

Logistic Regression (LOR)

זהו אלגוריתם קלסיפיקציה, בניגוד לרגרסיה ליניארית אשר מניב ערכים רציפים לכל דגימה, אלגוריתם זה מניב סיווג בדיד. כמו בכל מסווג,

הדאטה אימון במקרה הבינארי, מכילה וקטורי פיצורים ואת הלייבלים המתאימים (labels): $\{(\vec{x}^{(i)}, y^{(i)})\}_{i=1..m}$

אנחנו נלמד מודל שהוא יהיה פונקציה h המקבלת וקטור פיצורים x (בממד k למשל) כך ש- h מקיימת: $h: \mathbb{R}^k \rightarrow \mathbb{R}$.
מודל logistic regression חשוב לשיטות למידה מודרניות יותר וכן עבור deep learning.

אנחנו נרצה שהפונקציה שנמצא, h , תחזיר עבורנו הסתברויות כך שלמעשה היא תחזיר ערכים מהקטע $[0,1]$ ולא דווקא מ- \mathbb{R} .

ב-logistic regression אנחנו מניחים את ההנחה המרכזית הבאה:

נניח כי ההסתברות $P(Y=1|X=\vec{x})$ (ההסתברות להיות משוייך ללייבל מסוים Y בהינתן דגימה X) יכולה להיות משוערכת כפונקציית sigmoid

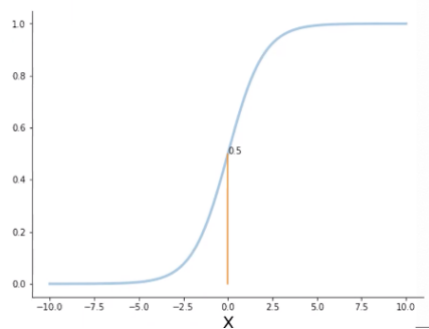
המיושמת על קומבינציה ליניארית על ה-input features. באופן מתמטי, עבור נקודת דאטה (\vec{x}, y) , אנו מניחים ב-LoR כי:

$\sigma(z) = \frac{1}{1+e^{-z}}$ כך ש- $P(Y=1|X=\vec{x}) = \sigma(z)$ (זוהי פונקציית ה-sigmoid), ו- $z = \theta_0 + \sum_{j=1}^k \theta_j x_j$ (קומבינציה ליניארית של הוקטור x המקורי

לכן z הוא למעשה סקלר, מספר ממשי, מפני שיש כאן מכפלה פנימית של הוקטור טטה עם הוקטור x).

ההצבה היא להלן, כאשר השוויון השני מתקיים מפני שאנחנו מכפילים מונה ומכנה ב $e^{\langle \theta^T, x \rangle}$

- הפונקציה הזו מונוטונית עולה ממש
- הטווח של הפונקצייה הוא $[0,1]$, הפונקציה ראשית כל חיובית, המכנה גדול מהמונה ולכן קטנה מ-1.
- הגבולות באינסוף: 0 במינוס אינסוף ו-1 באינסוף.



Derivative:
 $\sigma'(x) = \sigma(x)(1 - \sigma(x))$

תכונה נוספת של פונקציית סיגמויד שיכולה לסייע הינה:

מכאן, שנוכל להפוך את פונקציית ההיפוטזה שלנו / המודל שלנו לצורה הבאה:

$$h_{\theta}(\vec{x}) = P(y=1|\vec{x}) = \sigma(\langle \theta, (1, \vec{x}) \rangle)$$

מכפלה פנימית של טטה עם הרחבה של הוקטור איקס בתוך הפונקציה של סיגמויד. כאשר השוויון השני נובע מהנחת ה-logistic regression.

if $h_{\theta}(\vec{x}) > 0.5$ then 1
else 0
כאשר הפרדיקציה / הקלסיפיקציה תערך באופן הבא:

עלינו ללמוד את הטטות שמניבות את הניבוי הטוב ביותר – נבצע תהליך למידה עבורן. טטה הוא וקטור פרמטרים באורך $k+1$ (למען טטה 0), אנחנו נלמד את הערכים שלה תחת מסגרת maximum likelihood. לשם כך נשתמש ב-training data set שנסמן אותה ב- D שיש בה m דגימות. כל דגימה תכיל וקטור של פיצורים x מממד k וערך בינארי y שנקרא label.

