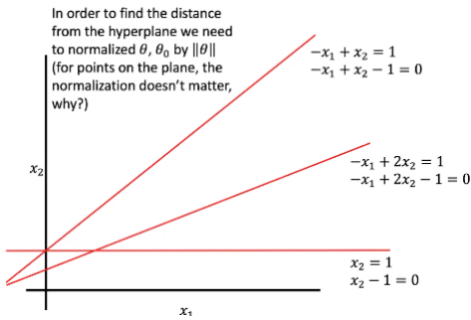


מסגרת תוכנית למידה: Linear Classifiers – 5 תרגול

מפרידים לינאריים בקלסיפיקציה –

- **היפר-מישור במרחב**: במרחב חד ממדי היפר מישור הוא נקודה, במרחב דו ממדי היפר מישור הוא קו ובמרחב תלת ממדי היפר מישור הוא משטח דו-ממדי. נשים לב שכביכול מפריד לינארי הוא ממרחב בממד אחד פחות (בייצוג הגיאומטרי) ממרחב הדאטה בו אנו מתעסקים.
- **הגדרת היפר-מישור במרחב**: מרחב ההיפר-מישור הוא $n-1$ (אם n הוא המרחב בו אנו עובדים). כל הנקודות על ההיפר-מישור פותרות את המשוואה הבאה: $\theta_1 x_1 + \dots + \theta_n x_n = b$ ($= \theta_0$) כאשר x הוא וקטור הקואורדינטות של הדגימה. היפר-מישור מפריד את המרחב לשני מרחבים, כאשר כל נקודה שמניבה תוצאה במשוואה שגדולה מ- b (ה-bias, טטה 0) נמצאת מעל המישור, וכל נקודה שמניבה תוצאה במשוואה שקטנה מ- b , נמצאת מתחת למישור.

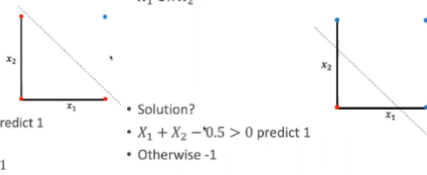


- ה-bias, b , טטה 0 – מייצג את התזוזה של המישור מראשית הצירים, והטטה (הכללית, המקדמים של x) בכלל מייצגת את הזווית של המישור במרחב.
- נשים לב שהדוט-פרודקט = המכפלה הפנימית, בין x לטטה, משמעותה **המרחק של הנקודה מהמישור**. על המישור יש נורמל שמאונך למישור, הדוט פרודקט הוא הטלה של האיקס על הנורמל. **מרחק חיובי משמעותו = מעל המישור, ומרחק שלילי משמעותו = מתחת למישור**. נשים לב שכדי לקבל את המרחק האמיתי על הוקטור טטה להיות נורמל (לחלק אותו בנורמל של עצמו).
- אם הטטה לא נורמלת = המרחק לא מייצג את המרחק האוקלידי מהמישור וכך יש משמעות רק לסימן של המרחק כפי שצוין לעיל.

אנחנו מחפשים מפריד לינארי כך שכל הנקודות המניבות תוצאה גדולה מ-0 יסווגו בקלאס +1 (או -1). וכל הנקודות עם תוצאה קטנה מ-0 יסווגו

לקלאס -1 (או +1). כלומר אנחנו מחפשים טטה $\theta \in R^{n+1}$ (n hyperplane weights & the bias θ_0) , שהיא וקטור משקולות או מקדמים ל- x , שמקיימת את המכפלה הפנימית המתוארת לעיל ומסווגת באופן הבא: $\theta \cdot x = \sum_{i=1}^n \theta_i x_i + \theta_0 > 0$ and -1 otherwise

• X_1 AND X_2



להלן כמה דוגמאות:

אלגוריתם הפרספטרון

אלגוריתם שמחפש מפריד לינארי במרחב, לכן מחפשים טטה באופן הבא: נתחיל עם טטה רנדומית (מישור) ובכל שלב נשפר אם יש צורך בשיפור (כלומר יש שגיאה = אאוטפוט פחות טרגט). השינוי בטטה הרלוונטית היא learning rate כפול גרדיאנט.

