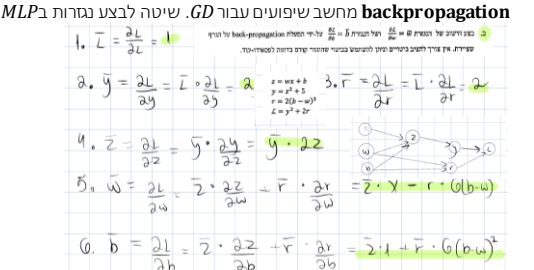


הגרדיאנט (הנגזרת שלילית). פלטי סיגמאיד
לא יכולים להיות 0 (לא זירו כנסטר)
גזרת: $\sigma'(x) = \sigma(x)(1 - \sigma(x))$
tanh טווח $[-1,1]$. הוא כן זירו כנסטר. כבד חישובית. הנגזרת שלילית נורוניהם
רונים הורגים את הגרדיאנט. גזרת: $\tanh'(x) = 1 - \tanh^2 x$
Relu קל לחשב. אינו רווי באזרי $z > 0$ הנגזרת לא מתאפסת. סיבוכיות יעילה.
בעיה שלא זירו כנסטר < לנתון עם מיני בטן. גזרת: $\text{Vrelu}(x) = \mathbb{I}(x > 0)$
Dead RELU עבור כל הקלטים כשהא שלילי ואז הנגזרת 0 וצעד העדכון
יבצע לא עדכון. הבירון לא יהיה פעיל.
Leaky Relu לא רווי כי לא מתאפסת. סיבוכיות יעילה, לא ימות. אין חסרונות.
ELU לעולם לא מת. קרוב ל0 **output**. סיבוכיות **exp**.



המערב הראשון מקלט לזרם
כל חישוב הנגזרות **Backward pass**
החישוב שעליו מדובר ברע הנוכח **Local**
תבניות התנהגות הגרדיאנט:



סנדור ייצוג כללי של ערכים ממדים כלשהו.
אקטיבציה מפעילים על כל איבר בנפרד לאותן ממדים.
נגזרת של סקלר לפי סקלר
נגזרת של סקלר לפי וקטור - הגרדיאנט = וקטור
נגזרת של וקטור לפי וקטור - מטריצה שגודלה $M \times N$
בחיובי סקלר וקטור, וקטור, מטריצה

אופטימיזציות הנדרשות לחישוב אקטיבציה, עיבוד מקדים, אתחול, גזירתיות, חישוב האקטיבציה, קבד למחר, היפר פרמטרים
ערכי סופית - לחישוב K -folds או validation set

batch normalization נרמול הקלט - הופך לקל להתכנסות.
כאשר D פצ'רים N נקודות. ממוצע זה כל בל קורדינטה.
center כלומר מרכזים את המידע.
בנקודת זמן או כשהמידע רחוק מהראשית.
normalized data לשנות את הגדלים

overfitting - מקרב אותם לאזור לינארי של פונקציית אקטיבציה.
2 מונע **exploding** - מקטין את הערכים. 3 מונע **vanishing** - הגרדיאנט של ערכים גדלים בפונקציית אקטיבציה. 4 מזהה התכנסות אימון והופך את האימון לחזק יותר על ידי בחירת היפר פרמטרים ואחלוחל משקולות.
סחמנות: המיני בטן שונה כל פעם, קשה לדבב כי לא יודעים ממוצע. חישוב לזכר את הערכים עם BN והסטטיסטי.
כנסים את FC לחישוב אחרי FC לפי אקטיבציה כי זה מונע חלקה 0.3.
אין למנוע overfitting. 1. שימוש בגזירתיות. 4. **Early stopping** עצירה מוקדמת. 2. **Data arguments** לשלל דוגמאות חדשות בעזרת הגדלת קלט. 3. להפחית את מורכבות הדיאלוג מזהר מזהר מספר העיבוד או היתויבים.
בעלררציה: למנוע מהמודל להסתכל **data** למנוע **overfitting**.
1. **early stopping** למנוע כשמיגיעים **overfitting** לא מפסיקים ללמוד את המודל.
2. **L2 regularization** למנוע שונה על w גדול מדי **Weight decay** עונש מנומרה 2. $L2(w) = \frac{\lambda}{2} \|w\|_2^2 + \sum (w, x^i, y^i)$
הגזרת $\frac{\partial L2(w)}{\partial w} = \lambda w + \nabla L(w, x^i, y^i)$
3. **L1 regularization** למנוע הערכים מהמודל להסתכל **data** למנוע **overfitting**.
יזירה את המאפיינים החשובים ביותר של **data** ומתעלמים מרעשים.

