

למידה עמוקה

סיכום קורס שהתקיים בבר-אילן בשנת 2014 ע"י פרופ' יעקב גולדברגר. סוכם והוקלד ע"י עמית מנדלבום,
Disclamier: יש כמה נושאים חשובים שחסרים כאן (Auto encoders, CNN (convent), RNN)
אך למי שאין רקע בתחום הסיכום מכסה את כל הנושאים הבסיסיים וכן קצת מעבר.

תוכן

2.....	למידה עמוקה
2.....	Single neuron/ logistic regression
4.....	Multi-Class logistic regression
5.....	Neural network with 1 hidden layer
7.....	Back propagation algorithm
9.....	Deep Neural network
9.....	אתחול פרמטרים
10.....	אופטימיזציה
10.....	Momentum method
11.....	Dropout
12.....	RBM – Restricted Boltzman machine
14.....	Gibbs Sampling

למידה עמוקה

הרצאה 1

למידה עמוקה הוא תחום מתוך למידת מכונה.

דוגמא: נתון $x = (x_1, x_2, x_3)$, $y \in \{0, 1\}$

אנו מחפשים מכונה לינארית $\hat{y} = a_1x_1 + a_2x_2 + a_3x_3 = \vec{a}\vec{x}$

נגדיר פונקציית מחיר: $S(\vec{a}) = \sum (\hat{y} - y)^2$ נרצה להביא למינימום. במקרה כזה ישנה פונקציה

סגורה כיוון שצמצמו את עצמנו למודל לינארי:



$$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \cdot \begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

עבור סיווג בינארי נבחר את המשערך $\hat{y} = \text{sign}(\vec{a}\vec{x})$ כלומר נקבל קו החלטה לינארי.

במקרה הלינארי קל למצוא פתרון, מציאת המשערך הלינארי הטוב ביותר היא פשוטה. הבעיה היא במישור הלא-לינארי שבו נרצה לעבוד, הוא ענק ויהיה קשה למצוא את המסווג (המכונה) הטוב ביותר.

טענה: כל פונקציה לא-לינארית ניתנת לקירוב ע"י רשת נוירונים. Logistic neural network.

הדרכים למציאת המסווג הטוב ביותר:

שנת 1986: Back propogation

שנת 1992: הומצא מודל יותר פשוט ה-SVM ונזנח ה-Neural Network

שנת 2006: חזרה ל-Neural network והנושא קיבל שם חדש Deep learning.

החל מ-2009 התברר כי deep learning מוכיחה את עצמה יותר משיטות אחרות.

הסיבות לרלוונטיות המחודשת של למידה עמוקה:

(1) שיפור כוח החישוב בצורה דרמטית

(2) גידול דרמטי במידע האימון (Training Data)

Single neuron/ logistic regression

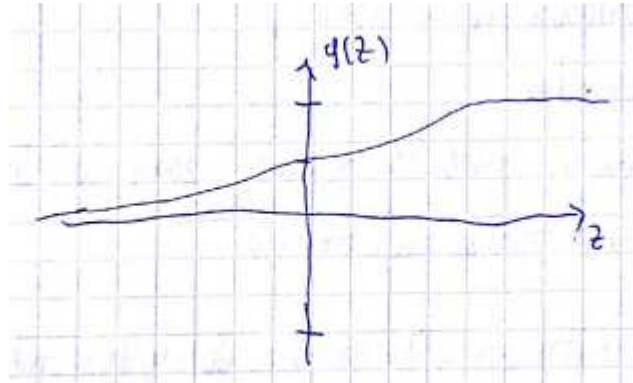
כניסת מידע (Training data): $X_1, \dots, X_n \in R^d$

מוצא (Labels): $y_1, \dots, y_n \in \{0, 1\}$

נחליט כי המסווג יסומן \hat{y} כך שמתקיים $\hat{y} = \begin{cases} 1 & \vec{a}\vec{x} > 0 \\ 0 & \vec{a}\vec{x} < 0 \end{cases}$ ונחפש את המסווג הטוב ביותר. הבעיה

שלנו עם מסווג זה היא אי-הרציפות שלו סביב 0. לכן נבצע "החלקה" של פונקציית המטרה ע"י למשל

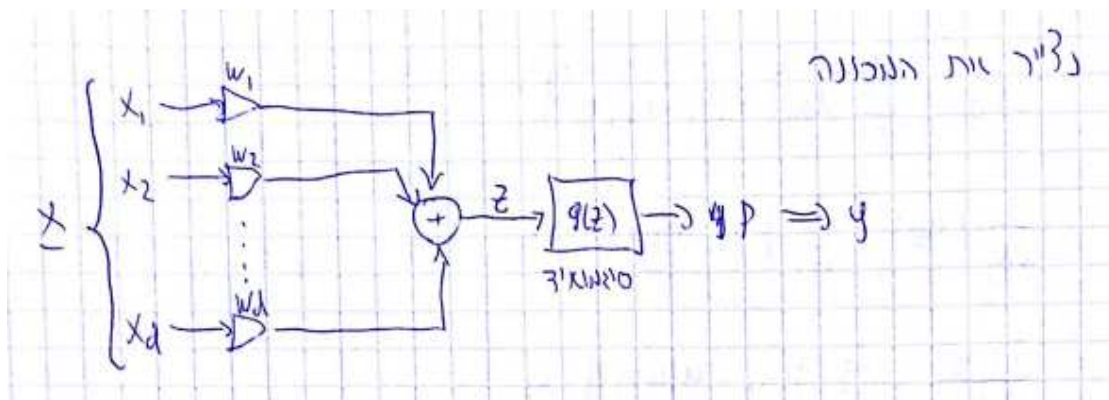
$$\text{sigmoid}(z) = \frac{1}{1+e^{-z}} = g(z), \quad \lim_{z \rightarrow \infty} g(z) = 1 \quad \lim_{z \rightarrow -\infty} g(z) = 0 : \text{Sigmoid}$$



כעת נבנה מסווג כך שבהינתן הוקטור x נכפיל במשקלים w ונכניס לסיגמואיד כך שייתקבל:

$$z = w_1 x_1 + \dots + w_d x_d$$

$$P(y=1|x) = g(z) = \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-\vec{w}\vec{x}}} \quad P(y=0|x) = 1 - g(z) = \frac{e^{-\vec{w}\vec{x}}}{1+e^{-\vec{w}\vec{x}}}$$



נגדיר $L(\vec{w}) = \sum_i \log\{p(y_i | x_i, w_i)\}$ ונחפש w כך L מקסימלי כלומר $\hat{w} = \arg \max \{L(\vec{w})\}$.