Maximum Likelihood

ראשית, נשים לב כי תחת ההנחה המרכזית שאנו מניחים ב-LoR, מתקיים עבור נקודת דאטה יחידה:

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

אבל זה מתקיים עבור נקודת דאטה אחת, לכן כעת עבור הדאטה אימון D עם m דגימות שנניח כי הן בלתי תלויות, נרצה למצוא את המודל

שממקס את ה-likelihood. עבור כל טטה, הלייקליהוד שלה בהינתן הדאטה, $P(D|\theta)$, היא ביחס ל-

נרצה למצוא את הטטה שמביאה זאת למקסימום.

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

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

לכן נגזור ונשווה ל-0. השורה הראשונה מועתקת מהעמודה הקודם. בשורה השניה אנחנו רוצים לבטל מכפלות כדי שיהיה קל יותר לכפול ולכן אנו מפעילים לן (מונוטוניות עולה ולכן איננה משנה את המקסימום). בשורה השלישית המכפלה הפכה לסכום והופכל לן נוסף על שתי החזקות כדי שנוכל להוריד את $y^{(d)}$ מהחזקה – שורה רביעית. ובשורה החמישית והאחרונה אנחנו הופכים סימן למינוס כדי למצוא את ה-minimum. כי אז ניתן לחשוב על זאת במונחי cost.

$$\begin{aligned} \operatorname{argmax}_{\theta} \prod_{d=1}^m \left(h_{\theta}(x^{(d)}) \right)^{y^{(d)}} \cdot \left(1 - h_{\theta}(x^{(d)}) \right)^{1-y^{(d)}} &= \\ \operatorname{argmax}_{\theta} \ln \left(\prod_{d=1}^m \left(h_{\theta}(x^{(d)}) \right)^{y^{(d)}} \cdot \left(1 - h_{\theta}(x^{(d)}) \right)^{1-y^{(d)}} \right) &= \\ \operatorname{argmax}_{\theta} \sum_{d=1}^m \ln \left(\left(h_{\theta}(x^{(d)}) \right)^{y^{(d)}} + \ln \left(\left(1 - h_{\theta}(x^{(d)}) \right)^{1-y^{(d)}} \right) \right) &= \\ \operatorname{argmax}_{\theta} \sum_{d=1}^m y^{(d)} \cdot \ln \left(h_{\theta}(x^{(d)}) \right) + (1 - y^{(d)}) \cdot \ln \left(1 - h_{\theta}(x^{(d)}) \right) &= \\ \operatorname{argmin}_{\theta} \sum_{d=1}^m -y^{(d)} \cdot \ln \left(h_{\theta}(x^{(d)}) \right) - (1 - y^{(d)}) \cdot \ln \left(1 - h_{\theta}(x^{(d)}) \right) & \end{aligned}$$

לכן נרצה להביא למינימום את: (למדא של טטה)

$$\Lambda(\vec{\theta}) = - \sum_{d=1}^m y^{(d)} \ln \left(\sigma(\theta, \vec{x}^{(d)}) \right) + (1 - y^{(d)}) \ln \left(1 - \sigma(\vec{\theta}, \vec{x}^{(d)}) \right)$$

נציין כי כאן לא עובד לחשב גרדיאנט ולהשוות ל-0 (חישוב אספליסיטי לא יעבוד כאן). לכן נשתמש בגרדיאנט דיסנט.

נוכל להראות כי: זהו הגרדיאנט (וקטור בגודל $k+1$)

$$\frac{\partial}{\partial \theta_j} \Lambda(\vec{\theta}) = \sum_{d=1}^m (\sigma(\vec{\theta}, \vec{x}^{(d)}) - y^{(d)}) x_j^{(d)}$$

הנגזרת הכיוונית לפי טטהן המחושבת בנקודה טטה שהיא וקטור $k+1$ ממדי.

מכאן, נוכל להתחיל להריץ גרדיאנט דיסנט מפני שכל מה שכתוב בגרדיאנט נמצא ברשותנו (בהנחה שיש לי ניחוש התחלתי לוקטור טטה).