צעד העדכון $\nabla L(x^i, y^i) < 0$ $\eta \nabla L(x^i, y^i)$
SGD כי זה רק נקודה אחת.

sign function 0 ל0. 1- לשלילי. 1 לחיובי.

Multi class לכל מחלקה $t \in \{1, \dots, k\}$ יש מסווג f_1, \dots, f_k
תבינה כללית: $f_j(x) = w_j^T x + b_j$; $y = \arg \max_k f_k$
 $L_p(w, x^i, t^i) = \sum_{j \neq t} \max(0, -f_j(x^i) - f_t(x^i))$
פרספרטור -
עדכון SGD - $w_j \leftarrow w_j - \eta x^i$ $w_t \leftarrow w_t + \eta x^i$
בפועל: להוריד ציונים שבגוהים מהמחלקה ולהעלות את המחלקה. יש נטרלים.

logistic regression $z = w^T x$ $y = \sigma(z) = \frac{1}{1 + e^{-z}}$
שיטה סטטיסטית לרגרסיה בינארית, תוחם את המשוואה לטווח σ . תכונות:
* מה הסיכוי ש1 $t = \sigma(w^T x)$
* $0 < y < 1$
* התבוננות $P_w(t = 1|x) = 1 - \sigma(-w^T x)$
* $P_w(t = 1|x) = \sigma(w^T x)$
* $P_w(t|x) = \sigma(tz)$ where $z = w^T x$ סיכום

Loss שימושי **ML** $L_{ML}(y, t) = -\log P_w(t|x) = -\log \sigma(tz)$
הנגזרת $\nabla L(w) = -\sigma(-tz)tx$
צעד העדכון עם SGD $w \leftarrow w + \eta x \sigma(-t^i z^i) t^i x^i$
מחיר $\epsilon(w) = \frac{1}{n} \sum_{i=1}^n -\log \sigma(t^i z^i)$ $\nabla \epsilon(w) = \frac{1}{n} \sum_{i=1}^n -\sigma(-t^i z^i) t^i x^i$

cross entropy מדד למדידת ההתאמה בין ההסתברויות שניתנות על ידי המודל וההסתברות האמיתית של התוצאות. ערך ההתאמה גדול - שופים. ערך קטן - קרובים.
הפונקציה מחזירה ערך לא שלילי ולכן ערכים גבוהים מצביעים על דרגת טעות גבוהה

עבור $H(p, \hat{p}) = -\sum p(c) \log \hat{p}(c) = -\sum p(c) \log \hat{p}(c|x) = -\log \hat{p}(c|x)$
עבור $L_{CE} = -\sum t \log(p_t)$ **SoftMax** אחרי p_t הוועיות האמיתיות
בביאור $H(p, \hat{p}) = -\sum p(c) \log \hat{p}(c|x) = -\log \hat{p}(c|x)$
 $L_{CE} = -\sum t \log(p_t)$ **SoftMax** אחרי p_t הוועיות האמיתיות
Binary cross entropy $y \in \{0, 1\}$; $\hat{y} = p = \sigma(z)$
אנטרופיה היא יחידת המידע המכונה חוסר סדר. כדי למצוא את פרמטרים טובים, הערך המינימלי של אנטרופיה צריך להיות נמוך. ההפסד מייצג מרחק בין התפלגות הוועיות הרצויה להתפלגות המתקבלת מהמודל. ערך המינימלי הוא 0 - מודל תואם בצורה מושלמת את ההתפלגות הרצויה.

