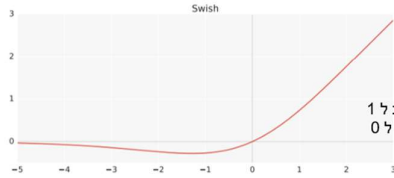


Swish Self Gated Activation

$$\sigma(z) = z \cdot \text{sigmoid}(z)$$

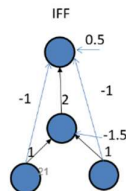
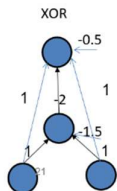
$$\text{Sigmoid}(z) = \frac{1}{1 + \exp(-z)}$$



- איזו סוג רשת היא הכוללת 2 ניוונים הראשון מחובר לשני והשני מחובר לראשון. האם ניתן לתת קלט לרשת כזו?
• RNN, כן: כל ניוון יכול לקבל עוד קלטים
- אילו סוגי רשתות ניתן לתאר כ DAG?
• FF
- איזו תכונות ביאולוגיות מבטאת בד"כ פונקציית האקטיבציה:
• קצב ירי כפונקציה של סכום הזרמים הנכנסים בדנדרטים
• הסתברות לירי כפונקציה של סכום הזרמים

– בנה רשת שמבצעת IFF

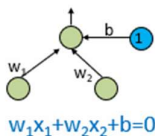
- בעזרת 4 ניוונים נסתרים (משפט הבניה)
- בעזרת 2 ניוונים נסתרים
- בעזרת ניוון נסתר בודד (כולל מעקף מהקלט ישירות לפלט)



שימו לב, הניוון הנסתר לא חייב לשנות את תפקידו. כדי להפוך את הפלט: משקולות הפלט מונפלות ב -1

הוכחה אלגברית:

לא ניתן לממש הפרדה לינארית ל XOR



00 → 0
01 → 1
10 → 1
11 → 0

– הוכחה בשלילה, נניח שניתן:

אזי קיימות משקולות שמפרידות את 4 הנקודות בסבלת האמת נציב 4 הנקודות במשוואת ה BTU, נקבל 4 אילוצים סותרים:

$$\begin{aligned} 0w_1 + 0w_2 + b &\leq 0 \\ 0w_1 + 1w_2 + b &> 0 \\ 1w_1 + 0w_2 + b &> 0 \\ 1w_1 + 1w_2 + b &\leq 0 \end{aligned}$$

$$\begin{aligned} w_1 + w_2 &\leq -b; w_1 > -b; w_2 > -b; b \leq 0 \\ w_1 + w_2 &\leq -2b \quad \text{וגם} \quad w_1 + w_2 > -2b \quad \leftarrow \text{סותר!!} \end{aligned}$$

Spike Time Dependent Plasticity-STDP

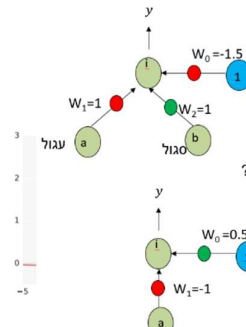
- נתונים 2 ניוונים A ו B. A מחובר דרך סינפסה A לניוון B ואילו ניוון B מחובר ל A דרך סינפסה B. מה יקרה לסינפסות אם ניוון B יורה וניוון A לא יורה? מה יקרה לסינפסה B אם ניוון A יורה וניוון B לא יורה?
• A תחליש B יתחזק
- באזור ה Paladus Ganglia מצאו כי כלל Anti Hebb עובד. מה יקרה לסינפסה בין 2 ניוונים שם אם הם יורים יחדיו?
• יחליש
- לניוון בקליפת המוח יש סף שלילי (ביאס חיובי) וסינפסות אינהיביטוריות והמסוגלות ליצר אך ורק זרמים שליליים (אינהיביטורים). הניוון יורה במידה וסך הזרמים המתקבלים בדנדרטים גדול מהסף. מה יקרה לניוון אם לא יתקבל שום קלט – z.a. שום זרם לא יוזקק לדנדרטים? מה יתקבל אם יתקבל זרם בעוצמה גוברת והולכת מהדנדרטים?
• הניוון יורה בקצב קבוע עד שיקבל זרם מהדנדרטים, ואז יקטין את קצב הירי עד להשתתות
- מה יקרה לניוון אם יגיעו זרמים דנדרטים אך סכום הזרמים קטן מהביאס?
• ימשיך לירות אבל בקצב יותר נמוך
- מה יקרה לניוון אם סכום הזרמים בערך מוחלט גדול מהביאס?
• הניוון ישתתק
- מה יקרה בניוון ביאולוגי אקסיטטורי כאשר השני (post) יורה לפני הראשון (pre)? מה יקרה לניוון מלאכותי כאשר מקטנים משקל משקל עוד ועוד?
• חוזק הסינפסה תיחלש עד שתגיע ל 0 (לא צעביה זרם) אולם לא תתפר לאינהיבטוריות
• בניוון מלאכותי, חוזק הסינפסה תרד תחצה שטח 0 ותתפר שלילית יוצר ויותר