אלגוריתם גרדיאנט דיסנט עבור Logistic Regression

- נתחיל את טטה להיות ערך התחלתי קטן / ניחוש
- נחזור על התהליך הבא עד אשר נגיע להתכנסות:

$$\theta_j^{New} = \theta_j^{Old} - \alpha \frac{\partial \Lambda}{\partial \theta_j}(\theta^{Old}, X, y)$$

נבצע עידכון של הווקטור טטה באופן הבא:

$$\theta_j^{New} = \theta_j^{Old} - \alpha \cdot \sum_{d=1}^m (\sigma(\vec{\theta}, \vec{x}^{(d)}) - y^{(d)}) x_j^{(d)}$$

או באופן יותר אקספליסיטי, אם נציב את הגרדיאנט שהראינו לעיל:

- כאשר אלפא הוא ה-learning rate

חישוב הגרדיאנט של למדא (שכבר ראינו לעיל) – הוספת סיגמויד (שמאל) ולאחר מכן גזירה לפי טטהן (ימין):

$$\begin{aligned} \Lambda(\vec{\theta}) &= - \sum_{d=1}^m y^{(d)} \ln \left(\sigma(\vec{\theta}, \vec{x}^{(d)}) \right) + (1 - y^{(d)}) \ln \left(1 - \sigma(\vec{\theta}, \vec{x}^{(d)}) \right) \\ \ln(\sigma(\vec{\theta}, \vec{x}^{(d)})) &= \ln \left(\frac{1}{1 + e^{-\theta^T x^{(d)}}} \right) = - \ln \left(1 + e^{-\theta^T x^{(d)}} \right) \\ \ln(1 - \sigma(\vec{\theta}, \vec{x}^{(d)})) &= \ln \left(1 - \frac{1}{1 + e^{-\theta^T x^{(d)}}} \right) \\ &= \ln \left(\frac{1 + e^{-\theta^T x^{(d)}}}{1 + e^{-\theta^T x^{(d)}}} - \frac{1}{1 + e^{-\theta^T x^{(d)}}} \right) \\ &= \ln \left(\frac{e^{-\theta^T x^{(d)}}}{1 + e^{-\theta^T x^{(d)}}} \right) \\ &= -\theta^T x^{(d)} - \ln \left(1 + e^{-\theta^T x^{(d)}} \right) \end{aligned}$$

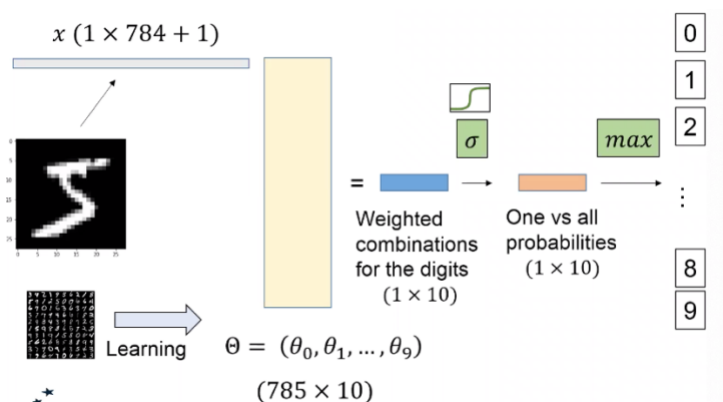
$$\begin{aligned} \frac{\partial}{\partial \theta_j} \left(- \sum_{d=1}^m y^{(d)} \theta^T x^{(d)} + \ln \left(1 + e^{-\theta^T x^{(d)}} \right) \right) \\ = - \sum_{d=1}^m y^{(d)} x_j^{(d)} + \sum_{d=1}^m \left(\frac{1}{1 + e^{-\theta^T x^{(d)}}} \right) x_j^{(d)} \end{aligned}$$

$$\frac{\partial}{\partial \theta_j} \Lambda(\vec{\theta}) = \sum_{d=1}^m (\sigma(\vec{\theta}, \vec{x}^{(d)}) - y^{(d)}) x_j^{(d)}$$

חזרה על עקרונות ה-Logistic Regressions שלדמנו עד כה:

- פונקציית ההיפותזה או המודל, על פי הנחת LoR הינה: $h_{\theta}(\vec{x}) = P(y = 1 | \vec{x}) = \sigma(\theta^T \vec{x})$
- פונקציית המחיר / השגיאה, שנובעת מחישוב ההסתברות $P(D | \theta)$ הינה: $\Lambda(\theta) = \frac{1}{m} \sum_{d=1}^m -y^{(d)} \cdot \ln \left(h_{\theta}(x^{(d)}) \right) - (1 - y^{(d)}) \cdot \ln \left(1 - h_{\theta}(x^{(d)}) \right)$
- נשתמש בגרדיאנט דיסנט במטרה למצוא את: $\operatorname{argmin}_{\theta} \Lambda(\theta)$