בוחנה ML שווה BCE
Let $t^{(i)}$ denotes the true class $t^{(i)} \in \{1, \dots, l\}$ for inputs $x^{(i)}$.
Then $\mathcal{L}_{CE}(x^{(i)}) = \mathcal{L}_{CE}(t^{(i)}) = -\log(\sigma(t^{(i)})) = -\log(\sigma(-t^{(i)})) = -\log(\sigma(-t^{(i)}))$
The maximum-likelihood (ML) maximizes w s.t.:
 $w^* = \arg \max_w \prod_i P_w(t^{(i)}|x^{(i)}) = \arg \max_w \sum_i \log P_w(t^{(i)}|x^{(i)}) = -\sum_i \log \sigma(t^{(i)} w^{(i)})$
Cross entropy minimizes $ce = -\sum_i \sum_c P(c|x^{(i)}) \log P_w(c|x^{(i)})$ where P is the true distribution and P_w is the approximated distribution.
In our case we have a binary class $t \in \{+1, -1\}$:
 $w^* = \arg \max_w \sum_i -t^{(i)} (1 - \log(\sigma(t^{(i)})) - \log(-1) \log(\sigma(-t^{(i)})) = -\sum_i \log \sigma(t^{(i)} w^{(i)})$
Q.E.D.
KL - מדד להפרש בין שתי התפלגויות הסתברות. אין קשר למודל.
 $KL(P||P_w) = \sum -P(c|x) \log \frac{P_w(c|x)}{P(c|x)} = \sum -P(c|x) \log P_w(c|x) - \sum -P(c|x) \log P(c|x)$
אנטרופיה $H(P)$ היא קבועה לא משפיעה על.

חוסר יציבות נומרית
 $\log 0 \rightarrow -\infty$ $\log \infty \rightarrow \infty$ $\log x < 0$ $\log x > 0$ $\log x = 0$ $x = 1$

SoftMax - ציונים לאחוזים.
 $\text{softmax}(f_1, \dots, f_k) = (\frac{\exp(f_1)}{\sum \exp(f_j)}, \dots, \frac{\exp(f_k)}{\sum \exp(f_j)})$
מבוסס על שכל ערך הוא בין 0 ל1 והמכנה 1.
 $L_{softmax} = -\log P_w(t^i|x^i) = -\log(y_i, y_i) = -\log \frac{e^{y_i}}{\sum e_j}$
נגזרת $\frac{\partial L_{softmax}}{\partial w} = -x^i + \frac{\exp(w, x^i)}{\sum \exp(w, x^j)}$ או $x^i - y_i$
SGD $w_j \leftarrow w_j - \eta (y_j - x_j) x_j^i$ for $j \neq t^i$ $w_t \leftarrow w_t + \eta x_t^i$ for $j = t^i$
בפועל הורה הרבה למי שבגוה ממנו וקצת ליתרילים. העלאה של המחלקה.

בעיות מסוגיות לינארי: אי אפשר להעביר ישיר. פתרון: אפשר לייצג אחרת את הערכים כדי שיוכלו להעביר קוד ולשקור על הצורה הקודמת שנוחר. **10 Mlp**
זה המחלקות. בעזרת אקטיבציה. $f(x) = W_2 \max(0, W_1 x + b_1) + b_2$
 w_2 הפך להיות המסווג. **MLP loss** $\frac{\lambda}{2} \|w\|_2^2 + L(w)$ לכל w ניתן עונש

פונקציות אקטיבציה לא לינאריות
sigmoid $[0,1]$. נשתמש סיבוכיים.
חישוב יקר נורוניהם רוים הורגים את