כלל העדכון של הפרספטרון

$$\Delta \theta_i = -\eta \sum_{d \in D} (o^{(d)} - t^{(d)}) x_i^{(d)}$$

זהו הכלל לעדכון כך שאם $O_d - T_d = 0$, אין טעות ולכן אין צורך בעדכון.

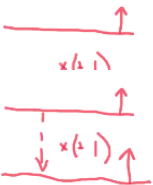
0	t	o-t	x_i	$\Delta \theta_i$	$x_i \cdot \theta_i$
-1	+1	<0	>0	>0	increased
-1	+1	<0	<0	<0	increased
+1	-1	>0	>0	<0	decreased
+1	-1	>0	<0	>0	decreased

עבור השורה הראשונה: נשים לב שיש לנו נקודה שהיא מהקלאס החיובי $t=+1$ אבל הפרדיקציה עבורה שלילית $o=-1$, על כן מבחינה

גאומטרית כך יראה המישור (איור). לכן ההפרש בין o ל- t הוא שלילי, ונניח שפיציר ה- xi הוא חיובי, אזי הדוט-פרודקט שלהם שלילי והדלתא טטה חיובית (מפני שאנחנו מכפילים במינוס אטה את הדוט-פרודקט). לכן, המישור "יורד" לכיוון הנקודה, לכיוון השלילי, ונקווה שהעדכון גרם לכך שהנקודה תהיה כעת מעל המישור – המרחק בין xi למישור גדל מפני שהדוט פרודקט בינו לבין טטה גדל.

עבור השורה השנייה: קורה דבר דומה, שוב יש להוריד את המישור לכיוון נקודה ולכן המרחק, המכפלה הפנימית בין xi לטטה, גדל.

לכן באופן כללי אם המכפלה הפנימית בין xi לטטה גדלה, המישור "יורד" ביחס לנקודה שאנחנו מדברים עליה. ואם המכפלה קטנה, המישור "עולה" ביחס לנקודה שאנחנו מדברים עליה.



Perceptron Algorithm

The algorithm:

- Initialize weights to some small random number
- Repeat until convergence (no error = no weight update):
 - For each $x^{(d)}$ in D compute: (* $x^{(d)} = \bar{x}_d$):
 - $o^{(d)} = \text{sgn}(\theta \cdot x^{(d)})$
 - For each θ_i do:
 - $\Delta\theta_i = -\eta \sum_{d \in D} (o^{(d)} - t^{(d)}) x_i^{(d)}$ for each i
 - Update $\theta_i = \theta_i + \Delta\theta_i$



Training data:

- $(-1, -1) \rightarrow +1, (-1, +1) \rightarrow +1, (+1, -1) \rightarrow +1, (+1, +1) \rightarrow -1$

Weight init:

- $\theta_0 = 0.1, \theta_1 = -0.2, \theta_2 = 0.15, \eta = 0.05$

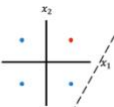


Check $(+1, -1) \rightarrow +1$

- $\text{sgn}(\theta \cdot x^{(d)}) = 0.1 - 0.2 * (+1) + 0.15 * (-1) = -0.25 < 0$
- $o = -1$

Since $t \neq o$ update required:

- $\theta_{0(\text{new})} = 0.1 - 0.05 * (-1 - 1) * 1 = 0.2$
- $\theta_{1(\text{new})} = -0.2 - 0.05 * (-1 - 1) * 1 = -0.1$
- $\theta_{2(\text{new})} = 0.15 - 0.05 * (-1 - 1) * (-1) = 0.05$