Swish Self Gated Activation

ניתן לבנות שערים לוגיים בעזרת רכיבי BTU

(בדרך הנדסית)

• שער AND: $y = a \text{ AND } b$



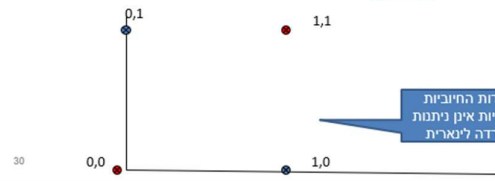
– לציפור ניוון שמגיב לצבע סגול וניוון שמגיב לצורה עגולה היא אוכלת רק דברים עגולים וסגולים (אונמניה-כן, כדור גולף-לא) איזה מעגל יציב יממש ניוון שמגיב לדברים אכילים?

• שער NOT: $y = \sim a$

Linear Inseparability of XOR Geometric View

• במרחב קלט דו מימדי, האם הדוגמאות ניתנות להפרדה לינארית?

00 → 0
01 → 1
10 → 1
11 → 0



הנקודות החיוביות והשליליות אינן ניתנות להפרדה לינארית

בקלט 2D, (יש 3 משקלים) כמה מלקרים בלתי תלויים לינארית צריך לכל הפחות על מנת שנוכל לבנות קבוצת אימון שאיננה ניתנת להפרדה? תשובה: 4

אלגוריתם: Perceptron Learning

- עבור על כל דוגמאות הקלט שבקבוצת האימון שוב ושוב (עד שאין שינוי)
1. חשב את y (החזוי) של הפרספטרון על קטור הקלט x_j
 2. חשב את השגיאה $(t_j - y)$
 3. דחף את קטור המשקולות כך שהשדה המשותף יתקדם לכיוון הנכון

$$\Delta w_j = \lambda(t_j - y_j)x_j$$

מה ההיגיון?

1. אם z_j קטן שווה 0, והמטרה $t_j=1$, השגיאה חיובית. אנחנו מתקנים את המשקל כך ש z יגדל בכיוון הנכון. אם הקלט חיובי מגדילים את המשקל ואם הקלט שלילי מקטנים.
 2. אם z_j גדול מ 0, והמטרה $t_j=0$, השגיאה שלילית. מתקנים את המשקל ומקטנים את z בכיוון הנכון. אם הקלט חיובי, מקטנים את המשקל ואם הקלט שלילי, מגדילים.
 3. האם התהליך יתכנס? האם לא יתכן שנתקרב למטרה בדוגמה אחת ונתרחק בדוגמה אחרת?
- הזנבולט הוכיח שאם יש פיתרון, האלגוריתם ימצא אותו, אם אין פיתרון, האלגוריתם לא יתכנס אף פעם ואף תיתכן התבדרות של המשקולות במה תלויה הצלחת הלמידה?
- אם דוגמאות האימון ניתנות להפרדה ליניארית!
- א.א. אם הקלט מכיל features "טובים" (הניתנים להפרדה ליניארית)

האם כלל הלמידה של פרספטרון מזכיר לנו את הכלל ההביאני?

$$\Delta w_{ij} \sim y_j x_i$$

למידה הביאנית מתבטאת בשינוי משקולות הסינפסות:

– הקשר מתחזק במידה 2 ו מיונים יורים ביחד

$$\Delta w_{ij} \sim -y_j x_i$$

למידה אנטי-הביאנית:

– הקשר נחלש אם המיונים יורים ביחד

$$\Delta w_{ij} \sim (t_j - y_j)x_i$$

האם כלל הפרספטרון הוא הביאני?:

- המשקולות תשתנה רק אם החזוי y לא זהה ל t
 - את כאשר הקלט x_i חיובי ו $y=0$, המשקל יתחזק (הביאני)
 - כאשר הקלט שלילי ו $y=1$ אזי המשקל יחלש (אנטי הביאני)
 - כאשר הקלט חיובי ו $y=1$, המשקל יחלש (אנטי הביאני)
 - כאשר הקלט שלילי ו $y=0$ המשקל יתחזק (הביאני)