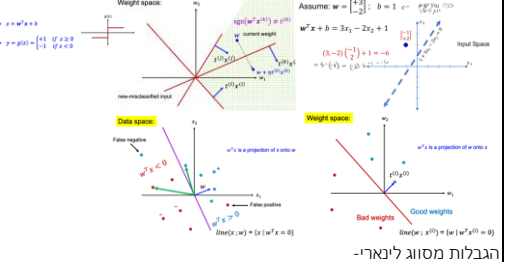
Adam - קומבינציה של RMSProp ו-momentum. $s_k = \frac{\partial L}{\partial w_k}$ נשאר זהה.
 $w_k \leftarrow w_k - \frac{\eta}{\sqrt{s_k} + \epsilon} s_k$ $s_k \leftarrow \alpha s_k + (1 - \alpha) \frac{\partial L}{\partial w_k}$
האדפטיב יותר טוב ממונטום כאשר ההתכנסות רגישה לגודל צעד. בבקורות אופן ובמישורים עם שיפועים קטנים מאוד. בבקורות צרים בעת ריצה על גיא. סיבות להתכנסות איטית:
1. קצב למידה נמוך. חפש קצב למידה טוב יותר **rate decay** ואדפטיב.
2. הרשת הגיעה לגא צר ותלכו בחלל הפרמטרים. ממונטום או אדפטיב.
3. נקודות אופן איזור שטוח במרחב הפרמטרים. ממונטום.

הערכה של w ראשונית:
בעיה 1: אם נאחלחל מספר קטן 0 ו1 לכל היחידות הנסתרות יהיו אות **bias** ואותו משקל ונלך אן אותו שיפוע (נגזרת זהה) ואותם עדכונים למשקלם, נהיה כלאים במרחב הפרמטרים. יחלץ אותם פצ'רים ולא ילמדו להיות שונות.
בעיה 2: אם נאחלחל שללילם כל האקטיבציות יהיו 0 **vanishing gradient** קטן ממסרצת ההזדה - מהר מתכנס ל0. אין עדכון של הרשת, אתחול לא טוב.
בעיה 3: 0.05 נקבל בין 1 ל1 ואין אימון. נאחלח עם מספר רנדומלי סביב 0 כדי שתלקי יהיו חזקים יותר. מסרצת ההזדה מתפוצצת **exploding gradient**. הנגזרות יתפוצצו כי משתמשים **tanh**. הנגזרת 0 ואין אימון.

אתחולים נכונים:
Xavier $\sigma = \frac{1}{\sqrt{D_{in}}}$ ניקח שורש של המימד לפני (השכבה לפני) המטרה = שיטות של הקלט של פלס $\text{var}(x_0) = \text{var}(x_1)$ ככל D_{in} גדול אז שינוי קטן מאוד יפסיק על הוועיות של **output**. המערב הוא $\text{var}(w, x) = D_{in}$ $\text{var}(y) = D_{in}$ וכן $\text{var}(y) = D_{in} \cdot \text{var}(w) \cdot \text{var}(x)$ (באם הוא אפס) $\text{var}(x_1) = \text{var}(w) \cdot \text{var}(x_0)$ נלקח את השונות היחידה של w $\text{var}(w) = \frac{1}{D_{in}}$ והכל מתמצים.

Glorot & Bengio $\sigma = \frac{2}{\sqrt{D_{in} + D_{out}}}$ אול השונות שונים אז עשים ממוצע ביניהם. אתחול בשניהם: $\sigma = \frac{2}{\sqrt{D_{in} + D_{out}}}$ כן לא פוגע.
Relu $\sigma = \frac{2}{D_{in}}$ מקרים שהערכים הם 0 בחיבורי ולכן נשתמש בהו.
None Linear regression $y = x^T w + b$ פרדיקציה מכמה ממדים

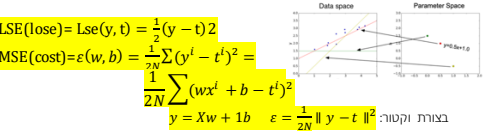
גרסיה פולינומיאלית - מגדירים פצ'רים במקום $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ המודל $y = w^T \phi(x)$ where $w^T = (w_0, w_1, w_2)$ $\phi(x) = \begin{pmatrix} 1 \\ x \\ x^2 \end{pmatrix}$
נבחר מודל שיתן לט ההתאמה טובה (לא יתירה ולא תת התאמה) **bias** טיה מהמכר **Variance** כמה הדגומות מפוזרות.
השיאה: $f = \text{bias}^2 + \text{var}$ ככל שמספר פרמטרים **var** נשאת **Bias** תתר. גדול המודל יעלה את **var** שלילי.