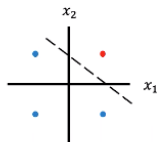


Check $(+1, +1) \rightarrow -1$

- $\text{sgn}(\theta \cdot x^{(d)}) = 0.2 - 0.1 * (+1) + 0.05 * (+1) = 0.15 > 0$
- $o = +1$

Since $t = o$ update required:

- $\theta_{0(\text{new})} = 0.2 - 0.05 * (+1 - (-1)) * 1 = 0.1$
- $\theta_{1(\text{new})} = -0.1 - 0.05 * (+1 - (-1)) * 1 = -0.2$
- $\theta_{2(\text{new})} = 0.05 - 0.05 * (+1 - (-1)) * 1 = -0.05$



We got the linear separator:

$$\theta \cdot \bar{x} = -0.1 * x_1 - 0.15 * x_2 + 0.2 = 0$$

Stochastic Perceptron

The algorithm:

- Set weights randomly
- Repeat until convergence:
 - Choose d randomly (or in some order)
 - Calculate $o^{(d)} = \text{sgn}(\theta \cdot x^{(d)})$
 - Calculate $\Delta\theta_i = -\eta(o^{(d)} - t^{(d)})x_i^{(d)}$ for each i
 - Then update $\theta_i = \theta_i + \Delta\theta_i$

להלן אלגוריתם הפרספטרון

אשר מתייחס לכל השגיאות (ועל

כל הסכום לפני האטה)

אבל נוכל להשתמש בפרספטרון

טוכסטי, שהוא משתמש כל פעם

בשגיאה של נקודה אחת.

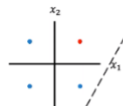
דוגמה עבור פרספטרון סטוכסטי:

שלב ראשון: אתחול הטטה ולהלן המישור ההתחלתי, כאשר "מעל" המישור מוגדר להיות שמאלה כלפי מעלה,

לכן יש שתי טעויות האדומה שמעל המישור והכחולה שמתחת למישור.

שלב שני: בדיקה עבור כל אחת מהנקודות, עבור 2 נקודות לא יתבצע עדכון ויתבצעו שני עדכונים עבור הנקודה

הכחולה שמתחת למישור ועבור הנקודה האדומה שמתחת למישור.

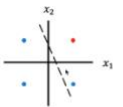


Check $(+1, -1) \rightarrow +1$

- $\text{sgn}(\theta \cdot x^{(d)}) = 0.1 - 0.2 * (+1) + 0.15 * (-1) = -0.25 < 0$
- $o = -1$

Since $t \neq o$ update required:

- $\theta_{0(\text{new})} = 0.1 - 0.05 * (-1 - 1) * 1 = 0.2$
- $\theta_{1(\text{new})} = -0.2 - 0.05 * (-1 - 1) * 1 = -0.1$
- $\theta_{2(\text{new})} = 0.15 - 0.05 * (-1 - 1) * (-1) = 0.05$



Check $(+1, +1) \rightarrow -1$

- $\text{sgn}(\theta \cdot x^{(d)}) = 0.2 - 0.1 * (+1) + 0.05 * (+1) = 0.15 > 0$
- $o = +1$

Since $t = o$ update required:

- $\theta_{0(\text{new})} = 0.2 - 0.05 * (+1 - (-1)) * 1 = 0.1$
- $\theta_{1(\text{new})} = -0.1 - 0.05 * (+1 - (-1)) * 1 = -0.2$
- $\theta_{2(\text{new})} = 0.05 - 0.05 * (+1 - (-1)) * 1 = -0.05$

נבצע בדיקה עבור הנקודה הכחולה, ומכיוון שהפרדיקציה לא

נכונה יתבצע עדכון – לפני העדכון מופיע משמאל ולאחריו

מופיע מימין. נשים לב כי הטטה גדל ולכן המרחק מראשי

הצירים גדל, שזה מצוין עבורנו כי כך הנקודה הכחולה אכן

תתמקם מעל המישור. ונשים לב שטטה 1 וטטה 2 שומרות על

הסימנים שלהן ועל כן ניתן להבין שהזווית לא משתנה מאוד.