38

Gradient Descent Algorithm for minimizing a loss function

- נתונה פונקציית עלות $loss_{D,h}(w_0, w_1, \dots)$ שהגרדיינט שלה ניתן לחישוב
- נתונה קבוצות אימון D עם דוגמאות: $\langle x_j, t_j \rangle$
- רוצים למצוא (w_0, w_1, \dots) שימזער: $Min_{w_0, w_1, \dots} \{loss_{D,h}(w_0, w_1, \dots)\}$
- העיקרון:

- התחל מ w_0, w_1, \dots משקולות אקראיים
- בעצ $Epochs$ שוב ושוב עד להתקיימות של תנאי העצירה.
- למשל: עד שהעלות מפסיקה לרדת, או שהגענו למקסימום $Epochs$

בכל Epoch:

- עבור D , חשב את הגרדיינט של פונקציית העלות בנקודה w
- שנה את w_0, w_1, \dots כך ש $loss(w_0, w_1, \dots)$ יקטן: חישוב השינוי יראה צעד כנגד כיוון השיפוע

$$\Delta w_i = -\lambda \frac{\partial loss_{D,h}(w)}{\partial w_i} = -\lambda \sum_{p \in D} \frac{\partial loss_D(y_p)}{\partial y_p} \frac{\partial h_w(x_p)}{\partial w_i} \quad y = h(w(x))$$

$$w_i = w_i + \Delta w_i$$

השיטה רלוונטית גם לפונקציות במימד גבוה וגם לפונקציות שאינן קמורות

למידת פרספטרון: המשך תרגיל AND

- נתון פרספטרון בעל 3 גורמי קלט: x_0, x_1, x_2 , בצע לשידה של הפונקציה AND בין 3 הקלטות כך ששיוון הפלט יהיה 1 אם כל הקלטות הם 1.
- המשקולות מתחילים ב $(-1.2, -0.3, -0.4)$ ויבצע הלמידה היא 0.9 .
- כלל העדכון:

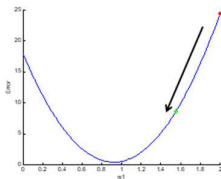
$$\Delta w_j = 0.9(t_j - y_j)x_j$$

	x_0	x_1	x_2	w_0	w_1	w_2	output	desired output
	1	1	0	-1.2	-0.3	-0.4	0	0
	1	0	1	-1.2	-0.3	-0.4	0	0
	1	0	0	-1.2	-0.3	-0.4	0	0
הוספת (0.9,0.9,0.9)	1	1	1	-1.2	-0.3	-0.4	0	1
הפחתת (0.9,0.9,0)	1	1	0	-0.3	0.6	0.5	1	0
	1	0	1	-1.2	-0.3	0.5	0	0
הוספת (0.9,0.9,0.9)	1	0	0	-1.2	-0.3	0.5	0	0
הפחתת (0.9,0.9,0)	1	1	1	-1.2	-0.3	0.5	0	1
	1	1	0	-0.3	0.6	1.4	1	0

צעד נוסף: עדכון המשקולות ע"פ ממוצע הגרדיינטים בנקודה החדשה

$$\Delta w_1 = -\lambda \frac{\partial loss}{\partial w_1} = -0.01 * 46.4 = -0.464$$

$$w_1^{new} = w_1 + \Delta w_1 = 2 - 0.464 \approx 1.54$$

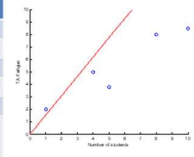


Gradient Descent-Single var with no bias

$$\Delta w_1 = -\lambda \frac{\partial loss}{\partial w_1} \quad \text{נבצע צעד GD}$$

$$\frac{\partial loss}{\partial w_1} = \frac{1}{m} \sum_{p=1}^m (y_p - t_p)x_p \quad \text{ערך התחלתי } w_1 = 2$$

y_p	t_p	y_p	$(y_p - t_p)x_p$
10	8.5	$2 * 10 = 20$	$(20 - 8.5) * 10 = 125$
5	3.8	$2 * 5 = 10$	$(10 - 3.8) * 5 = 31$
1	2	$2 * 1 = 2$	$(2 - 2) * 1 = 0$
8	2	$2 * 8 = 16$	$(16 - 8) * 8 = 64$
4	5	$2 * 4 = 8$	$(8 - 5) * 4 = 12$
$\frac{1}{m} \sum_{p=1}^m (y_p - t_p)x_p$:			46.4

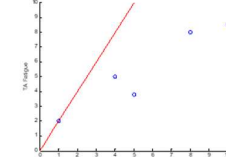


Gradient Descent-Single var with no bias

$$\Delta w_1 = -\lambda \frac{\partial loss}{\partial w_1} \quad \text{הנגזרת החלקית לפי } w_1$$

$$= -\lambda \frac{\partial loss}{\partial w_1} = \frac{\lambda}{m} \sum_{i=1}^m (t_p - y_p)x_p$$

נגדיר $\lambda = 0.01$ ונתחיל מ $w_1 = 2$



Stochastic GD: מדוע צריך לשפר את GD?