לאחר שנמצא w אופטימלי נקבל מכונה שמוציאה soft-decision על סמך x חדש (לא ידוע)

נבצע כמה חישובי עזר:

$$\begin{aligned} \frac{d}{dz} \log[g(z)] &= 1 - g(z) & \frac{d}{dz} \log[1 - g(z)] &= -g(z) \\ \frac{\partial}{\partial w_i} \{\log[P(y=1|x, w)]\} &= \frac{\partial}{\partial w_i} \{\log[g(wx)]\} = [1 - g(wx)]x_i \\ \frac{\partial}{\partial w_i} \{\log[P(y=0|x, w)]\} &= \frac{\partial}{\partial w_i} \{1 - \log[g(wx)]\} = -g(wx)x_i \end{aligned}$$

כך שלסיכום נקבל:

$$\vec{w}, \vec{x} \in R^d, \quad \vec{w}\vec{x} \in R, \quad g(\vec{w}\vec{x}) \in R$$

$$L(\vec{w}) = \sum_t \log\{p(y_t | \vec{x}_t, \vec{w})\} \quad \frac{\partial}{\partial w_j} L(\vec{w}) = \sum_t [y_t - g(\vec{x}_t \vec{w})] x_{tj}$$

$$\frac{\partial L(\vec{w})}{\partial \vec{w}} = \sum_t [y_t - g(\vec{x}_t \vec{w})] \vec{x}_t = \sum_t [y_t - p(y_t = 1 | \vec{x}_t, \vec{w})] \vec{x}_t = \sum_t [y_t - \hat{y}_t] \vec{x}_t$$

נשים לב כי $\frac{\partial L(\vec{w})}{\partial \vec{w}}$ הוא וקטור.

כעת כדי למצוא את w האופטימלי (כלומר להביא את הנגזרת ל-0) נשתמש בשיטת **Gradient Ascent**:

כאשר ε נקרא קצב הלמידה (learning rate) $\vec{w}_{t+1} = \vec{w}_t + \varepsilon \frac{\partial}{\partial \vec{w}} L(\vec{w}_t) = \vec{w}_t + \varepsilon \sum_t [y_t - \hat{y}_t] \vec{x}_t$

ניתן להוכיח (יש בדפים) כי $L(w)$ הינה פונקציה קעורה ולכן עם מקסימום בודד. בסיום התהליך נקבל w אופטימלי ואז קו החלטה לינארי כמו שראינו קודם.

הרצאה 2

Multi-Class logistic regression

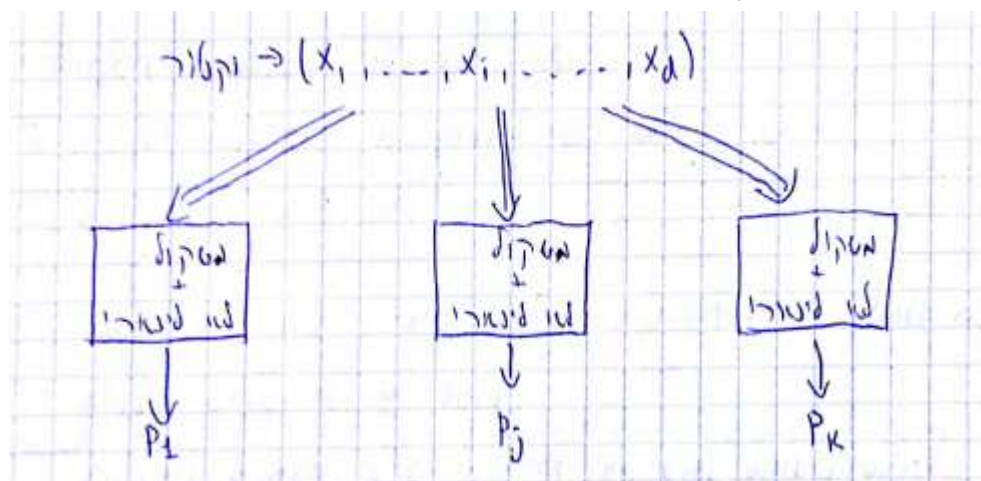
במקרה הזה ל- γ יש יותר משתי אפשרויות $\{0,1\}$:

Input : $x_1 \dots x_n \in R^d$ Output : $y_1 \dots y_n \in \{1, \dots, k\}$

$$P(y = k | x) = \frac{e^{\vec{w}_k \vec{x}}}{\sum_j e^{\vec{w}_j \vec{x}}} = \text{soft max}(w_1 x, w_2 x \dots w_k x) \quad \text{ואז:}$$

זאת לעומת : $\arg \max \{p(y = i | x)\} = \arg \max \{w_i x\}$

המודל הגרפי של המערכת:



תהליך החלטה:

(1) לוקחים x עם d רכיבים

(2) מכפילים את הוקטור x במשקלים w_j ומקבלים k ערכים

(3) מכניסים את k הערכים ל $p(y = i | x)$ שראינו קודם, ונקבל התפלגות על k ערכים.

(4) מחליטים על הערך של γ לפי המקסימום מביניהם.

נגדיר פונקציית מחיר $S(w) = \sum_t \log[p(y_t | x_t, w)]$ שוב נחפש w אופטימלי (כלומר שמביא

פונקציה זו למקסימום), כלומר נגזור את ההסתברות לקבל j לפי הוקטור w_i :

:

$$\begin{aligned} \frac{\partial}{\partial \vec{w}_i} \log \left(\frac{e^{\vec{w}_j \cdot \vec{x}}}{\sum_j e^{\vec{w}_j \cdot \vec{x}}} \right) &= \frac{\partial}{\partial \vec{w}_i} [\vec{w}_j \cdot \vec{x} - \log(\sum_j e^{\vec{w}_j \cdot \vec{x}})] = \delta_{ij} \cdot \vec{x} - \frac{e^{\vec{w}_j \cdot \vec{x}}}{\sum_j e^{\vec{w}_j \cdot \vec{x}}} \vec{x} \\ &= \left(\delta_{ij} - \frac{e^{\vec{w}_j \cdot \vec{x}}}{\sum_j e^{\vec{w}_j \cdot \vec{x}}} \right) \vec{x} = [\delta_{ij} - p(y = i | x, w)] \vec{x} \end{aligned}$$

כאשר δ_{ij} הוא 1 רק אם $i=j$.

