

הרצאה 9 - גישה

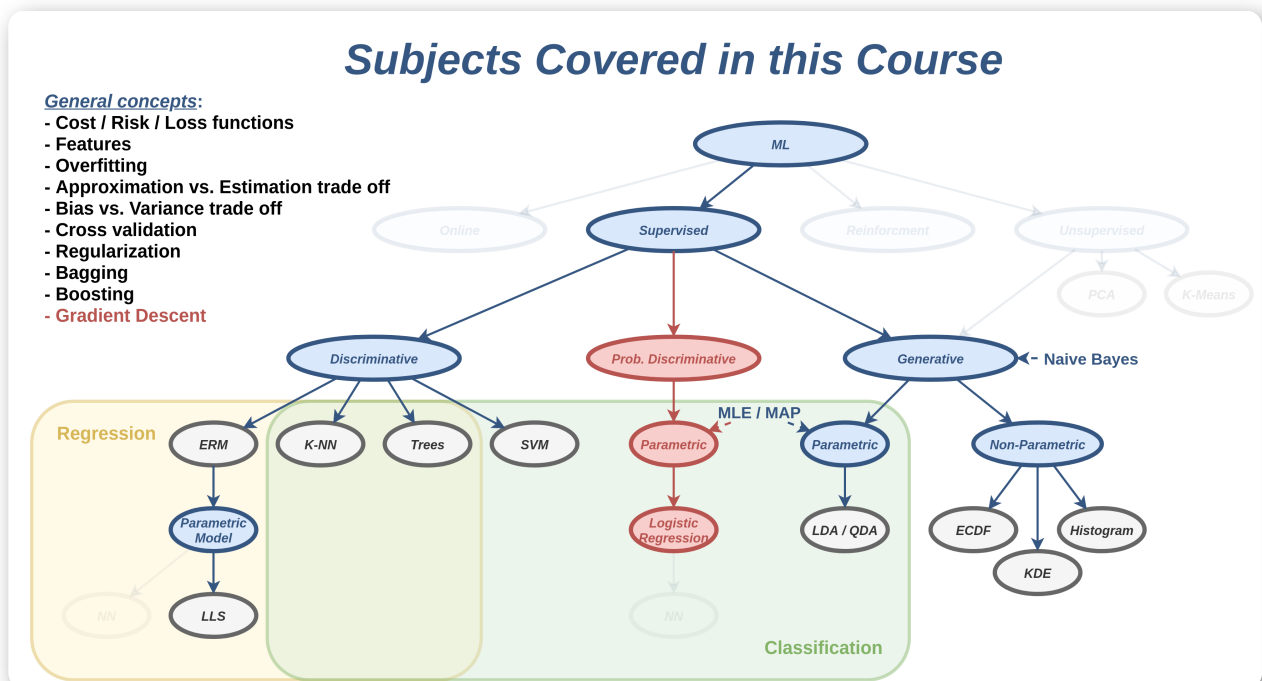
דיסקרימינטיבית הסתברותית

Slides

PDF

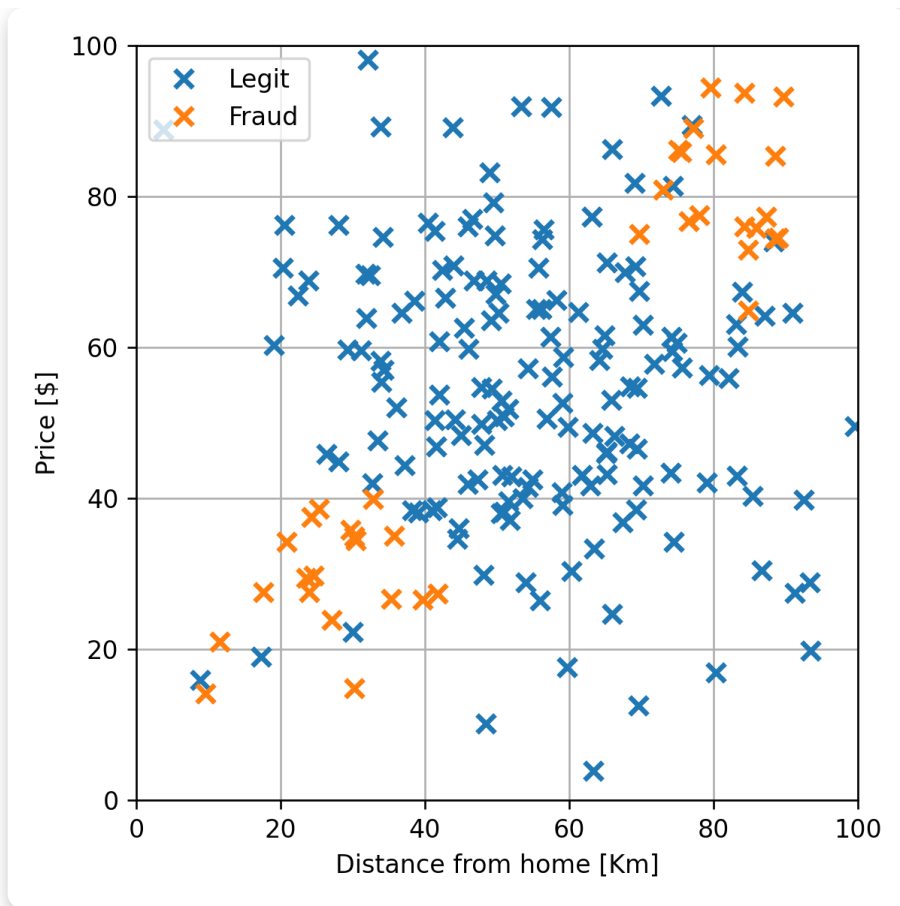
Code

מה נלמד היום



דוגמא להמחשת הבעיה בגישה הגנרטיבית פרמטרית

נסתכל שוב על הבעיה של חיזוי עסקאות שחשודות כהונאות אשראי:



ננסה להשתמש ב QDA על מנת להתייחס מודל לדגימות במדגם. נשתמש בנוסחאות לפרמטרים של מודל ה QDA:

$$p_y(0) = \frac{|\mathcal{I}_0|}{N} = 0.81$$

$$p_y(1) = \frac{|\mathcal{I}_1|}{N} = 0.19$$

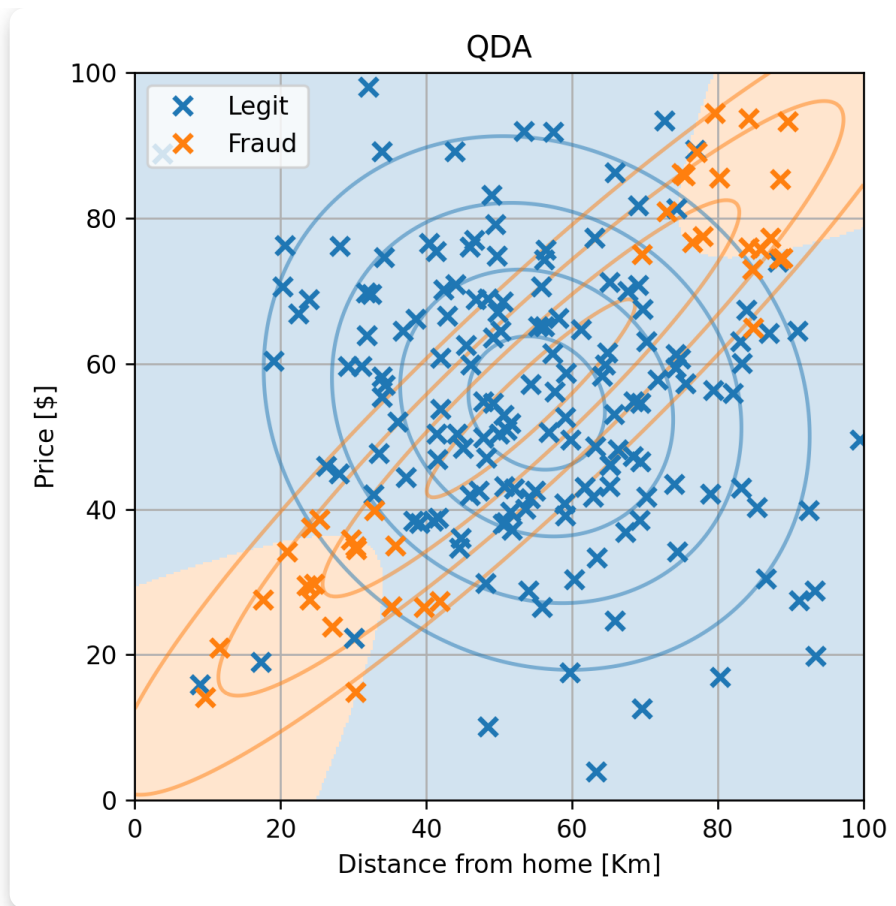
$$\boldsymbol{\mu}_0 = \frac{1}{|\mathcal{I}_0|} \sum_{i \in \mathcal{I}_0} \mathbf{x}^{(i)} = [55.1, 54.6]^T$$