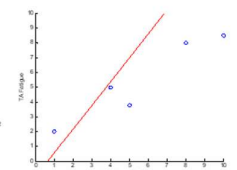
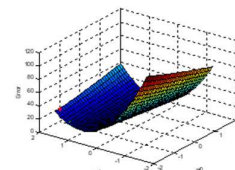
בשיטת GD מבצעים EPOCHS שוב ושוב בכל EPOCH: מחשבים גרדיינט לכל דוגמה ובסוף ממוצע של הגרדיינטים כדי לבצע צעד שינוי משקולות קטן אחד

- GD הוא אלגוריתם off-line: משתמשים בכל ה training set (Epoch) כדי לבצע צעד אחד ב Batch אם קבוצת האימון גדולה, צריך למצע המוני גרדיינטים לפני שעושים צעד בודד
- On-line Gradient Descent: זוראצייה של GD ל on-line - מבצעת צעד לאחר כל דוגמה

עדכון המשקולות וערך השגיאה

$$\Delta w_0 = -\lambda \frac{\partial loss}{\partial w_0} = -0.01 * 4.74 = -0.0474 \quad \Delta w_1 = -\lambda \frac{\partial loss}{\partial w_1} = -0.01 * 38.8 = -0.388$$

$$w_0^{new} = w_0 + \Delta w_0 = -1 - 0.0474 \approx -1.05 \quad w_1^{new} = w_1 + \Delta w_1 = 2 - 0.388 \approx 1.61$$



כלל ה LMS של SGD מול כלל הלמידה של פרסטרון

$$\Delta w_i = \frac{\lambda}{m} \sum_{p \in D} (t_p - y_p) x_{ip}$$

GD ב
וב mini-batch

$$\Delta w_i = \lambda (t_p - y_p) x_{ip}$$

SGD ב

במה שונה מכלל הפרסטרון?

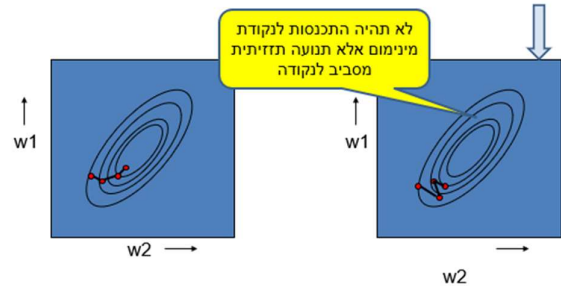
- γ לינארי (ולא פונקציית מדרגה)
- GD משתמש בממוצעים ואילו פרסטרון עובד דוגמא דוגמא
- Batch vs. on-line
- SGD מקרב ירידה עם גראדיינט שממוצע פונקציית Loss.
- יש משמעות למינימום גם כשאין הפרדה לינארית מלאה
- למידת פרסטרון לא תעזור אם אין הפרדה לינארית

בפרסטרון, y היא פונקציית המדרגה
ברגרסיה y לינארי
On line של רגרסיה לינארית נהיה דומה לכלל הפרסטרון. ההבדל הוא: איך מחושב y
הירידה היא לא steepest אלא כל דוגמא, מושכת לכיוון שלה (יותר רועש, ואיטי)

Online (SGD) vs Batch (GD) walking in loss space

פונקציית השיגאה בנירן לינארי היא קערה n מימדית- מינימום יחיד

- GD נתקדם במאונך לקווי הגובה של הקערה.
- SGD נתקדם בזיג-זג.



O/L SGD vs. GD יתרונות וחסרונות

יתרונות: SGD

- מהיר יותר כאשר יש מיליונים של דוגמאות באימון.
- כל דוגמה מקדמת אותנו ולא צריך לסכם מיליונים לפני עדכון.
- כאשר יש כמה מינימומים מקומיים, התזזיתיות עוזרת לפעמים לברוח ממנימום מקומי
- ניתן להשתמש ב on-line learning (המשך הלמידה גם במצב production)

חסרונות:

- מסלול תזזיתי אינו מביל ישירות למטרה (ב GD המסלול חמדני- הכי תלול),
- אין התכנסות (יש הסתובבות אינסופית סמוך למינימום)

תרגיל: איטרציה אחת של משקולות ב SGD
 $m=1$, עבור סיגמואיד בודד במרחב קלט 1D.