כעת באותה דרך נגזור את פונקציית המחיר ונקבל:

$$\frac{\partial S}{\partial \vec{w}_i} = \sum_t [1_{\{y_t=i\}} - P(y_t = i | x_t)] \vec{x}_t$$

כאשר $1_{\{y=i\}}$ הוא 1 רק אם $y=i$ ואפס במקום אחר (פונקציית

אינדיקטור). גם כאן ניתן להוכיח כי קיים מקסימום גלובלי אחד.

וכעת באותה דרך נבצע Gradient Ascent רק שעכשיו יש לבצע K פעמים עבור כל אחד מווקטורי ה- w .

הערה: במקרה הבינארי ראינו כי ניתן להשתמש בוקטור w אחד. נראה איך זה יוצא מכאן:

$$P(y = 1 | x) = \frac{e^{w_1 x}}{e^{w_1 x} + e^{w_0 x}} = \frac{1}{1 + e^{-(w_1 - w_0)x}} = \text{sigmoid}[(w_1 - w_0)x]$$

מקרה רציף

במקרה הרציף נשתמש בשיטה אחרת: **Linear Regression** (בניגוד ל Logistic regression שראינו עד כה):

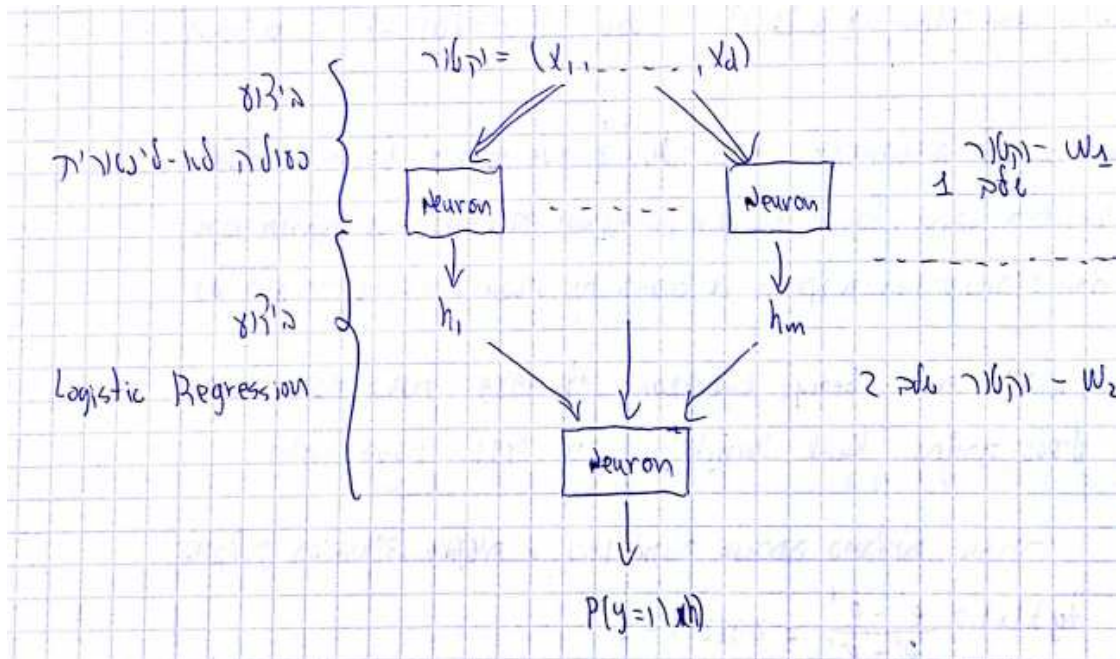
$$\text{Input} : x_1 \dots x_n \in R^d \quad \text{output} : y_1 \dots y_n \in R$$

$$\hat{y} = \vec{w} \cdot \vec{x} \quad S(\vec{w}) = \sum_t (\vec{w} \cdot \vec{x}_t - y_t)^2$$

קיימת כאן הנחה מובלעת כי γ מתפלג נורמלית עם תוחלת w_x .

Neural network with 1 hidden layer

נשים לב כי המסווג שראינו לעיל הוא לינארי וכמו שאמרנו זוהי קבוצה קטנה של מסווגים. אנו רוצים לחזור לעולם הלא לינארי. הדרך לעשות זאת היא להעביר את x דרך פעולה לא-לינארית ואז להשתמש במסווג לינארי כמו שראינו למעלה, אם נחזור למקרה הבינארי הרשת תיראה כך:



כך ש: $P(y=1|x) = \frac{1}{1+e^{-\bar{w}_2 \cdot h}}$ כאשר $h_i = \frac{1}{1+e^{-\bar{w}_1 \cdot x}}$ קבוצת הווקטורים בשלב הראשון

הדבר הזה נקרא רשת נוירונים

מודל רשת הנוירונים

- (1) שלב הפעולה הלא-לינארית בו משתמשים ב- LR אך לא כמסווג
- (2) שלב הסיווג בו ניתן להשתמש בכל מסווג (בינארי, Multiclass וכו')

משפט הקירוב האוניברסאלי

תהי $\varphi(x)$ פונקציה רציפה מונוטונית עולה וחסומה לדוגמא סיגמואיד. ותהי $f(x)$ פונ' רציפה אז לכל

$$\varepsilon \text{ קיים } \hat{f}(x) = \sum_i a_i \varphi(w_i x + b_i) \text{ כך ש } |f(x) - \hat{f}(x)| < \varepsilon$$

במילים אחרות כיוון שהשלב הראשון ברשת הנוירונים מורכב מסיגמואידים אזי הניורון האחרון הוא \hat{f} , ומכאן שבעזרת רשת נוירונים ניתן ליצור קירוב לכל פונקציה שנרצה.