דוגמה: דאטה שנקרא MNIST Digits ושמשמעותו ספרות שכותבים בני אדם. כל ספרה מוצגת בגודל 28×28 פיקסלים ונרצה למצוא מודל שזוהה את הספרה הכתובה. אנחנו ניקח כל פיקסל ונשטח אותו לוקטור אחד, **אינפוט וקטור** x , $(28 \times 28 = 784)$, שהוא שורה אחת ו-784 תאים כלומר 785. לכל ספרה למדנו מודל = טטה, ששואל את השאלה, האם זה 0 או לא? האם זה 1 או לא? ... כלומר כל אחת מהטטות שואלת שאלה בינארית האם הספרה שמופיעה בתמונה היא i לכל i בין 0 ל-9 או שלא. לכן קיבלנו 10 טטות. התוצאה היא **מטריצה (הצהובה)** שגודלה 785 שורות ו-10 עמודות, 785 שורות מפני שכל טטה מקבלת אינפוט שהוא 785. נבצע מכפלה פנימית בין הוקטור x של התמונה ה"משוטחת" לבין המטריצה הזו



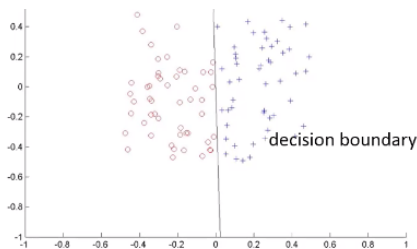
ונקבל **וקטור בעל 10 ערכים (1×10)** ונכניס אותו לתוך פונקציית סיגמויד (המעבר בין הכחול לכתום) ונקבל וקטור כתום שהוא גם כן בגודל 1×10 אבל הערכים בו הם בין 0 ל-1, אלה הן ההסתברויות: כמה המודל חישב שהסיכויים שמדובר במספר 0, הסיכויים שמדובר ב-1, הסיכויים שמדובר ב-2 וכו... ומתוך הוקטור הזה ניקח את המיקום בו ההסתברות היא מקסימלית ולפיכך נסווג את התמונה לקלאס מ-0 עד 9. הוקטור הכתום הוא למעשה כבר ה-posterior, האלגוריתם LoR מגיע ישירות ל-posterior. גישה זו נקראת one vs all. ומספר הטטות שנלמד הוא כמספר הקלאסים (אם עלינו ללמוד אותיות אנגליות קטנות היינו לומדים 26 טטות)

לסיכום,

- LoR היא גישה לקלסיפיקציה בעלת שם מטעה (כי בגרסיה לינארית לא ביצענו סיווג).
- האלגוריתם המבצע (או המודל) מסווג בהתבסס על וקטור בעל פיצורים נומריים.
- הגישה מבוססת על הנחת ה-LoR שהינה $P(y = 1|\vec{x}) = \sigma(\theta^T \vec{x})$
- תהליך הלמידה משתמש בגרדיאנט דיסנט כדי למצוא טטה המניבה מקסימום לייקליהוד / מינימום טעות בהינתן ה-observed data, תחת הנחת ה-LoR.
- אין אפשרות להשתמש ב-pseudo inverse ב-LoR.
- ניתן להשתמש באלגוריתם זה עבור קלאסים רבים (כפי שראינו בדוגמה לעיל one vs all).

Linear Classifiers – קלסיפיקציה לינארית

הערה: מסווג בייסיאני יכול לסווג גם משתנים קטגוריאליים וגם רציפים, LoR לעומת זאת מסווגת לפי פיצורים נומריים בלבד. **מפרידים לינאריים** = נניח כי מרחב הדגימות שלנו הוא R^2 (כלומר 2 פיצורים). כל דגימה היא נקודה. אם יש קו המפריד את שני הקלאסים (או קו המפריד בין קונספטים) אזי שני הקלאסים הם **ניתנים להפרדה לינארית** (או: הקונספט ניתן להפרדה לינארית).