רוצים לסווג גידולים על פי גודלם.

נניח שהתחלנו ממשקולות: $w_1=-2, w_0=1$

קצב הלמידה $\lambda=1$

בהינתן דוגמה $\langle x=2, t=1 \rangle$

כיצד ישתנו המשקולות?

$$\Delta w_i = -\lambda \frac{\partial C}{\partial w_i} = \lambda(t - y)x_i$$

$$y = g(w_1 x_1 + w_0) = g(-2 \cdot 2 + 1) = g(-3) = 1/(1+e^3) \approx 0.05$$

$$\Delta w_0 = 1(t - y)1 = (1 - 0.05) = 0.95$$

$$\Delta w_1 = 1(t - y)x_1 = 1(1 - 0.05)2 = 1.9$$

$$w_0 = w_0 + \Delta w_0 = 1 + 0.95 = 1.95$$

$$w_1 = w_1 + \Delta w_1 = -2 + 1.9 = -0.1$$

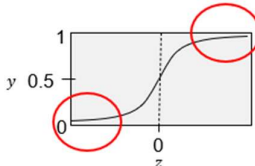
58

בעיית הגרדיאנטים הנעלמים

$$\frac{\partial MSE}{\partial w_i} = -\frac{1}{m} \sum_p y^p (1 - y^p) (t^p - y^p) x_i^p$$

הגרדיאנט כופל את הפונקציה הלוגיסטית עצמה y^p וכן את הביטוי $1 - y^p$.

מה יקרה כאשר y^p קרוב לאפס או אחד??



מה יקרה לפרוצדורת ה Gradient Descent
מדוע ב CE לא נעלמים הגראדיינטים?

גזירת ה- CEntropy

$$y = g(z) = \frac{1}{1 + e^{-z}} \quad z = wx = w_0 + \sum_i w_i x_i \quad \text{נוירן לוגיסטי:}$$

$$\frac{\partial y}{\partial z} = \frac{\partial g(z)}{\partial z} = y(1-y)$$

$$\frac{\partial z}{\partial w_i} = x_i$$

CEntropy על דוגמא אחת:

$$C(y, t) = -t \ln(y) - (1-t) \ln(1-y)$$

$$\frac{\partial C}{\partial y} = -\left(\frac{-t}{y} + \frac{1-t}{1-y}\right) = \frac{-t(1-y) + y(1-t)}{y(1-y)} = \frac{-t+y}{y(1-y)}$$

$$\frac{\partial C}{\partial w_i} = \frac{\partial C}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial w_i} = \frac{(y-t)}{y(1-y)} y(1-y) x_i = (y-t) x_i$$

$$\text{On-line SGD: } \Delta w_i = -\lambda \frac{\partial C}{\partial w_i} = \lambda(t - y)x_i$$

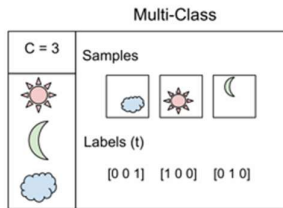
$$\text{Mini-Batch SGD: } \Delta w_i = \lambda/m \sum_{p \in D} \lambda(t_p - y_p)x_{ip}$$

עבור SGD/GD נרצה לחשב גראדיינט של ה
loss כאשר $y = g(wx+b)$

• אם משתמשים ב MSE מקבלים:

- פונקציה שאינה קמורה (יתכנו הרבה נקודות מינימום)
- גראדיינטים שנעלמים
- אין הצדקה ביאזיבית ל MSE עבור קלסיפיקציה

Multi-Class Cross Entropy



$$MCC(y, t) = - \sum_i^{C=k} (t_i \log(y_i) + (1 - t_i) \log(1 - y_i))$$

עבור דוגמא בודדה

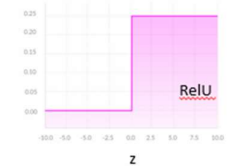
$$MCCE(y, t) = - \frac{1}{m} \sum_p \sum_i t_{pi} \log(y_{pi}) + (1 - t_{pi}) \log(1 - y_{pi})$$

עבור כל קבוצת האימון:

70

RelU: מהו כלל עידכון המשקולות עבור

יחיד



$$\frac{\partial C}{\partial w_i} = \frac{\partial C}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial w_i} = \frac{(y - t)}{y(1 - y)} \frac{\partial y}{\partial z} x_i =$$

$$\frac{(y - t)}{y(1 - y)} x_i \quad \text{כאשר } z > 0$$

אחרת 0

$$\Delta w_i = \frac{\lambda}{m} \sum_p (t_p - y_p) y_p (1 - y_p) x_{pi}$$

73

בבית: מהו כלל עידכון המשקולות עבור tanh ?

