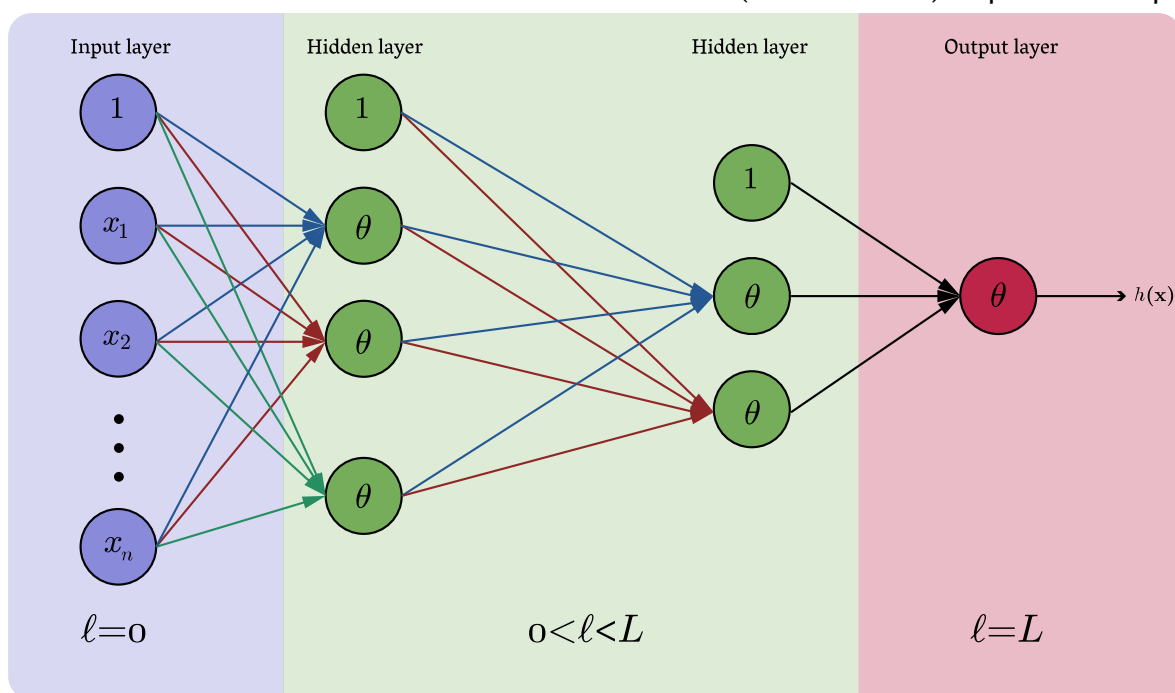
Neural Networks - 'softened' MLP 2.3

2.3.1 ארכיטקטורת הרשת:

- Input layer - שכבת הכניסה (קלט): השכבה הראשונה ברשת, לא נספרת מבין השכבות - $\ell = 0$
- Hidden layers - שכבות חביויות: כלל השכבות שנמצאות בין שכבת הכניסה לשכבת המוצא - $0 < \ell < L$
- Output layer - שכבת המוצא: ממנה יוצא החזאי - $\ell = L$
- בנוסף לפיצורים מוסיפים את האיבר ההטיה (1), אין לו כניסות והוא מהווה כניסה לפונקציית האקטיבציה שלנו

2.3.2 שיטת הפעולה

- כל node בכל שכבה מוציא ערך
- כל קשת בין node ל-node מכילה משקל שמוכפל בערך היוצא מ-node המקור אל יעדו
- ערך הכניסה ל-node הוא חיבור תוצאות המכפלה שנוצרה במשקולות
- מפעילים את פונקציית האקטיבציה (כמו \tanh) על התוצאה שהתקבלה מהחיבור
- נקבל במוצא פונקציה (לרוב לא ליניארית) מורכבת של x ושל w



2.3.3 שימוש ברשתות נוירונים לבעיות למידה שונות

נוכל להשתמש ברשתות נוירונים עבור בעיות למידה שונות

- **רגרסיה** - נחליף את θ בשכבת המוצא בפונקציית הזהות
- **הסתברות/רגרסיה לוגיסטית** - נחליף את θ בשכבת המוצא ב sigmoid
- **סיווג** - נחליף את הסיגמויד θ בשכבת המוצא בפונקציית sign

2.3.4 אופטימיזצית הרשת

ברגע שקבענו את הארכיטקטורה (מספר השכבות ומספר ה-nodes בכל שכבה), נרצה לחפש פרמטרים אופטימליים שיתאימו את המודל לנתונים בצורה הטובה ביותר.