$$\boldsymbol{\mu}_1 = \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} \mathbf{x}^{(i)} = [54.4, 55.2]^T$$

$$\Sigma_0 = \frac{1}{|\mathcal{I}_0|} \sum_i (\mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}}) (\mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}})^T = \begin{bmatrix} 350.9 & -42.9 \\ -42.9 & 336 \end{bmatrix}$$

$$\Sigma_1 = \frac{1}{|\mathcal{I}_1|} \sum_i (\mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}}) (\mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}})^T = \begin{bmatrix} 817.9 & 730.5 \\ 730.5 & 741.7 \end{bmatrix}$$

פרמטרים אלו מתארים כאמור את שני הגאוסיאנים שמתאימים לכל אחת משתי המחלקות בבעיה. נשרטט את הגאוסיאנים על המדגם יחד עם החזאי המתקבל מהמודל:



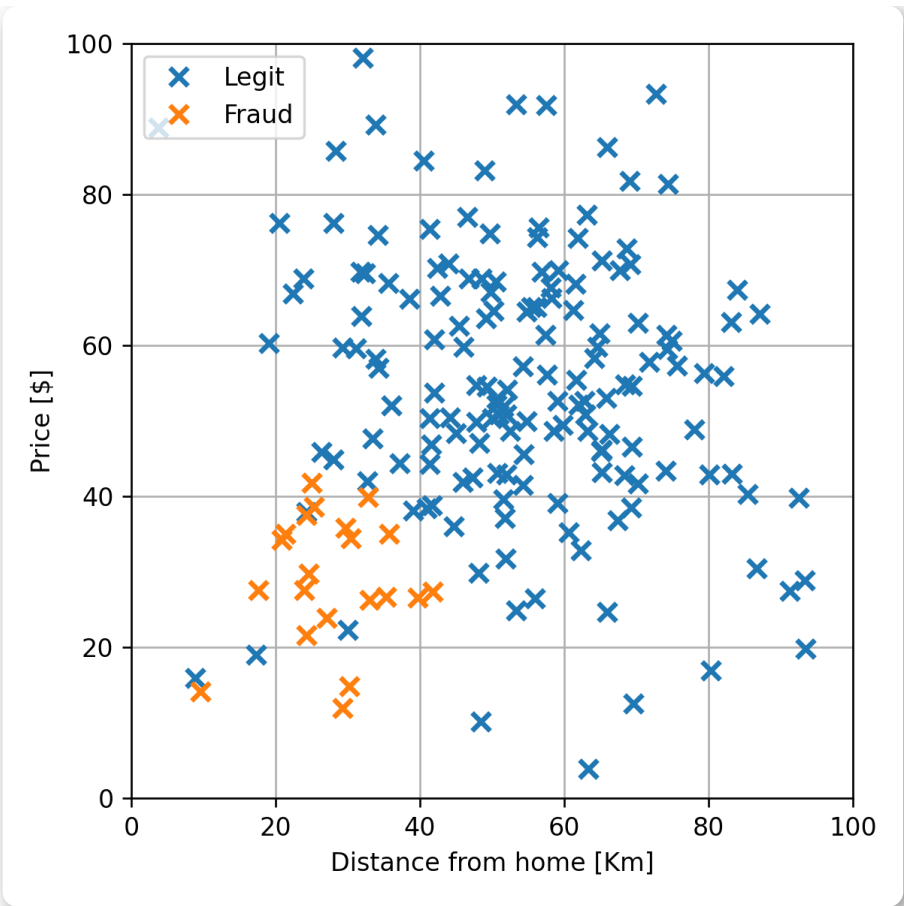
שגיאת החיזוי (miscalssification rate) על ה test set הינה 0.08.