Multi class learning with Softmax units

$$y_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

בשכבת הפלט כל נורון מחשב:

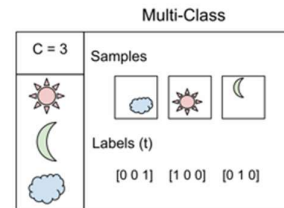
$$\frac{\partial y_i}{\partial z_i} = y_i(1 - y_i) \quad \text{הנגזרת זהה לזו של הנורון הלוגיסטי}$$

שיטה מקובלת לקלסיפיקציה פרמטרית Multi Category:

- מספר הפלטים זהה למספר הסיווגים (קטגוריות). מציגים הסתברויות
- סכום האקטיבציות בשכבת הפלט הוא 1
- מתנהג בדומה לפונקציית MAX (אבל רק יותר)
- כאשר ערך y של יחידת softmax אחת גדל, ערכי שאר היחידות קטנים

80

Variation: Categorical Cross Entropy-



one-hot, רק השגיאה של הקטגוריה הנכונה נלקחת בחשבון. בכל פעם שהחזיון טועה, מגדילים את ההסתברות רק לנורון שאמור להיות 1 מה החיסרון?

$$CCE(y, t) = - \sum_i^{C=k} t_i \log(y_i)$$

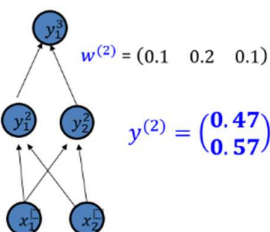
במיקרה הבינארי: יש זהות CE=CCE

$$- \sum_i^{C=2} t_i \log(y_i) = -t_1 \log(y_1) - (1 - t_1) \log(1 - y_1)$$

78

Feed forward – חישוב ערך השכבה הנסתרת

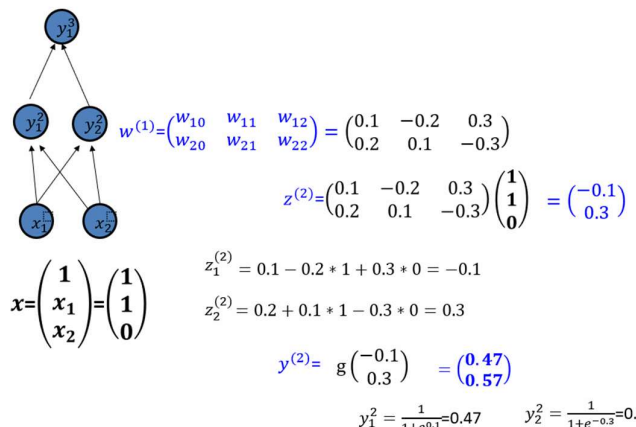
$$y_1^{(3)} = 0.56$$



$$z^{(3)} = (0.1 \quad 0.2 \quad 0.1) \begin{pmatrix} 1 \\ 0.47 \\ 0.57 \end{pmatrix}$$

$$z_1^{(3)} = 0.1 + 0.2 * 0.47 + 0.1 * 0.57 = 0.25$$

$$y_1^{(3)} = \frac{1}{1 + e^{-z_1^{(3)}}} = \frac{1}{1 + e^{-0.25}} = 0.56$$



$$w^{(1)} = \begin{pmatrix} w_{10} & w_{11} & w_{12} \\ w_{20} & w_{21} & w_{22} \end{pmatrix} = \begin{pmatrix} 0.1 & -0.2 & 0.3 \\ 0.2 & 0.1 & -0.3 \end{pmatrix}$$