- **רגרסיה** - שימוש בשגיאה ריבועית: $e(\mathbf{x}) = (h(\mathbf{x}) - y)^2$ + מיצוע ביחס לכלל הדוגמאות
- **סיווג** - שימוש ב cross entropy (ידובר בהמשך) או במדדי שגיאה אחרים

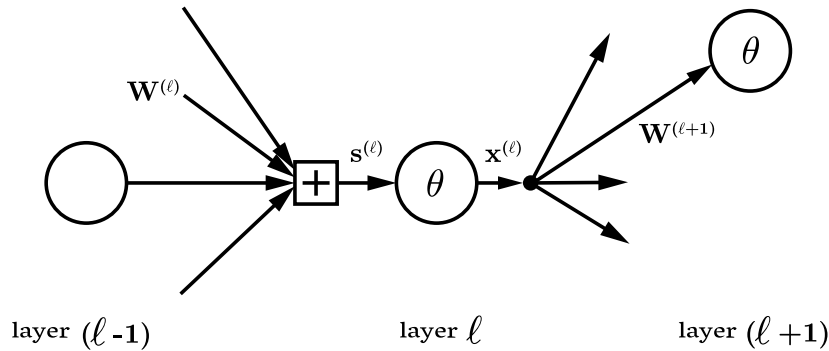
2.4 ארכיטקטורת הרשת

כפי שנאמר ארכיטקטורת הרשת מוגדרת על ידי:

- מספר השכבות L
- ממד כל שכבה (מספר הנוירונים/nodes)

ההיפותזה מוגדרת על ידי בחירה של משקלים ברשת - משמע כל בחירה של משקלים נותנת היפותזה אחרת.

2.4.1 סימונים



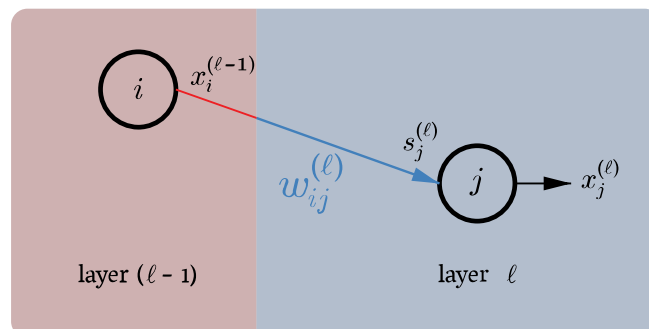
משמעות	שם	סימון	ממד
כניסה לאקטיבציה בשכבה ℓ (לאחר סכימה)	signals in	$\mathbf{s}^{(\ell)}$	$d^{(\ell)}$
יציאה מפונקציית האקטיבציה בשכבה ℓ	outputs	$\mathbf{x}^{(\ell)}$	$d^{(\ell+1)}$
מטריצה של משקולות הנכנסות לשכבה ℓ	weights in	$W^{(\ell)}$	$(d^{(\ell-1)} + 1) \times d^{(\ell)}$
מטריצה של משקולות שיוצאות משכבה ℓ	weights out	$W^{(\ell+1)}$	$(d^{(\ell)} + 1) \times d^{(\ell+1)}$

- אין ערכים שנכנסים לשכבת הכניסה ולכן אין $W^{(0)}$ או $s^{(0)}$
- אברי המטריצה $W^{(\ell)}$ הם מהצורה: $W_{ij}^{(\ell)}$ כך ש:

כל שורה i שייכת ל-node בשכבה $\ell - 1$ וכל עמודה j שייכת ל-node בשכבה ℓ

- התוספת של $+1$ לממדים היא כתוצאה מהוספת ההטיה (bias)

zoom-in למעבר בין נוירונים i ו- j :



- מ- i יוצא המוצא $x_i^{(\ell-1)}$
- הקשת בין i ל- j בעלת המשקולת $w_{ij}^{(\ell)}$
- המכפלה בין $w_{ij}^{(\ell)}$ ל- $x_i^{(\ell-1)}$ מסומנת ב- $s_j^{(\ell)}$

Forward Propagation 2.5

כפי שהסברנו עד כה, אנחנו "מעבירים קדימה" את המכפלות שלנו בצורה הבאה:

- בכל שכבה מכפילים את $\mathbf{x}^{(\ell)}$ ב- $W^{(\ell+1)}$ ומקבלים $\mathbf{s}^{(\ell+1)}$:

$$\mathbf{x} = \mathbf{x}^{(0)} \xrightarrow{W^{(1)}} \mathbf{s}^{(1)} \xrightarrow{\theta} \mathbf{x}^{(1)} \xrightarrow{W^{(2)}} \mathbf{s}^{(2)} \xrightarrow{\theta} \mathbf{x}^{(2)} \dots \rightarrow \mathbf{s}^{(L)} \xrightarrow{\theta} \mathbf{x}^{(L)} = h(\mathbf{x})$$

- האלגוריתם:

אלגוריתם Forward Propagation

(1) התחל את $\mathbf{x}^{(0)}$ בערך \mathbf{x}

(2) לכל $\ell = 1$ ועד L בצע:

(א) הגדר את $\mathbf{s}^{(\ell)}$ בתור:

$$\mathbf{s}^{(\ell)} \leftarrow \left(W^{(\ell)}\right)^T \mathbf{x}^{(\ell-1)}$$

(ב) הגדר את $\mathbf{x}^{(\ell)}$ בתור:

$$\mathbf{x}^{(\ell)} \leftarrow \begin{bmatrix} 1 \\ \theta(\mathbf{s}^{(\ell)}) \end{bmatrix}$$

(3) החזר $h(\mathbf{x}) = \mathbf{x}^{(L)}$

○ מדוע $\mathbf{s}^{(\ell)}$ מוגדר כך?