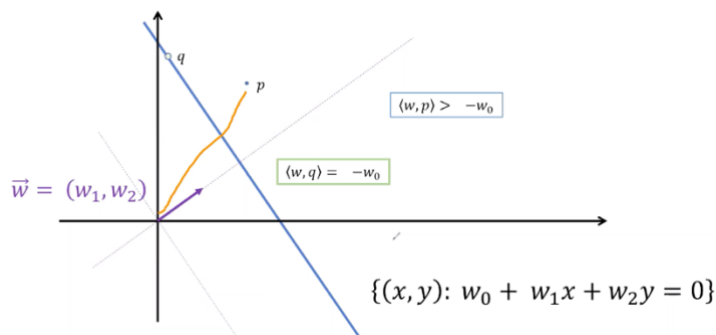
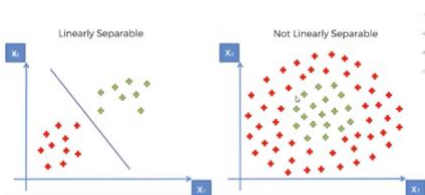
2D Dichotomy Example

ב- R^2 נייצג קו על ידי משוואת קו ישר: $f(x,y) = w_1x + w_2y + w_0 (=0)$

(דיכוטומיה = דאטה עם עיגולים ופלוסים כפי שרואים באיור)

דיכוטומיה ב- R^2 ניתנת להפרדה ע"י קו אם: $f(x,y) > 0 \Rightarrow +1$, $f(x,y) < 0 \Rightarrow -1$

המשמעות הגיאומטרית עבור המשוואה שבה בחרנו לייצג קו ישר:



Linear Discriminant Functions

$$g(x) = w^T x + w_0 = 0$$

בהכללה כללית יותר, ב- R^n , כאשר $g(x)$ היא לינארית היא מגדירה היפר-מישור ע"י ההשוואתה ל-0: w יקרא וקטור המשקולות ו- w_0 הינו ה-bias. נשים לב כי נוכל לחשוב על כך כמכפלה פנימית של (w_0, w) עם $(1, x)$. והסיווג יתבצע באופן הבא:

$$\begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise} \end{cases}$$

בהינתן שקיים דאטה כזה שהוא ניתן להפרדה לינארית, נרצה ללמוד את w שיקנה עבורנו מסווג. לכן עלינו למזער איזשהי פונקציית טעות ולמצוא את ה- w ים אשר ממזערים אותה.

$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} (\vec{w} \cdot \vec{x}_d - t_d)^2$$

נסה ללמוד את ה- w ים אשר מביאים למינימום את השגיאה הריבועית על כל דגימות הדאטה (D הוא הדאטה סט): פונקציית הטעות מקבלת וקטור w ומה שהיא עושה היא לסכום את: המכפלה הפנימית של הוקטור w עם וקטור הפיצורים x_d פחות הלייבל t_d של הנקודת דאטה הזו (זה למעשה y_d וסתם סומן באופן שונה) בריבוע. זו פונקצייה המזכירה סכום ריבועים פחותים. זהו למעשה בדיוק MSE לרגרסיה לינארית! לקחנו את שאלת הקלסיפיקציה שלי והפכתי אותה לשאלת רגרסיה לינארית. ולכן גם נשתמש כאן בגרדיאנט דיסנט. ומכאן נקבל את האלגוריתם - LMS Classification Algorithm.

LMS Classification Algorithm

סימונים: אטה (נראית כמו η) תסמן את ה-learning rate שלנו (למשל 0.05), ו- D יסמן את הסט של דגימות ה-training. כל דגימת אימון הינה זוג מהצורה $\langle x, t \rangle$ כאשר x הוא וקטור הפיצורים ו- t הוא ה-label או ה-target output value.