$$L(w) = \sum_i \log P(y_i | \bar{x}_i, \bar{w}) \text{ נגדיר את המודל.}$$

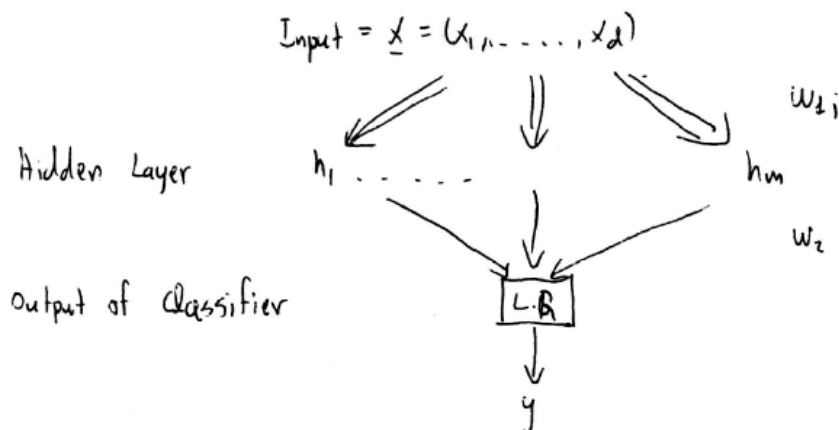
הבעיה היא שכעת יש לנו הרבה נקודות מקסימום מקומיות למראות זאת עדיין נשתמש ב -

$$\text{Gradient Ascent, כלומר } w_{t+1} = w_t + \varepsilon \frac{\partial L}{\partial w}(w_t)$$

הבעיה היא שצריך לעבור על כל הרשת כל פעם עבור כל אחד מהפרמטרים (כל איבר בכל וקטור w) לכן המציאו שיטה מהירה יותר שעושה הכל במכה (כלומר גוזרת לפי כל ה-w בכל מעבר על הרשת)

הרצאה 3

תזכורת ברשת נוירונים עם שכבה נסתרת אחת:



$$h_i = \frac{1}{1 + e^{-(\vec{w}_{1i} \vec{x} + b_{1i})}} \quad p(y=1 | h) = \frac{1}{1 + e^{-(\vec{w}_2 \vec{h} + b_2)}}$$

כעת שוב נחפש w ו- b המוצאים מקסימום לפונקציית L בעזרת Gradient ascent:

$$L(w, b) = \sum_t \log[P(y_t | \vec{x}_t, \vec{w}, b)]$$

$$w_{t+1} = w_t + \varepsilon \frac{\partial S}{\partial \vec{w}}(w_t) \quad b_{t+1} = b_t + \varepsilon \frac{\partial S}{\partial b}(b_t)$$

כעת כמו שהזכרנו, זה לא יעיל לגזור עבור כל איבר בכל אחד מווקטורי ה- w בנפרד לכן נלמד שיטה:

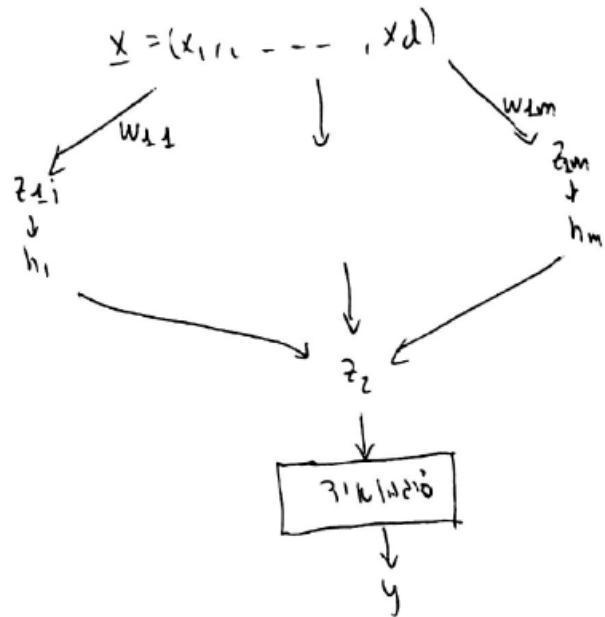
Back propagation algorithm

Feed forward

נעבור על הרשת מלמעלה למטה כרגיל כך ש:

$$z_{1i} = \vec{w}_{1i} \vec{x} + b_1 \quad z_2 = \vec{w}_2 \vec{h} + b_2$$

$$h_i = g(z_{1i}) \quad \hat{y} = p(y=1 | x, w, b) = g(z_2)$$



Back propagation

כפי שראינו בהרצאה הראשונה: $\frac{\partial S}{\partial z_2} = (y - \hat{y})$ ולכן

$$\frac{\partial S}{\partial w_2} = \frac{\partial S}{\partial z_2} \frac{\partial z_2}{\partial w_2} = (y - \hat{y})h = [y - p(y=1 | x, w, b)]h$$

$$\frac{\partial S}{\partial b_2} = \frac{\partial S}{\partial z_2} \frac{\partial z_2}{\partial b_2} = (y - \hat{y})$$

כעת עבור השלב הראשון:

$$\frac{\partial S}{\partial w_{1i}} = \frac{\partial S}{\partial z_2} \frac{\partial z_2}{\partial h_i} \frac{\partial h_i}{\partial z_{1i}} \frac{\partial z_{1i}}{\partial w_{1i}} = (y - \hat{y})w_{2i}g'(z_{1i})\bar{x} = (y - \hat{y})w_{2i}h_i(1 - h_i)\bar{x}$$

$$\frac{\partial S}{\partial b_{1i}} = \frac{\partial S}{\partial z_2} \frac{\partial z_2}{\partial h_i} \frac{\partial h_i}{\partial z_{1i}} \frac{\partial z_{1i}}{\partial b_{1i}} = (y - \hat{y})w_{2i}h_i(1 - h_i)$$

דוגמא נוספת – Non-linear regression

Input : $x_1 \dots x_n \in R^d$ Output : $y_1 \dots y_n \in R$

Feed forward

$$h_i = \frac{1}{1 + e^{-(\vec{w}_{1i}x + b_{1i})}} \quad \hat{y} = \vec{w}_2 \vec{h} + b_2$$

$$S = (y_i - \hat{y}_i)^2 \quad L(w, b) = \sum_t (y_t - \hat{y}_t)^2 \quad \text{נגדיר:}$$