★ כיוון שאנחנו מחשבים את הסכום הממושקל שכולל פלטים מהשכבה הקודמת ומשקולות בין השכבה הקודמת לנוכחית

$$s_j^{(\ell)} = \sum_{i=0}^{d^{(\ell-1)}} w_{ij}^{(\ell)} x_i^{(\ell-1)}$$

• זוהי בעצם מכפלה סקלרית בין העמודה ה- j -ית של W לבין הוקטור \mathbf{x} !

★ לפיכך נוכל פשוט לכתוב בכתוב וקטורי:

$$\mathbf{s}^{(\ell)} = \left(W^{(\ell)}\right)^T \mathbf{x}^{(\ell-1)}$$

- לבסוף נקבל את ההיפותזה h ונבדוק את השגיאה:

$$E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n; \mathbf{w}) - y_n)^2$$

$$\stackrel{3}{=} \frac{1}{N} \sum_{n=1}^N \left(\mathbf{x}_n^{(L)} - y_n\right)^2$$

○ נבחין כי נצטרך לאפסם את השגיאה על פי \mathbf{w} וזוהי משימה לא פשוטה, לשם כך נשתמש ב Back propagation

2.6 Backpropagation Algorithm

אלגוריתם ה-Backpropagation עוזר לנו לאפסם את פונקציית השגיאה שראינו ממקודם. כיוון ש $h(\mathbf{x})$ היא פונקציה "חלקה" (גזירה), נוכל להפעיל **Gradient Descent** על פונקציית השגיאה.

נתחיל במוטיב:

דוגמה. The sigmoidal perceptron

נתון פרספטרון שבמקום פונקציית ה- sign נתונה פונקציית ה- \tanh :

$$\begin{aligned} h(\mathbf{x}) &= \tanh(\mathbf{w}^T \mathbf{x}) \\ E_{in}(\mathbf{w}) &= \frac{1}{N} \sum_{n=1}^N (\tanh(\mathbf{w}^T \mathbf{x}_n) - y_n)^2 \end{aligned}$$

נבדוק את הגרדיאנט:

• תחילה נבדוק את נגזרת \tanh ונראה כי ישנו פתרון סגור:

$$\frac{\partial \tanh}{\partial s} = \frac{\partial \sinh(s)}{\partial s \cosh(s)} = \frac{\cosh^2(s) - \sinh^2(s)}{\cosh^2(s)} = \boxed{1 - \tanh^2(s)}$$

$$\nabla E_{in}(\mathbf{w}) = 2 \cdot \frac{1}{N} \sum_{n=1}^N (\tanh(\mathbf{w}^T \mathbf{x}_n) - y_n) \cdot (1 - \tanh^2(\mathbf{w}^T \mathbf{x}_n)) \cdot \mathbf{x}_n$$

- במידה ונשתמש ב MLP (רשת נוירונים) לא יהיה פתרון פשוט וסגור, לשם כך נשתמש באלגוריתם ה-Backpropagation
- אלגוריתם זה משתמש בשיטת ה Gradient Descent שנמצאת בסיכום שיעור 4, אך להלן חזרה קצרה:

○ זוהי טכניקה למציאת מינימום של פונקציה גזירה

○ בהתחשב במשקולות ההתחלתיות, הנתיב שהאלגוריתם יניב יביא אותנו למינימום מקומי

○ נרצה תמיד לקחת צעד כנגד כיוון הרדיאנט

○ נהוג להתחל עם משקולות קטנות בצורה אקראית

○ נהוג לסיים על פי:

★ מספר צעדים/איטרציות

★ בחינת הנורמה והשגיאה

★ חסם עליון לשגיאה (נרצה שיגיע מתחת לערך מסוים)

- כיוון שיש מס' נקודות אופטימום, נהוג להתחיל ממספר נקודות שונות ואז לבחור את הנקודה שהניבה את הביצועים הטובים ביותר

אלגוריתם Gradient Descent

(1) **התחל** את המשקולות של צעד $t = 0$ ל $\mathbf{w}(0)$, ואת גודל הצעד η

(2) **לכל** $t = 0, 1, 2, \dots$ **בצע:**

(א) **חשב** את הגרדיאנט: $\mathbf{g}_t = \nabla E_{in}(\mathbf{w}(t))$

(ב) **קבע** את כיוון התזוזה: $\mathbf{v}_t = -\mathbf{g}_t$

(ג) **עדכן** את המשקולות: $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \mathbf{v}_t$

(ד) **המשך** לצעד הבא אם לא הגענו לתנאי עצירה

(3) **החזר** את סט המשקולות הסופי

מה שנרצה לעשות באלגוריתם ה-Backpropagation הוא דבר פשוט:

- נרצה להסתכל על מיצוע השגיאות $E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N e_n$
- ולמצוא מינימום בהתאם למטריצת המשקולות:

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

- האלגוריתם מתבסס על שימוש חוזר בכלל השרשרת במהלך מעבר מסוף רשת הנוירונים להתחלתה.