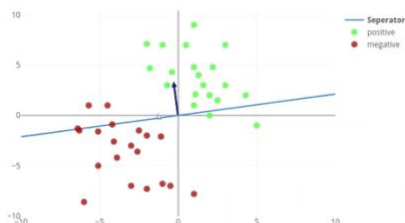
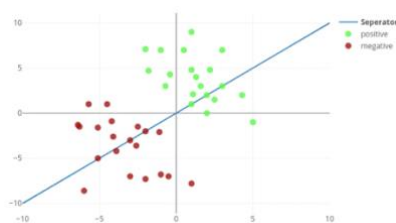
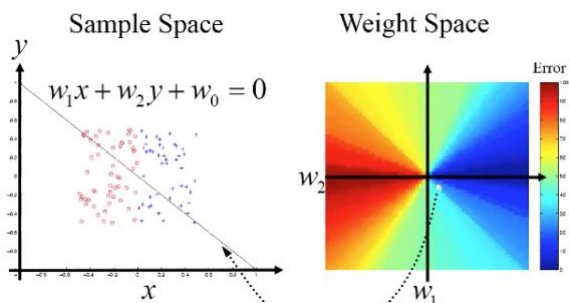
- נאתחל ערכים עבור וקטור המשקולות w
- נחזור עד התכנסות על הצעד הבא:

$$o_d = \vec{w} \cdot \vec{x}_d \quad \text{-- לכל וקטור פיצורים } x_d \text{ בדאטה-סט } D \text{ נחשב:}$$

$$\Delta w_i = -\eta \sum_{d \in D} (o_d - t_d) x_{id}$$

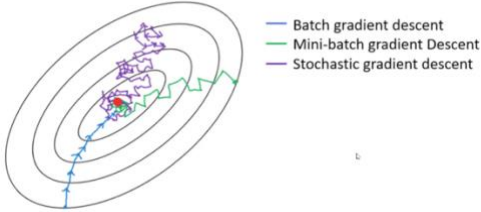
$$w_i = w_i + \Delta w_i$$

-- ועבור כל יחידת משקל לינארית w_i , נבצע:



דוגמה לתהליך הלמידה: משמאל הדאטה שיהיה עלינו ללמוד (הירוקים מול האדומים) וראינו כיצד תהליך הלמידה משנה את כיוון המפריד, כך שהחץ המאונך בדוגמה מימין (שלב מסוים מתהליך הלמידה) מייצג את וקטור המשקלים הלינאריים w .

Stochastic Vs. Standard Gradient Descent



Batch

Initialize each w_i to some small random number.

Until termination do

For each \vec{x}_d in D compute

$$o_d = \vec{w} \cdot \vec{x}_d$$

For each linear unit weight w_i , Do

$$\Delta w_i = -\eta \sum_{d \in D} (o_d - t_d) x_{di}$$

$$w_i = w_i + \Delta w_i$$

Stochastic

Initialize each w_i to some small random number.

Until termination do

For each \vec{x}_d in D compute

$$o_d = (\vec{w} \cdot \vec{x}_d)$$

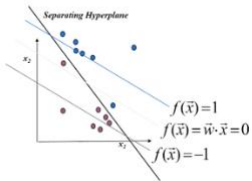
For each linear unit weight w_i , Do

$$\Delta w_i = -\eta (o_d - t_d) x_{di}$$

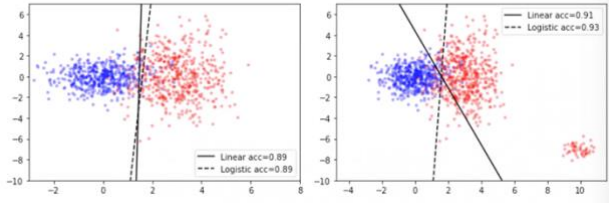
$$w_i = w_i + \Delta w_i$$

1st Yakhini IDC

$$E[\vec{w}] = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 = \frac{1}{2} \left[\sum_{d \in D} (\vec{w} \cdot \vec{x}_d - 1)^2 + \sum_{d \in D} (\vec{w} \cdot \vec{x}_d + 1)^2 \right]$$



Given the hypersurface (derived from attributes of objects *linear* and *logistic*), we found the line $h(x, y) = 0.5$ (the threshold that separates the data, in the LoR case)



לכן, נרצה לשפר את LMS ולכפות עליו למזער טעויות: נרצה למצוא פונקציית טעות שאכן סופרת טעויות. הפונקציה מימין: m הוא הגודל של D , ונרצה לספור את כמות הטעויות. הסכום על פונקציית הסימן אשר מבצעת מכפלה בין הלייבל ובין הפרדיקציה שלנו (כאשר אנו מסווגים בין 1 ל-1), תסכום ל- m אם אין טעויות מפני שמינוס מינוס יניב פלוס, ופלוס פלוס יניב פלוס, וכל מיס קלסיפיקציה תיספר כמינוס. לכן ברגע שנבצע m פחות הסכום הנ"ל, ככל שהסכום יותר קרוב ל- m , יש פחות טעויות, וככל שנתרחק מ- m , האלגוריתם שלנו יבצע יותר מיס-קלסיפיקציות. השאלה היא איך ניתן לעבוד ולמזער את הפונקציה הזו (שמחזירה ערכים בין 0 ל-1) ל- m טעויות) שכן לגזור את פונקציית הסימן לא תניב תוצאה יפה. מה שמביא אותנו ל-**Perceptron**.