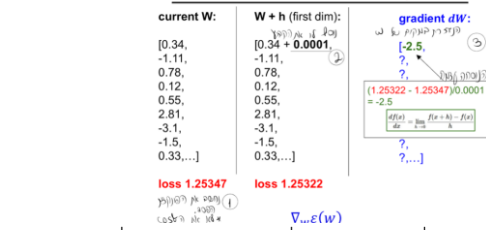
patterns - ניקח ממוצע שפטרן **A**: 0.25 לכל כניסה (הכניסות השחורות). הווקטור הממוצע לכל ההזדה 0.25, 0.25, 0.25... נניח שהוא מצליח לסווג ומה שמקבל ערך גדול מסה **A**. אז, $\sum w_i x_i^{(i)} > 0$, $\text{so then } \sum w_i x_i^{(i)} > 0$.
אפשר לעשות אותו דבר גם על קטן מסה כי ההתבנית הממוצעת של **B** זהה ולכן סתירה. הפתרון פצ'רים לא לינארים בהתחלה ואחר כך רשתות עמוקות.
XOR - עם הטבלה ליצור משוואות ששוות לז
עם האיקסים-סתירה.

losses:
loss 01 $L_{01}(z, t) = \mathbb{I}(tz < 0)$ אם התנאי יתקיים 1 אחרת 0.
כאשר $z = wx + b$ וגם $t = \text{true label}$
 $L_{01}(y, t) = \begin{cases} 0 & \text{if } y = t \\ 1 & \text{if } y \neq t \end{cases}$
העלות $\frac{1}{n} \sum \mathbb{I}(t^i x^i < 0)$ $\epsilon(w) = \frac{1}{n} \sum \mathbb{I}(t^i x^i < 0)$
לפי הגרף הנגזרת שווה ל0 - לא יעדכון w - נגזרת מתאפסת.
perceptron $L_p(w, x, t) = \max(0, -tz)$; where $z = w^T x$
 $t = c$ היא האמית', z הוא הפרדיקציה
הנגזרת $\Delta L_p(w) = -\mathbb{I}(tz < 0)tx$

Linear regression - בעיית רגרסיה לערכים רציפים. נרצה למצוא מודל $y = wx + b$ = הקו הלינארי שיתן לנו מרחק מינימלי מכל הנקודות.



די למער **cost** נמצא אופטימיזציה פרמטרי המודל:
direct solution $\frac{\partial}{\partial w} f(w, x, y) = 0$
שיטה נומרית - קל לביצוע, מהר, אטית, לא מדויק (גשאת עיגול)
Gradients calculated numerically



שיטה אנליטית - גזירת סימבולית. מהיר ומדויק, מיועד לשגיאות.
 $\frac{\partial \epsilon(w, b)}{\partial w} = \frac{1}{N} \sum (x^T w + b - t^i) x^i$ $\frac{\partial \epsilon(w, b)}{\partial b} = \frac{1}{N} \sum (x^T w + b - t^i)$
הבעיה - גרסה לינארית בלבד. **T** אמד מודל (פרמטרים רבים) ובסופו גם אי אפשר לעשות תמיד נגזרות עדיף לרדת בשיפוע.
GD iterative - נקודות מינימום לגלובלי בטורגריקה. קל ליישום ותמיד נותן פתרון גם אם נתקל בלוקלי. המטרה למצוא את פונקציית ההפסד המודדת את ההבדל בין התפוקה החזויה לבין הפועל. ברוב המקרים המודל הוא מורכב ובעל ממדים, עם מינימום מקומי ואופן. מציאת מינימום לגלובלי אינה אפשרית לכן משתמשים ב**GD** להפחית הפסד.