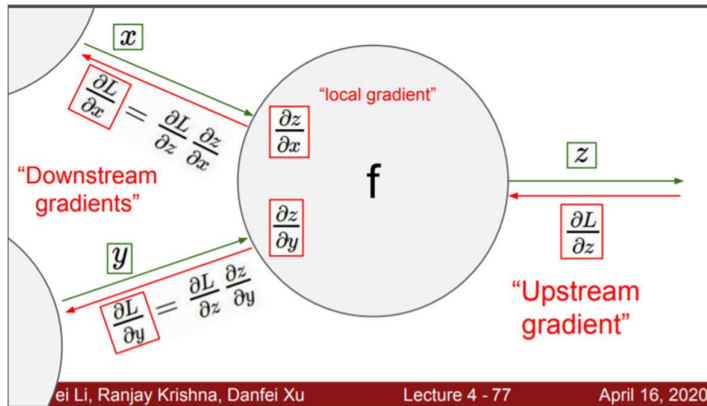
$$z^{(2)} = \begin{pmatrix} 0.1 & -0.2 & 0.3 \\ 0.2 & 0.1 & -0.3 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} -0.1 \\ 0.3 \end{pmatrix}$$

$$z_1^{(2)} = 0.1 - 0.2 * 1 + 0.3 * 0 = -0.1$$

$$z_2^{(2)} = 0.2 + 0.1 * 1 - 0.3 * 0 = 0.3$$

$$y^{(2)} = g \left(\begin{pmatrix} -0.1 \\ 0.3 \end{pmatrix} \right) = \begin{pmatrix} 0.47 \\ 0.57 \end{pmatrix}$$

$$y_1^{(2)} = \frac{1}{1 + e^{-0.1}} = 0.47 \quad y_2^{(2)} = \frac{1}{1 + e^{-0.3}} = 0.57$$



Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Upstream gradient (Local gradient)

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2 x_2)}}$$

Diagram illustrating the forward and backward passes for a sigmoid function.

Forward pass values:

- $w_0: 2.00, x_0: -1.00 \rightarrow -2.00$
- $w_1: -3.00, x_1: -2.00 \rightarrow 6.00$
- $w_2: -3.00, x_2: -2.00 \rightarrow 6.00$
- Sum: $-2.00 + 6.00 + 6.00 = 10.00$
- Sigmoid output: $1 / (1 + e^{-10.00}) \approx 0.98$

Backward pass values (gradients):

- Output gradient: 0.73
- Gradients for x_0, x_1, x_2 : $-0.20, -0.20, -0.20$
- Gradients for w_0, w_1, w_2 : $-0.20, -0.20, -0.20$

Equations shown:

$$\frac{df}{dx} = e^x$$

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$$

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2 x_2)}}$$

Diagram illustrating the forward and backward passes for a sigmoid function.

Forward pass values:

- $w_0: 2.00, x_0: -1.00 \rightarrow -2.00$
- $w_1: -3.00, x_1: -2.00 \rightarrow 6.00$
- $w_2: -3.00, x_2: -2.00 \rightarrow 6.00$
- Sum: $-2.00 + 6.00 + 6.00 = 10.00$
- Sigmoid output: $1 / (1 + e^{-10.00}) \approx 0.98$

Backward pass values (gradients):

- Output gradient: 0.73
- Gradients for x_0, x_1, x_2 : $-0.20, -0.20, -0.20$
- Gradients for w_0, w_1, w_2 : $-0.20, -0.20, -0.20$

Equations shown:

$$f(x) = e^x \rightarrow \frac{df}{dx} = e^x$$

$$f_a(x) = ax \rightarrow \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \rightarrow \frac{df}{dx} = 1$$

AutoEncoders using BP

- Unsupervised learning •
- מטרות: •
- Automatic feature discovery –
- Dimensionality Reduction - דחיסה •
- יצוג אובייקטים דומים בעזרת וקטורים דחוסים
- המקימים "דימיון" וקטורי התואם "דימיון" אפליקטיבי
- Semi-supervised learning –

מדוע להשתמש בקלט OneHot בזבזני ולא ביצוג בינארי מקובל עם הרבה פחות מימדים

- במקום ליצג קלט אובייקט מסוים כ: $(0,0,...,0,1,0,0)$ •
- Unit- לכל אובייקט: דרושים n units •
- מדוע לא ליצג בצורה חסכונית יותר כ $(0,0,1,1)$ •
- דרושים רק $\log(n)$ units תשובה: •
- היצוג הבזבזני הוא בוודאות Linear separable ולכן היצוג הפנימי (המבוזר) בנקל יכול ליצור features רצויים לו •
- כאשר מקצים ביט לכל אובייקט לא מניחים שום ידע קודם. •
- קידוד דחוס יותר חייב להניח דימיון בין חלק מהאובייקטים המיוצגים. •
- היצוג הבזבזני, כל האובייקטים (לא) דומים במידה שווה. •
- האם ניתן ללמוד יצוג דחוס המייצג שבו יצוגים דומים מייצגים אובייקטים דומים? •