$$E[\vec{w}] = \frac{1}{2} \sum_{d \in D} (o_d - t_d)^2$$

בגרדיאנט דיסנט סטנדרטי טעויות נסכמות על כל הדגימות לפני העידכונים: לפעמים ה- training set הוא מאוד גדול ואנחנו לא נרצה לחכות לעדכון עד שנרץ לולאה על כל הדגימות. בגרסה הסטוכסטית של האלגוריתם, משקלים מתעדכנים על פי **דגימה אחת** רנדומלית או על פי **קבוצה קטנה של דגימות mini-batch**. המשמעות היא שדרושים פחות חישובים לפני כל עידכון, אבל גם דרושים צעדים קטנים יותר מכיוון שלא עושים שימוש בגרדיאנט האמיתי. הגישות הסטוכסטיות יכולות להתמודד יותר טוב עם הבריחה ממנימום לוקאלי וכן מאפשרות מיקבול. מה שראינו עד עכשיו היא גישת ה-Batch (חישוב הטעות על כל הדאטה).

$$\mathbf{X} \cdot \vec{w} = \vec{t}$$

$$\begin{pmatrix} x_{10} & x_{11} & \dots & x_{1n} \\ x_{20} & x_{21} & \dots & x_{2n} \\ \vdots & \vdots & & \vdots \\ x_{m0} & x_{m1} & \dots & x_{mn} \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{pmatrix} = \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_m \end{pmatrix}$$

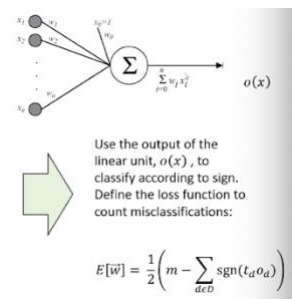
הצבה נוספת: גישת ה-pseudo inverse (או pinve) גם כן עובדת עבור LMS

מחן הבעיות באלגוריתם LMS: נוסה למנות אותן באמצעות הדוגמה הבאה:

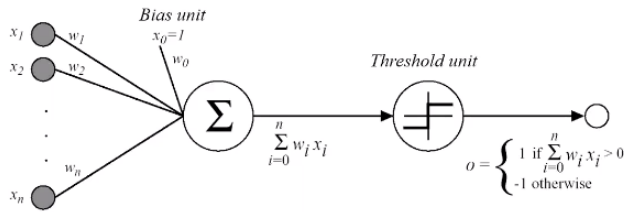
נפעיל את האלגוריתם על הדאטה הכחול-סגול שלנו ונקבל את הקו האפור הבהיר. ולמעשה הוא איננו מפריד מושלם! יש טעויות, ישנה נקודה סגולה בין כחולים ונקודה כחולה בין סגולים! אבל הדאטה אכן ניתן להפרדה לינארית – להלן הקו השחור! אבל למה ההפרדה איננה מושלמת? זאת מכיוון שהגדרנו פונקציית טעות שאיננה סופרת מיס-קלסיפיקציות! אלא רק מגיעה לסוג של איזון, רגרסיה, מרחקים "טובים" מספיק בין הנקודות להפרדה.

נשווה בין LMS לבין LoR:

נפעיל את שני האלגוריתמים על שני דאטה-סטס, כך שהקלסיפיקציה עבור שני קבוצות הדאטה היא כחול-אדום, ונחשב את הדיוק של כל אחד מהאלגוריתמים (בכמה אחוזים צדקנו). בדאטה השמאלי שנוצר משני גאוסיאנים אחוז הדיוק זהה 89% עבור שני האלגוריתמים. לעומת זאת בדאטה הימני אשר נוצר מגאוסיאן כחול אחד ושני גאוסיאנים אדומים, אחוז הדיוק של LoR גובר על אחוז הדיוק של המסווג הלינארי LMS. הפונקציה הלינארית רואה את הענף האדום הקטן כ-1 או -1, לכן היא לא "מסתובבת" יותר מידי ולכן מתווספות אליה טעויות.