Back propagation

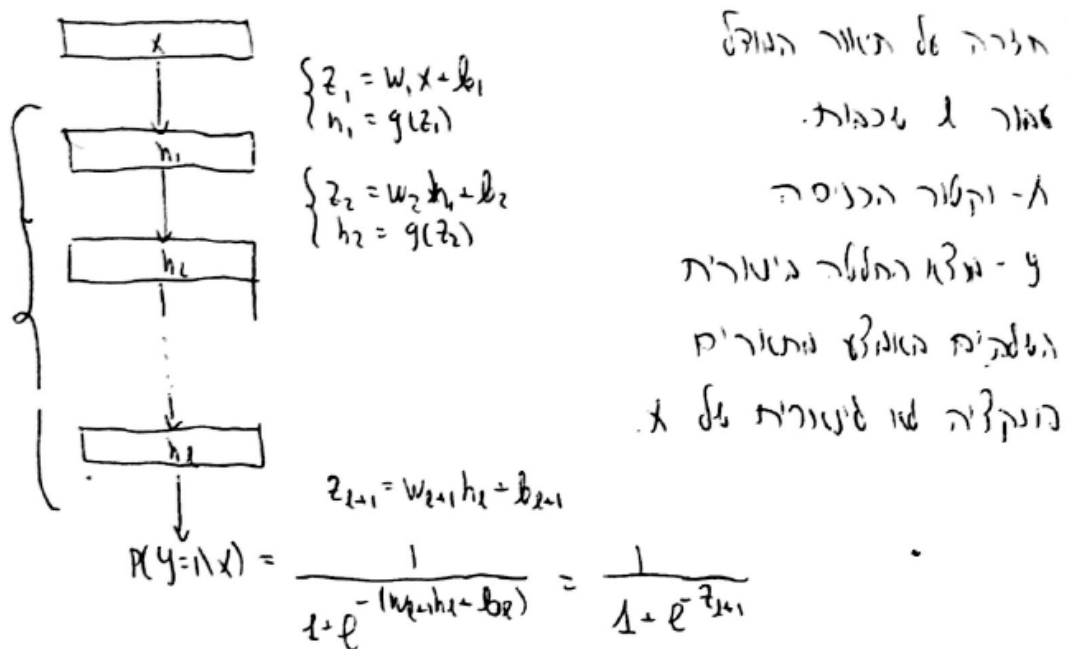
$$\begin{aligned}\frac{\partial S}{\partial \hat{y}} &= -2(y - \hat{y})h = 2(\hat{y} - y) & \frac{\partial S}{\partial w_2} &= \frac{\partial S}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_2} = 2(\hat{y} - y)h & \frac{\partial S}{\partial b_2} &= \frac{\partial S}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial b_2} = 2(\hat{y} - y) \\ \frac{\partial S}{\partial h} &= \frac{\partial S}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h} = 2(y - \hat{y})w_2 & \frac{\partial S}{\partial z_{1i}} &= \frac{\partial S}{\partial h_i} \frac{\partial h_i}{\partial z_{1i}} = [2(y - \hat{y})w_2]g'(z_{1i}) = [2(y - \hat{y})w_2]h_i(1 - h_i) \\ \frac{\partial S}{\partial w_{1i}} &= \frac{\partial S}{\partial z_{1i}} \frac{\partial z_{1i}}{\partial w_{1i}} = [2(y - \hat{y})w_2]h_i(1 - h_i)x \\ \frac{\partial S}{\partial b_{1i}} &= \frac{\partial S}{\partial z_{1i}} \frac{\partial z_{1i}}{\partial b_{1i}} = [2(y - \hat{y})w_2]h_i(1 - h_i)\end{aligned}$$

תהליך ה- Feed forward אינו בהכרח לינארי
תהליך ה Back propogation הינו תמיד לינארי

הרצאה 4

Deep Neural network

כשאנו מדברים על DNN הכוונה היא לרשתות נוירונים עם יותר מ-hidden layer אחד:



גם כאן משתמשים בדיוק באותו FF ו BP- שראינו ב- Single layer אך כעת נשים לב כי W בנגזרות הינו וקטור ולא סקלר כמו מקודם. כמו כן נזכיר כי h_L ו w_{L+1} הינם בעלי אותו מימד
ראה כל הנוסחאות בדפים המצורפים להרצאה 3

אתחול פרמטרים

עד כאן דיברנו רק על הדרך המהירה למצוא W אופטימלי, כעת נדבר על אתחול.

תחילה נשים לב שברשת נוירונים רגילה אם משקולות מסוימים יהיו זהים גם המוצא יהיה זהה, לכן נרצה להגריל תנאי התחלה כך שתאים שונים יקבלו אתחול שונים. כמו כן יש לעדכן את המשקלים כך שערכי הסיגמואיד ייפלו באיזור הלינארי של הסיגמואיד (בו הנגזרת גבוהה יחסית) ולכן נרצה להגריל באיזור הרלוונטי (קרוב יחסית ל-0) **כך ש:**

$$w \sim U\left[-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\right] \quad \text{כאשר } n \text{ מספר האבירים בווקטור } x, \text{ וכן } b = 0 \text{ (אם כי גם אותו מומלץ}$$

להגריל אקראית בין 0 ל-1)

מאותה סיבה **נרצה לנרמל את X** כך שגם יהיה קטן (למשל אם כל איבר בווקטור X מקבל ערכים בין 0 ל-1 אזי נחלק את X ב-n).

אופטימיזציה

$$S(\vec{w}) = -\frac{1}{n} \sum_i \log[p(y_i | x_i, w)]: \text{Cross Entropy}$$

$$S(\vec{w}) = \frac{1}{n} \sum_i S_i \quad \text{כך ש} \quad S_i(\vec{w}) = -\log[p(y_i | x_i, w)]$$

את הפונקציה הזו נרצה להביא למינימום לכן נעבוד עם **Gradient Descent**:

$$w_t = w_{t-1} - \varepsilon \frac{\partial S}{\partial w}(w_{t-1})$$

כעת נגדיר **SGD (Stochastic gradient descent)**:

$$w_t = w_{t-1} - \varepsilon \frac{\partial S_i}{\partial w}(w_{t-1})$$

הדבר הזה יותר מהיר מאשר לעבור על כל הטור כל פעם.