עדכון data שיטות: עוצורים אופטימום נקודות ה0 או לא מתעדכן.
batch * $w \leftarrow w - \frac{\eta}{N} \sum_{i=1}^N (x^i T^i - t^i) x^i$ החזרה $w^{t+1} = w^t - \frac{\eta}{N} \sum_{i=1}^N (x^i T^i - t^i) x^i$
מעבדת את המשקולות לאחר עיבוד כל נקודות התנהגות בכל מחזור אימון (השיפוע מחושב על פני כל נקודות הנתונים בכל איטרציה). סיבוכיות $O(n)$
SGD * $w \leftarrow w - \eta (w^T x^i - t^i) x^i$ החזרה $w^{t+1} = w^t - \eta (w^T x^i - t^i) x^i$ מהיר וחוש.

עבור כל נקודה במערך הוא מעדכן את המשקולות ומחיל תיקון לכיוון השיפוע ותרבות על פני **Batch *** [זמן אימון מהיר יותר כי מעדכנים בצורה קטנה את המשקולות. מאפשר לעבוד עם סט גדול של ערכים. *יעילות טובה יותר עבור מערכי נתונים עם נתונים לא אחידים. **batch** ניתקל ברעש ששיפוע על הכל יאילו עם **SGD** נעבן עבור כל נקודה באופן עצמאי ונמנע מבעיות חסרונות על פני **Batch *** התכנסות איטית ושונות מוגבר יוביל להתנהגות לא יציבה. *דרוש כוונת היפר פרמטר זהיר כדי להשיג ביצועים טובים. סיבוכיות $O(1)$
* **batch *** $w \leftarrow w - \frac{\eta}{B} \sum_{i=1}^B (w^T x^i - t^i) x^i$ החזרה $w^{t+1} = w^t - \frac{\eta}{B} \sum_{i=1}^B (w^T x^i - t^i) x^i$ סיבוכיות $O(n)$

קצב הלמידה: משמש להחלפת גודל צעד העדכון בכל איטרציה.
קטן מדי - לעול להישלל ולעלות על המינימום הגלובלי.
קטן מדי - יתקן הרבה זמן להתכנס, יקר מבחינה חישובית.
1. **עדין ידני** - ניסוי וטעיה ע"י ניחוש.
2. **עדין ידני** - אסטרטגיות בתנומה.
3. **ממונטום** - תנפה. כאשר נאחלח נמצאים בבקו צר ואורך תוך זיגוג ע פני גיא או שנתקנע במינימום לוקלי או אזור שטוח עקב **vanishing** או נקודת אופן (הנגזרת 0) או בהתכנסות איטית כאשר צעד העדכון קטן מדי.
עדכון המומנטום $v \leftarrow \alpha v - \eta \frac{\partial L(w)}{\partial w}$ $w \leftarrow w + v$
 α צריך להיות 0.9 או 0.99 פחות מ1 כי אם יהיה גדול מ1 המומנטום גדל שונה ויכול **exploding** ולא פחות מינימום. קרוב ל0 זה יעצמצם לטסטדרט.
Nestrov - טכניקה נוספת: $v \leftarrow \alpha v - \eta \frac{\partial L(w)}{\partial w}$ $w \leftarrow w + v$ עובקים להשתמש במקום הנכחי כדי לחשב את השיפוע. מחשבת תחילה מקום ביניים על ידי צעד בכיוון המומנטום הקודם וכן לאחר מכן מחיל את המומנטום על השיפוע. תיקון וקטר המומנטום לפי חישוב השיפוע. התכנסות מהירה יותר ולחלח אופטימיזציה יציב.
עדכון אדפטיב - גודל צעד שמתאים לכל וקטור בהתאמה. שיפועים גדולים יקבלו מדרגות קטנות ושיפועים קטנים יקבלו גדולים. עוזר להשיג התכנסות מהירה בבקוים צרים כשמזג מעל גיא. עוזר באזורים עם שיפועים נעלמים או מתפוצצים. **RMSProp** $w_k \leftarrow w_k - \frac{\eta}{\sqrt{s_k} + \epsilon} s_k$ $s_k \leftarrow \beta s_k + (1 - \beta) \frac{\partial L}{\partial w_k}$
אפסילון כדי שלא תחלק באפס. ככל שז גדול - קטן יותר.