נבצע בדיקה עבור הנקודה האדומה, וגם כאן יתבצע עדכון. לפני

העדכון משמאל ולאחריו מימין. נשים לב כי טטה קטנה ולכן

התקרבו לראשית, וכן השתנה הסימן של טטה 2 ולכן נוצר שינוי

בזווית של המישור. כעת הנקודה האדומה אכן מתחת למישור.

אבל נוצרה שגיאה עבור הנקודה הכחולה $(+1, -1)$ שאמורה להיות

מעל המישור אך מופיעה מתחת למישור.

לאחר עדכון נוסף עבור טעות זו נקבל את המישור הסופי שלנו.

איך יודעים מה זה מעל ומה זה מתחת למישור? מציבים את ראשית הצירים וכך אם מתקבל ערך

חיובי, אנו מעל המישור אחרת מתחת. אם נציב את ראשית הצירים כאן נקבל (חיובי) 0.2 ולכן

לכיוון הכחולים נחשב מעל המישור.

מה הבעיה באלגוריתם הפרספטרון ואיך נוכל לפתור אותה?

לאלגוריתם אין תנאי עצירה למעט התכנסות ל-perfect classification, מה שלא תמיד יתאפשר. נפתור אותה על ידי הגבלת כמות הריצות (הבעיה

כאן היא שאנחנו לא נדע אם יכולנו לשפר את מה שהגענו אליו עם "עוד קצת" ריצות). או על ידי הגבלת מספר הטעויות אבל גם כאן נוצרת אותה

הבעיה. ולבסוף הפתרון האלגנטי ביותר היא להפוך את בעיה זו לבעיית אופטימיזציה – להשתמש בגרדיאנט-דיסט ולמצוא מינימום טעות.

אלגוריתם זה נקרא LMS.

LMS – Least Mean Squares

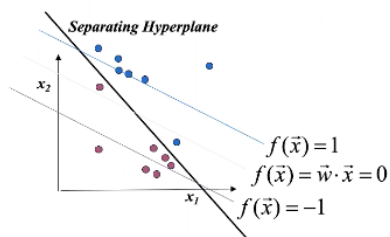
The algorithm:

- Initialize weights to some small random number
- Repeat until convergence (no error = no weight update):
 - For each $x^{(d)}$ in D compute: $\{ * x^{(d)} = \tilde{x}_d \}$:
 - $o^{(d)} = (\theta \cdot x^{(d)})$
 - For each θ_i do:
 - $\Delta \theta_i = -\eta \sum_{d \in D} (o^{(d)} - t^{(d)}) x_i^{(d)}$
 - Update $\theta_i = \theta_i + \Delta \theta_i$

$$E[\tilde{\theta}] = \frac{1}{2} \sum_{d \in D} (o^{(d)} - t^{(d)})^2 = \frac{1}{2} \left[\sum_{d \in D^+} (o^{(d)} - 1)^2 + \sum_{d \in D^-} (o^{(d)} + 1)^2 \right]$$

Minimize the distance between the positive instances and the +1 iso-line of the function

Minimize the distance between the negative instances and the -1 iso-line of the function



What if instead of predicting a continues value we try to predict a class?

Hypothesis function :

$$h_{\theta}(x) = \theta^T x$$

Cost function (MSE):