נגדיר שיטות בהקשר הזה:

Full batch: גוזרים כל פעם לפי כי מידע הלמידה

Online: גוזרים כל פעם רק לפי יחידת מידע אחת (אקראית)

Mini batch: מסדרים את המידע בקבוצות בגודל M ואז עבור כל קבוצה M מבצעים Full batch.

הערה: שיטת Online היא Mini batch עבור M=1. ושיטת Full batch הינה שיטת Mini batch עבור M=n.

הערה נוספת: כאשר עובדים עם Mini batch חשוב "לערבב" את המידע לפי שמחלקים לקבוצות כך שלא נקבל בקבוצה מסוימת תווית אחת (כלומר מידע מסוג מסוים אחד) אלא נקבל מכל התוויות.

Learning rate

אין חוקים ברורים פה, עושים לפי מה שעובד. ניתן אולי רק לומר שככל שמתקרבים למקסימום\מינימום אפשר להקטין אותו

Momentum method

בשיטת GA או GD אנחנו זזים בכיוון הגרדיאנט אך ייתכן שהגרדיאנט לכיוון המינימום מקסימום הוא קטן ולכן ייקח לנו זמן להתכנס אליו. שיטה זו אמורה לעזור לנו.

נגדיר $\Delta_t = -\varepsilon \frac{\partial s}{\partial w}(w_{t-1}) + \alpha \Delta_{t-1}$ כאשר $\alpha \sim [0,1]$ נקרא מקדם המומנטום הדבר הזה מאפשר

לנו למצוא את "הכיוון הנכון" של הגרדיאנט יותר בקלות זאת כי כיוון שאנו מצרפים את הנגזרות מזמנים קודמים. כך שהגרדיאנט בכיוון הנכון, גם אם הוא קטן, הוא קבוע ולכן מצטבר, ואילו הגרדיאנטים הלא נכונים אולי גדולים אך בכיוונים משתנים ולכן מתבטלים עם הזמן (אם מצרפים אותם).

$$w_t = w_{t-1} - \varepsilon \frac{\partial s}{\partial w}(w_{t-1}) + \alpha \Delta_{t-1} \quad \text{כלומר נקבל}$$

הדבר הזה כמובן עוזר לאוסלציות אך עדיין לא יימנע מאתנו להיכנס למינימום מקומי.

Nesteror momentum method

$$w_t = w_{t-1} + \alpha \Delta_{t-1} - \varepsilon \frac{\partial s}{\partial w}(w_{t-1} + \alpha \Delta_{t-1}) \quad \text{שינוי קל לשיטה הקודמת כך ש:}$$

הדבר נובע מכך שאני יודע שבכל מקרה לפי שיטת המומנטום אני ממילא זז לכיוון $w_{t-1} + \alpha \Delta_{t-1}$ לכן אפשר כבר מראש לעשות את הנגזרת שם.

5 הרצאה

Dropout

הרעיון כאן הוא שבזמן הלמידה אנחנו מאפסים באופן אקראי חלק מהניורונים (חלק מה-h) ניתן טיפה רקע:

נניח כי נתונים לנו מספר מסווגים כך ש:

$$P_{j1} = P(y=1|x) \quad \text{משתמש במסווג j.}$$

אפשרויות למיצוע:

$$P(y=1|x) = \frac{1}{k} \sum_j P_{j1} \quad \text{מיצוע חשבוני:}$$

$$\frac{1}{k} \sum_j \log \left(\frac{P_{j1}}{P_{j0}} \right) \quad \text{או בכתוב לוגריתמי} \quad \frac{P(y=1|x)}{P(y=0|x)} = \left(\prod_j \frac{P_{j1}}{P_{j0}} \right)^{\frac{1}{k}} \quad \text{מיצוע גאומטרי:}$$

נגיד שיש לנו m כמות של h, אזי כל אחד מהם מקבל ערך של 0 או 1 כך שמקבלים 2^m סוגים של רשתות (בכל אחת מהרשתות מאפסים h אחרים) וכל פעם מאמנים אחת מתוך הרשתות האלה מה שייקרה בסוף הוא שנקבל 2^m מסווגים ובסופו של דבר בזמן Testn פשוט נעשה מיצוע

$$P_A(y=1|x) = \frac{1}{1 + e^{-(\sum w_{2i} h_i a_i + b_2)}} \quad P_T(y=1|x) = \frac{1}{1 + e^{-(\frac{1}{2} w_2 h + b_2)}}$$

כלומר בזמן הלמידה כל פעם מאפסים h אחר, כלומר מכפילים במשתנה אקראי a שהינו 1 או 0 עבור כל אחד מה-h ובזמן Testn פשוט עושים מיצוע על ידי הכפלה בחצי.

נשים לב כי בשיטה זו כאשר עושים Back propagation אין צורך לבצע גזירה עבור המשקלים של ה- h שאותם איפסנו (הנגזרת תהיה שווה לאפס פשוט), נשים לב כי מדובר הן על המשקלים שהולכים אל אותו ה- h וכן על המשקלים שיוצאים ממנו לעבר y .
(ראה דפים של dropout לנוסחאות המעודכנות)

(ראה הוכחה לנכונות ה- Dropout בדפים של אלון הרצאה חמש דף אחרון)
סיכום: בכל מעבר על הרשת (בין אם עבור כל מידע או mini-batch וכו') נאפס חלק מהנוירונים כך שתתקבל תת רשת אקראית. נמשיך כך עד שנעבור על כל האפשרויות של הרשתות האקראיות (ייתכן שיידרשו כמה שימושים בכל המידע בשביל זה).
בסיום נקבל סט פרמטרים להם ניתן פקטור $1/2$ ובהם נשתמש בזמן ה-Test
הדבר הזה שקול לאימון על 2^m רשתות שונות

הרצאה 6

RBM – Restricted Boltzman machine

אלגוריתם זה קשור לנושא ה- Unsupervised pre-training , מטרתו היא למצוא אתחול בצורה מושכלת לפרמטרים.