התוצאה שקיבלנו אומנם סבירה, אך ניתן לראות מתוך השרטוט שהניסיון לייצג את הדגימות של ההונאות (הכתומות) על ידי גאוסיאן לא מאד מתאימה לפילוג שלהם בפועל. היינו רוצים לבחור במודל אשר יכול לייצג בנפרד את שני האיזורים השונים של הדגימות של ההונאות. לצערנו המבחר של המודלים בהם אנו יכולים להשתמש בגישה הגנרטיבית הפרמטרית לא גדול. כפי שציינו בהרצאה הקודמת המגבלה הזו נובעת מהצורך שהמודל הפרמטרי ייצג תמיד פילוגים חוקיים.

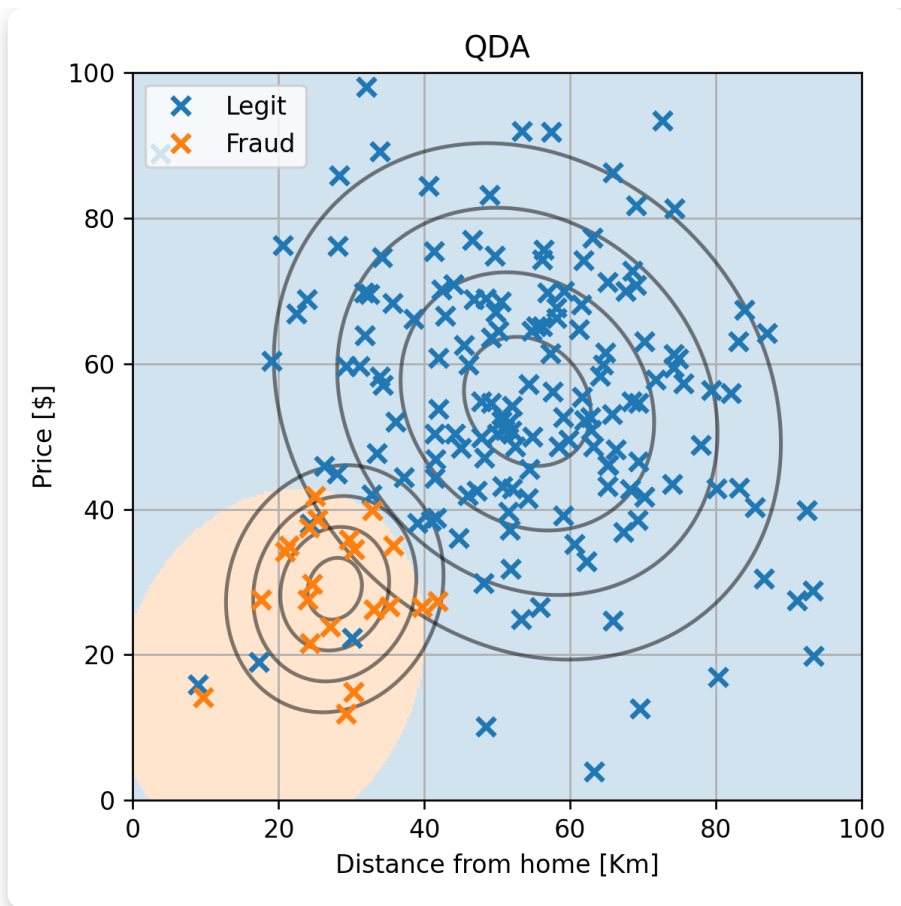
הערה: נושאים אלו לא מכוסים בקורס אך ישנו מודל פרמטרי נפוץ אשר מאד מתאים לתיאור פילוגים מסוג זה, המכילים מספר איזורים שונים בהם מרוכזים הדגימות. המודל נקרא Gaussian Mixture Model (GMM) והוא בנוי מקומביניציה של מספר גאוסיאנים. את ההתאמה של הפרמטרים של המודל הזה ניתן לעשות באמצעות שיטה אשר נקראת Expectation Maximization (EM).