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

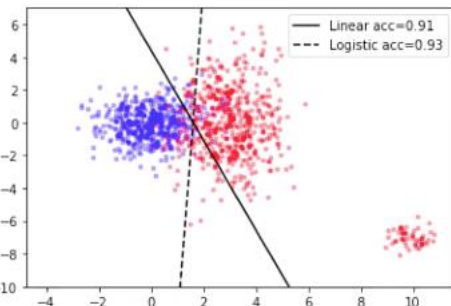
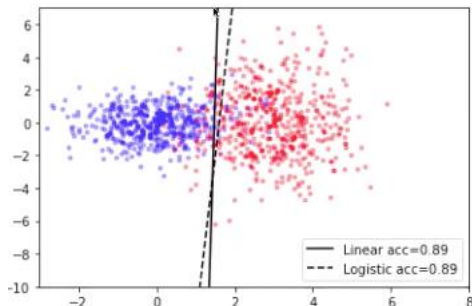
Goal :

$$\min_{\theta} J(\theta)$$

x_1	y
1	1
1	1
0	0
1	1
0	0
1	1

Predict

if $h_{\theta}(x) > 0.5$ then 1
else 0



נשים לב שלסיכום,

- פרספטון לא תמיד טוב מפני שלא תמיד יתכנס
- וקלסיפיקציה על פי רגרסיה לינארית = LMS לא תמיד תניב מפריד אופטימום (כפי שניתן לראות מימין ובדוגמה לעיל).

ולכן נשאלת השאלה איך מוצאים מישור באלגוריתם שמתכנס גם אם הדאטה לא מופרד לינארית? <<

Logistic Regression

$$h_{\theta}(x) = P(1|x)$$

נסה לחזות את ההסתברויות של דגימה להשתייך לקלאס 1. פונקציית ההיפתזה היא הפוסטריור: $h_{\theta}(x) = P(1|x)$. כולומר איך נעבור מרחק ממישור להסתברות? ככל שנקודה נמצאת בצד הנכון ורחוקה יותר מהמישור, ההסתברות שלה גבוהה יותר. כלומר עלינו לקחת למרחקים ולהפוך אותם להסתברויות ביחס למישור. מתמטית נעביר סקאלה של מרחקים לסקאלה של 0-1. אחת הפונקציות שיועזות לעשות זאת נקראית

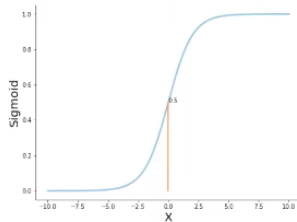
$$S(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}} = \frac{e^{\theta^T x}}{1 + e^{\theta^T x}}$$

סיגמויד sigmoid : ונשתמש בה ב-LoR

כאשר מכפלה פנימית של טטה עם איקס היא בדיוק המרחק של הנקודה מהמישור ולכן אנחנו לוקחים את המרחק הזה מכניסים לסיגמויד ומקבלים תוצאה בין 0-1.

x_1	y
1	1
1	1
0	0
1	1
0	0
1	1

$$S(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x}$$



Score	$-\infty$	-2	0	+2	$+\infty$
Sigmoid (Score)	$\frac{1}{1 + e^{-\infty}}$ = 0		$\frac{1}{1 + e^0}$ = 0.5		$\frac{1}{1 + e^{-\infty}}$ = 1

לכן, נציב את הסיגמויד (מסומן ב-S) בפונקציית ההיפותזה שלנו: $h_\theta(x) = S(\theta^T x)$.
והפרדיציה שלנו תהיה 1 אם $h_0(x) > 0.5$ ואחרת, 0.

נשתמש ב-maximum likelihood כדי להגדיר את פונקציית העלות. נרצה לחזות את הלייבל y בהינתן הדאטה והמישור, ולקבל עבורו את ההסתברות הגבוהה ביותר. ובאופן מתמטי נרצה למקסם את ההסתברות

$$P(y|x, \theta) \triangleq (h_\theta(x))^y \cdot (1 - h_\theta(x))^{1-y}$$

הבאה (הלייקליהוד):

נזכור כי $h_0(x)$ היא ההסתברות של האינסטנס x להשתייך לקלאס 1.

$$\begin{aligned} P(0|x, \theta) &= 1 - h_\theta(x) \\ P(1|x, \theta) &= h_\theta(x) \end{aligned}$$