משמעות ה- Unsupervised היא שאנחנו משתמשים ב- Data בלבד ללא Labels. בדור"כ המידע זמין

יותר לעומת labels שיותר קשה למצוא. נניח $\vec{x} \in \{0,1\}^n$

נתבונן בפונקציית אנרגיה $E: \vec{x} \rightarrow \mathbb{R}$, אם נמיר למודת הסתברותי נקבל:

$$P(\vec{x}) = P(x_1, \dots, x_n) = \frac{e^{-E(\vec{x})}}{\sum_{\vec{x}} e^{-E(\vec{x})}}$$

Partition factor

$$P(\vec{x}) = P(x_1, \dots, x_n) = \frac{e^{-7}}{\sum_{\vec{x}} e^{-7}} = \frac{1}{2^n} \quad \text{אז} \quad E(x_1, \dots, x_n) = 7$$

נניח X מתפלג אחיד, ונניח $E(x_1, \dots, x_n) = 7$

המכנה נובע מכך שיש 2^n אפשרויות עבור הוקטור x . כלומר לפונקציית אנרגיה יש 2^n התאמות, אחת עבור כל אפשרות של X .

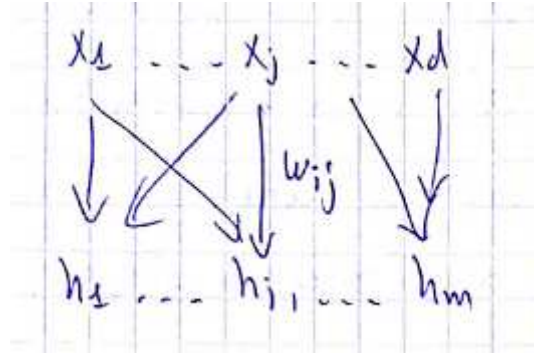
נתבונן על פונקציית אנרגיה ריבועית: $E(\vec{x}) = \sum_{ij} w_{ij} x_i x_j - \sum_i b_i x_i$ הדבר הזה נקרא הסתברות

בולצמן (Boltzman machine) כעת נחלק את הוקטור x לשתי קבוצות: $\vec{x} = (\vec{x}, \vec{h})$ וייתקבל:

$$E(\vec{x}) = E(\vec{x}, \vec{h}) = -\sum_{ij} w_{ij} h_i x_j - \sum_j a_j x_j - \sum_i b_i h_i = -\vec{h}^T \vec{w} \vec{x} - \vec{a} \vec{x} - \vec{b} \vec{h}$$

זהו מודל RBM, כך ש $m+d=n$ ו $h \in \{0,1\}^m$ $x \in \{0,1\}^d$

בצורה גרפית זה ייראה כך:



כלומר דואגים לכך שבחלוקה ל- x ו- h לא תתקבל הכפלה של x ב- x ושל h ב- h , לכן זה נקרא restricted.

$$P(x, h) = \frac{e^{-E(\vec{x}, \vec{h})}}{\sum_{x, h} e^{-E(\vec{x}, \vec{h})}} = \frac{e^{\vec{h}^T \vec{w} \vec{x} + \vec{a} \vec{x} + \vec{b} \vec{h}}}{\sum_{x, h} e^{\vec{h}^T \vec{w} \vec{x} + \vec{a} \vec{x} + \vec{b} \vec{h}}} \quad \text{כעת}$$

$$P(h | x) = \prod_{i=1}^m P(h_i | x) \quad \text{טענה: } P(h_i = 1 | x) = \text{sigmoid}(\sum_j w_{ij} x_j + b_i)$$

הוכחה:

$$P(\vec{h} | \vec{x}) \stackrel{\text{Bayes}}{=} \frac{P(\vec{x}, \vec{h})}{P(\vec{x})} = \frac{e^{-E(\vec{h}, \vec{x})}}{\text{Const}} \propto e^{-E(\vec{h}, \vec{x})}$$

נשים לב כי הקבוע במכנה הוא בגלל ש-

x נתון.

$$e^{\vec{h}^T \vec{w} \vec{x} + \vec{a} \vec{x} + \vec{b} \vec{h}} = \exp(\sum_i \vec{h}_i \vec{w}_i^T \vec{x} + b_i h_i) = \prod_i \exp[(\vec{h}_i \vec{w}_i^T \vec{x}) + \vec{b}_i \vec{h}_i] \quad \text{כעת}$$

כאשר את

$$P(h | x) = \prod_{i=1}^m P(h_i | x) \quad \text{הוכחנו ש-} P(h_i | x) \text{ הוא קבוע (} x \text{ נתון ו-} a \text{ וקטור קבוע כלשהו).}$$

נשים לב כי הצלחנו לפרק את פונקציית צפיפות ההסתברות למכפלה של גורמים שונים ולכן אותם גורמים הם בת"ס, כך שעד כדי נרמול קיבלנו את ההתפלגות הבאה:

$$P(h_i | x) = \exp(\vec{h}_i \vec{w}_i^T \vec{x} + \vec{b}_i \vec{h}_i)$$

$$P(h_i = 0 | x) = e^0 = 1$$

$$P(h_i = 1 | x) = e^{\vec{w}_i^T \vec{x} + b_i} \quad \text{אם ניזכר כי } h \in \{0, 1\} \text{ נקבל כי:}$$

ונזכיר שהסתברויות אלה ללא נרמול (לכן סכומן גדול מ-1)

$$\frac{P(h_i = 1 | x)}{P(h_i = 1 | x) + P(h_i = 0 | x)} = \frac{e^{\vec{w}_i^T \vec{x} + b_i}}{1 + e^{\vec{w}_i^T \vec{x} + b_i}} = \frac{1}{1 + e^{-(\vec{w}_i^T \vec{x} + b_i)}} \quad \text{לאחר נרמול נקבל כי:}$$

כלומר קיבלנו **סיגמואיד!** בדיוק כמו ב-Logistic regression. והוכחנו את הטענה.