דוגמא למדגם שמתאים למודל של QDA

רק לצורך הדגמה נסתכל על גירסא מעט שונה של המדגם שבה יש רק איזור אחד שבו נמצאות הדגימות של ההונאות:

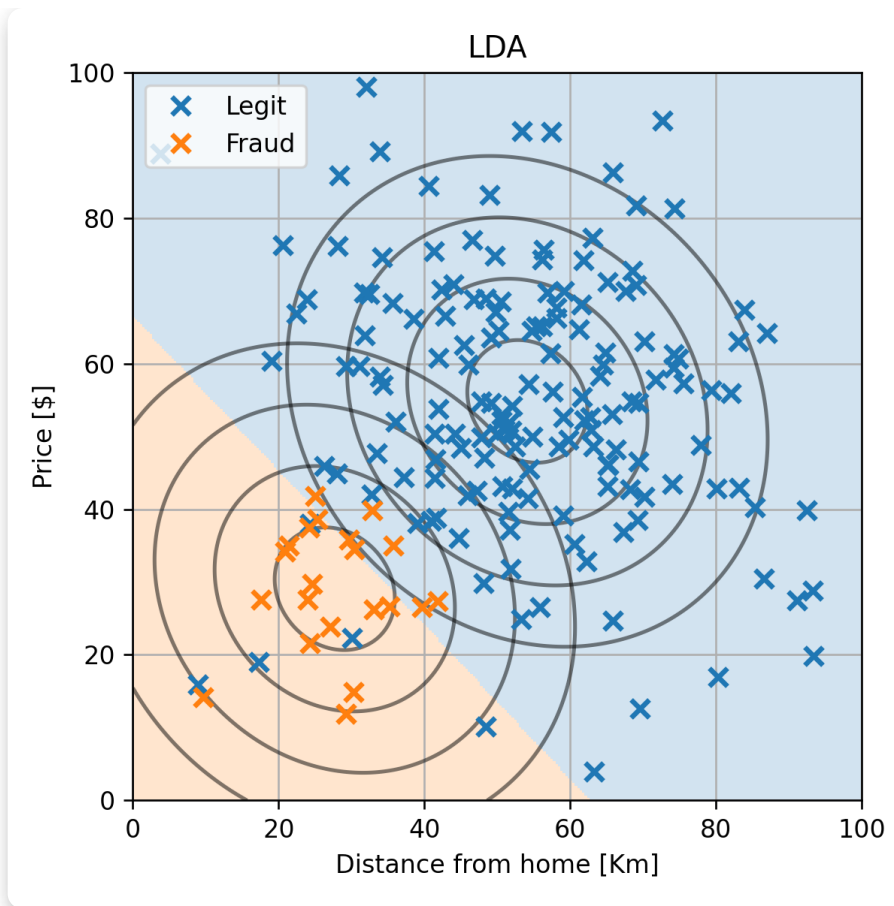


למדגם כזה המודל של QDA יתאים בצורה טובה:



שגיאת החיזוי (miscalssification rate) על ה test set במקרה הזה הינה 0. (זה לא שמודל זה לא יעשה טעויות, הם כנראה יהיו פשוט נדירות ויצא במקרה שעל ה test set הוא לא טועה).