The Perceptron



$$o(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise} \end{cases}$$

Or in vector notations:

$$o(\vec{x}) = \begin{cases} 1 & \text{if } \vec{w} \cdot \vec{x} > 0 \\ -1 & \text{otherwise} \end{cases}$$

נרצה לקחת את האינפוטס x_1, \dots, x_n ולבצע מכפלה פנימית עם וקטור המשקלים w , לבדוק אם הוא מניב תוצאה חיובית או שלילית ולסווג בהתאמה לקלאס 1, -1. עלינו ללמוד w -ים שיביאו למינימום את הפונקציה הזו. להלן האלגוריתם הלומד:

The Perceptron Algorithm

אנחנו נניח כי הלייבל, ה-target value, עבור הקלסיפיקציה יהיו 1 או -1. ואלגוריתם זה יהיו דומה מאוד ל-LMS – אך מטרתו היא למצוא את המפריד הליניארי ללא טעויות! וכן פרספטרון הוא תמיד סטוכסטי, אין פרספטרון אשר משתמש בדגימה אחת או קבוצה קטנה של דגימות (אין batch). האם הדבר תמיד מתאפשר?

- נאתחל את ערכי וקטור המשקלים w
- נבצע עד התכנסות:

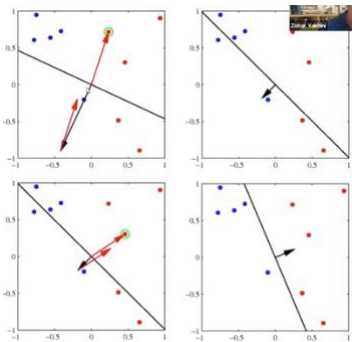
$$o_d = \text{sgn}(\vec{w} \cdot \vec{x}_d) \quad \text{-- לכל וקטור פיזורים } x_d \text{ ב-D נחשב:}$$

For any misclassified (at the present iteration) training instance, x , with $C(x) = +1$ we update the weights as:
 $w = w + \eta x$

$$\Delta w_i = -\eta(o_d - I_d)x_{id}$$

-- ולכל יחידת משקל לינארית w_i נבצע: $w_i = w_i + \Delta w_i$

תעבה: כאשר בחישוב הדלטא w נשים לב כי o_d שווה ל-0 רק כאשר אין טעות! ולכן אין לנו למה לבצע שינוי / לתקן. למעשה העידכון הזה 0 כאשר אנחנו צודקים, כאשר אנחנו טועים ולמשל התקבל 2, אזי w כולו כאשר אנחנו טועים אנחנו מוסיפים "חתיכה" של x_d . אז איך בעצם פונקציית הטעות משתפרת בכל איטרציה?



Red dots are positive ($t = +1$)
 The marked one is initially mis-classified
 Recall that we then update by:
 $w^{(j+1)} = w^{(j)} + \eta x$
 We therefore get, for a misclassified x with $C(x) = +1$:

$$w^{(j+1)} \cdot x > w^{(j)} \cdot x$$

נדגים כיצד פרספטרון משפר את פונקציית הטעות:

נשים לב ששילבי השיפור הם: השלב הראשון הוא הגרף השמאלי העליון, השלב השני הוא הגרף הימני העליון, השלב השלישי הוא הגרף השמאלי התחתון והשלב הרביעי בו הגענו למסוג מושלם הוא הימני התחתון. כאשר הוקטור w והשינוי שלו מיוצג על ידי החץ השחור. והחצים האדומים מייצגים את הוקטורים x_d אשר מצביעים לנקודת דאטה, ו"חלק" מהם על פי הטעות מתווספים ל- w .

- Note: we also need to control η to really guarantee convergence (if its too big we may overshoot the perfect classifier)
- Some results on the rate of convergence were proven and can be useful in the context of ANNs
- The Perceptron itself is not a practical learning approach but is an important component of many modern learning approaches.

משפט רוזנבלאט לאלגוריתם הפרספטרון: אלגוריתם הלמידה פרספטרון מתכנס למסוג מושלם (ללא טעויות על ה-training data) אם ורק אם ה-training data, D, ניתנת להפרדה לינארית.