לסיכום חלק זה: אם נסתכל על המודל ההסתברותי בצורה גרפית אזי בהינתן פונקציה אנרגיה של x ו- h ומודל RBM, אזי ההסתברות של h בהינתן x היא סיגמואיד, בדיוק כמו ברשת נוירונים

כעת נחשב את ההסתברות השולית של x כאשר $1/z$ גורם הנרמול:

$$\begin{aligned} P(\vec{x}) &= \frac{1}{Z} \sum_{\vec{h}} e^{\vec{h}\vec{w}\vec{x} + \vec{a}\vec{x} + \vec{b}\vec{h}} = \frac{1}{Z} e^{\vec{a}\vec{x}} \sum_{\vec{h}} e^{\vec{h}\vec{w}\vec{x} + \vec{b}\vec{h}} = \frac{1}{Z} e^{\vec{a}\vec{x}} \sum_{\vec{h}} e^{\sum_i \vec{h}_i \vec{w}_i \vec{x} + b_i h_i} \\ &= \frac{1}{Z} e^{\vec{a}\vec{x}} \sum_{\vec{h}} \left(\prod_i e^{\vec{h}_i \vec{w}_i \vec{x} + b_i h_i} \right) = \frac{1}{Z} e^{\vec{a}\vec{x}} \prod_i \left(\sum_{h_i \in \{0,1\}} e^{\vec{h}_i \vec{w}_i \vec{x} + b_i h_i} \right) = \\ &= \frac{1}{Z} e^{\vec{a}\vec{x}} \prod_i \left(1 + e^{\vec{w}_i \vec{x} + b_i} \right) \end{aligned}$$

(חסר פה את המשך ההרצאה שקשור לאימון RBM ראה 2 עמודים אחרונים ברצאה שש אצל אלון או בדפים, בגדול זה פחות רלוונטי למבחן)

הרצאה 7

המטרה שלנו במודל RBM היא למצוא את הפרמטרים w ו- b ולהשתמש בפרמטרים אלו לטובת Neural network. נדרש כמובן לאתחל את מודל ה-RBM

Gibbs Sampling

נרצה לדגום את המ.א. (x, h) דהיינו לדגום n דגימות מ- x ו- m דגימות מ- h כך שהם ייתפלגו לפי

$$\begin{aligned} P(x, h) &= \frac{1}{Z} e^{-E(x, h)} \quad \text{נזכיר כי במקרה שלנו} \\ P(h | x) &= \prod_i P(h_i | x) \rightarrow P(h_i = 1 | x) = g\left(\sum_j w_{ij} x_j + b_i\right) \\ P(x | h) &= \prod_j P(x_j | h) \rightarrow P(x_j = 1 | h) = g\left(\sum_i w_{ij} h_i + a_j\right) \end{aligned}$$

כמו כן:

כעת מה שנעשה הוא שבהינתן x_1 כלשהו ראשוני (נבחר אותו אקראית) נוכל להגריל את h_1 לפי פונקציית צפיפות הפילוג המותנה ומתוך אותו h_1 להגריל את x_2 לפי פונקציית צפיפות הפילוג המותנה השנייה וכן הלאה מ- x ל- h :

$(x_1 \rightarrow h) \rightarrow (x_2 \rightarrow h_2) \rightarrow \dots (x_t \rightarrow h_t)$ כך עבור t גדול מספיק $(x_t, h_t) \sim P(x, h)$.

כעת בהינתן שזורקים את הדגימות הראשונות (שאינן מתפלגות עדיין לפי פונקציית הפילוג, למשל את 100 הראשונות) אזי משם והלאה כל הדגימות שניקח ניתן להניח כי הן מתפלגות לפי P

אם נרצה למשל לחשב תוחלת $\eta[f(x, h)] = \frac{1}{T} \sum_t f(x_t, h_t)$ כאשר t גדול מספיק

או אם למשל נרצה לחשב תוחלת של איבר ספציפי $\eta(x[3]) \approx \sum_t x_t[3]$

נניח כעת שנתונים לנו n דגימות של x כך ש: $x_1, \dots, x_n \in \{0, 1\}$ אזי:

מתוכה את $P(x,h)$ ומתוכה את $S(\vec{w}, \vec{a}, \vec{b}) = \frac{1}{n} \sum_t \log P(x_t; w, a, b)$ ידועים ניתן למצוא התפלגות w, a, b אם

$P(x)$ כמו שמצאנו להלן. כעת בעזרת Gibbs sampling ניתן לדגום $X-h$ כלומר, בהינתן מידע למידה

$$x_t \text{ אזי } h_{t1,i} \sim p(h_i | x_t) \rightarrow x_{t2,j} \sim P(x_j | h_{t1})$$

ומתוכם ה- X שהתקבלו לחשב את S , וכמובן שנרצה למקסם את S ע"י w, a, b ונעשה זאת ע"י

Gradient Ascent

(ראה נגזרות ואלגוריתם בדפים המצורפים להרצאה 6)

כלומר נעשה כך:

- (1) נגדיל w, a, x התחלתיים (כמו ברשת נוירונים)
- (2) נגדיל h בהינתן x ואז נגדיל x בהינתן h
- (3) נעשה נגזרות לפי הנוסחאות ולפי ה- X וה- h שמצאנו ונתקדם לפי גרדיאנט אסנט
- (4) נחזור לצעד 2 ונמשיך כמה שרוצים

כלומר מדובר פה על אימון של רשת RBM, לאחר שנאמן אותה נמצא למעשה w, b המאפשרים להגיע לייצוג טוב יותר של X ואותם ניקח אתנו אל רשת הנוירונים.

לסיכום: RBM הינה שיטת Pre-training בה אנו רוצים לאתחל את הפרמטרים w, b בצורה אופטימלית.

הקשר בין RBM לרשת נוירונים:

- (1) בהינתן וקטור X
- (2) מוצאים w, a, b בעזרת RBM
- (3) מוצאים $P(x, h)$
- (4) מוצאים $P(h_i = 1 | x) = g(\sum_j w_{ij} x_j + b_j)$ למעשה חזרנו לרשת נוירונים

הערה: בהינתן כמה שכבות של רשת בהינתן X נגדיל h ואז h הוא ה- x של השכבה השנייה וכן הלאה.

כמו כן נזכיר כי ברשת RBM $h \in \{0, 1\}$ לעומת רשת נוירונים בה $h \in R$.

השימוש ב-RBM נועד כדי למצוא פרמטרים התחלתיים טובים יותר עבור w, a, b . הדבר שימושי במיוחד כאשר יש לנו מעט מידע ללמידה (כלומר מעט מידע עם labels)