AutoEncoders for Unsupervised Learning

- סוג של למידה unsupervised, שמנסה לדחוס את המידע, למצוא ייצוג פנימי ועל ידי כך לגלות features כלליים וטובים. אובייקטים עם מאפיינים דומים יקבלו ייצוג וקטורי דומה

- למידה של פונקציית הזהות: הקלט שווה לפלט הרצוי.
- בגלל צוואר בקבוק בשכבות הנסתרות, אין שינוי של המידע אלא מתגלים features חשובים
- אין מורה: למשל: ניתן לקחת תמונות (או הקלטות קול) ללא תוויות (או ללא תמלול) וללמוד להוציא מהם את המאפיינים החשובים שיאפשרו שיחזור של הקלט לאחר דחיסה

AutoEncoders for Internal representation

הרשת מוצאת ייצוג פנימי מבוצר
: מסבה Onehot encoding לקוד דחוס יותר

לאחר אימון. מקבלים ייצוג נוסתר שונה לכל קלט. מה נקבל אם נעגל את ערכי היחידות הנסתרות ל 0/1?

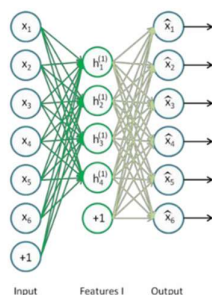
Learned hidden layer representation:

Input	Hidden Values	Output
10000000 →	.89 .04 .08 →	10000000
01000000 →	.01 .11 .88 →	01000000
00100000 →	.01 .97 .27 →	00100000
00010000 →	.99 .97 .71 →	00010000
00001000 →	.03 .05 .02 →	00001000
00000100 →	.22 .99 .99 →	00000100
00000010 →	.80 .01 .98 →	00000010
00000001 →	.60 .94 .01 →	00000001

Auto-encoders

האם ניתן לאמן רשת עמוקה ביעילות?

באמצע שנות האלפיים חיפשו טכניקות יעילות לאמן רשתות עמוקות (הימנעות ממינימום מקומי, האצת זמן האימון)



כאשר יש הרבה שכבות, משטח השגיאה מפותל מאוד, כאשר יש רק שכבה אחת, פחות נתקעים במינימום מקומי.
הרעיון: לאמן כל שכבה בנפרד, לקבל איתחול למשקולות, ואז לעשות Fine tuning

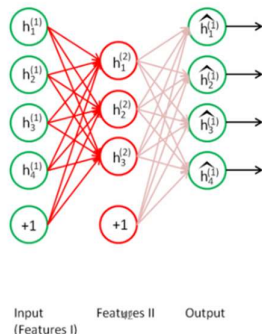
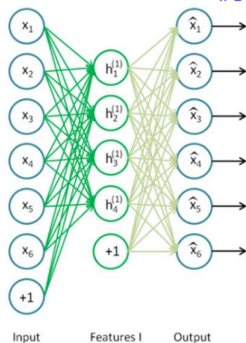
Bottle-neck layers in AutoEncoders

- דחיסה: השכבה האמצעית קטנה- צוואר הבקבוק מכביד על יכולת ההתאמה.
- Sparse Auto-encoder: מכבידים על הרשת בעזרת "עונשים" על האקטיבציה בשכבה הנסתרת (ולאו דווקא באמצעות דחיסה)
- Dropout הכבדה באמצעות הרג של נוירונים בשכבה הנסתרת

Stacked Auto-Encoders, Bengio (2007)

– After Deep Belief Networks (2006)

- שימוש חמדי של הרבה auto-encoders בטור.
- כל auto-encoder הוא בעל שיכבה נסתרת אחת. תוצאות ה encoding של רשת i הם הקלט (והפלט) לרשת i+1
- לאחר שאימנו auto-encoder i, נזרק ממנו את ה decoder ושאר רק אם ה encoder ואותו נשרשר לפלט של ה encoder ה i-1.



Stacked Auto-Encoders for Deep network initialization

- לאחר שיש לנו את שיכבת ה features האחרונה, נסיף שכבת softmax (לקלסיפיקציה) ונאמן Supervised רק את המשקולות בשכבה זו.
- לאחר שאומנו כל השכבות בנפרד, מצרפים את ה encoders של כל שכבה לרשת עמוקה, ועל השכבה העליונה, את שכבת ה softmax
- בסוף, נבצע אימון supervised על כל הרשת (fine-tuning לכל המשקולות)