ולכן נקבל:

$$P(y|x, \theta) = (h_\theta(x))^y \cdot (1 - h_\theta(x))^{1-y}$$

Assuming independent instances

$$P(D|\theta) = \prod_{d=1}^m P(y^{(d)} | x^{(d)}, \theta) =$$

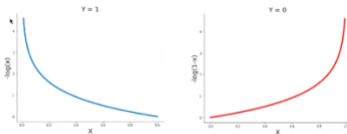
$$\prod_{d=1}^m (h_\theta(x^{(d)}))^{y^{(d)}} \cdot (1 - h_\theta(x^{(d)}))^{1-y^{(d)}}$$

פיתוח מתמטי: (לאחר שנניח אי תלות על האינסטנסים נרצה למצוא את המקסימום של מכפלת ההסתברויות המוקפת באדום, לכן נפעיל לן ונגזור ונקבל את הגרדיאנט – פיתוח מלא ומפורט יותר מופיע בהרצאה)

$$\underset{\theta}{\operatorname{argmin}} \sum_{d=1}^m -y^{(d)} \cdot \ln(h_\theta(x^{(d)})) - (1 - y^{(d)}) \cdot \ln(1 - h_\theta(x^{(d)}))$$

Cost Function Intuition

$$\text{Cost}(x, \theta) = \begin{cases} -\log(h_\theta(x)) & y = 1 \\ -\log(1 - h_\theta(x)) & y = 0 \end{cases}$$



$$-\sum_{d=1}^m y^{(d)} \theta^T x^{(d)} - \ln(1 + e^{\theta^T x^{(d)}})$$

לאחר הצבת הסיגמויד זו הפונקציה שנרצה למזער:

נמזער את הפונקציה הזו על ידי גרדיאנט דיסנט כאשר זו פונקציית העלות:

$$\text{cost}(\vec{\theta}) = -\sum_{d=1}^m y^{(d)} \theta^T x^{(d)} - \ln(1 + e^{\theta^T x^{(d)}})$$

$$\frac{\partial}{\partial \theta_i} \text{cost}(\vec{\theta}) = -(y - S(\vec{\theta}, \vec{x})) x_i$$

הנגזרת עבור כל אינסטנס כאשר i הוא הפיצ'ר ה-i תניב:

$$\frac{\partial}{\partial \theta_i} \text{cost}(\vec{\theta}) = \sum_{d=1}^m (S(\vec{\theta}, \vec{x}^{(d)}) - y^{(d)}) x_i^{(d)}$$

ולכל m נקודות הדאטה מתקיים:
וכעת נפעיל גרדיאנט דיסנט.

ולסיכום, להלן ההיפותזה, העלות ונרצה למצוא טטה אשר ממזערת את העלות.

(כאן אין פתרון על ידי pinv)

• Hypothesis function :

$$h_\theta(x) = S(\theta^T x)$$

• Cost function:

$$J(\theta) = \frac{1}{m} \sum_{d=1}^m -y^{(d)} \cdot \ln(h_\theta(x^{(d)})) - (1 - y^{(d)}) \cdot \ln(1 - h_\theta(x^{(d)}))$$

• Goal :

$$\min_{\theta} J(\theta)$$

איך נוכל לפתור בעיה שיש בה יותר משני קלאסים עם מפרד ליניארי? גישת one vs all.

נניח כי יש לנו 3 קלאסים, נמצא לכל אחד מהם מפרד ליניארי עבור כל אחד מהקלאסים, סך הכל 3 מפרדים ליניארים. כאשר תגיע נקודה חדשה נסווג אל הקלאס שביחס לאחרים ההסתברות של נקודה חדשה להשתייך אליו תסווג אליו – ההיפותזה של כל קלאס אומרת מה ההסתברות שאנחנו משתייכים לקלאס זה. עבור הנקודה העליונה ברור שההיפותזה של קלאס האיקסים תהיה גבוהה מההיפותזה של הריבועים והעיגולים ולכן היא תסווג כאיקס.