רק לשם השוואה, נציג גם את התוצאה המתקבלת ממודל ה LDA:



בהרצאה זו נציג גישה אלטרנטיבית חדשה, אשר דומה לגישה הגנרטיבית אך מאפשרת חופש בבחירה גדול יותר של המודל הפרמטרי.

הגישה הדיסקרימינטיבית הסתברותית

עד כה למדנו על הגישה הדיסקרימינטיבית והגישה הגנרטיבית לבניית חזאים. נציג כעת גישה ביניים אותה נכנה גישה דיסקרימינטיבית הסתברותית.

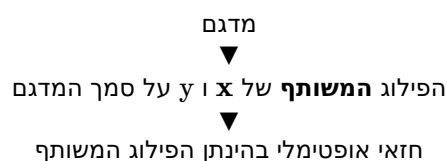
בגישה הגנרטיבית ניסינו ללמוד את הפילוג המלא של הדגימות במדגם. זאת אומרת, את הפילוג המשותף של x ו y . ראינו אבל שעבור רוב פונקציות המחיר (ספציפית עבור כל אלה שמוגדרים בפונקציית סיכון) החזאי האופטימלי יהיה תלוי רק בפילוג המותנה של y בהינתן x . בגישה הדיסקרימינטיבית הסתברותית ננסה ישירות ללמוד את הפילוג של y בהינתן x .

נציג את ההבדל הבין שלושת הגישות בעזרת השרטוט הבא:

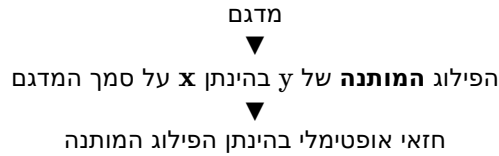
הגישה הדיסקרימינטיבית



הגישה הגנרטיבית



הגישה הדיסקרימינטיבית הסתברותית



ההתייחסות לגישה זו במקורות אחרים

בדומה לגישה הדיסקרימינטיבית הרגילה, גם גישה זו מוכוונת ישירות למציאת החזאי והיא לא מנסה ללמוד את התכונות המלאות של המדגם כפי שעושה הגישה הגנרטיבית. עקב כך גישה זו נחשבת לגישה דיסקרימינטיבית.

השם גישה דיסקרימינטיבית הסתברותית לא מופיע במקומות אחרים ואנו בחרנו לעשות שימוש בשם זה בקורס בכדי להבדיל את השיטה הזו מהגישה הדיסקרימינטיבית הרגילה. במרבית המקומות מציינים שלגישה הדיסקרימינטיבית אכן יש שתי תתי גישות, אך לא נותנים שם שונה לכל אחת משתי התתי הגישות האלה.

שימוש במודלים פרמטריים

הלמידה של הפילוג המותנה יעשה בגישה זו על ידי שיערוך פרמטרי (בדומה לשיערוך הפרמטרי בגישה הגנרטיבית). זאת אומרת שאנו נבחר מודל פרמטרי אשר יתאר את הפילוג המותנה, לרוב את הפונקציה $p_{y|x}(y|x)$, ונסה לשערך את פרמטרים של המודל בשיטות דומות לאלו של הגישה הגנרטיבית, כדוגמת MLE ו MAP.

נראה כעת שעל מנת למצוא את הפרמטרים של הפילוג המותנה בעזרת MLE ניתן פשוט לרשום את בעיית ה MLE הרגילה ולהחליף את הפילוג המשותף בפילוג המותנה.

נניח שאנו רוצים לשערך את הפרמטרים של מודל פרמטרי מסויים $p_{y|x}(y|x; \theta)$ על פי מדגם כל שהוא $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}$ בעזרת MLE. נשתמש בעובדה שבהינתן המודל הפרמטרי ניתן לרשום את הפילוג המשותף באופן הבא:

$$p_{\mathbf{x},y}(\mathbf{x}, y; \theta) = p_{y|x}(y|\mathbf{x}; \theta)p_{\mathbf{x}}(\mathbf{x})$$

עם פילוג שולי כל שהוא $p_{\mathbf{x}}(\mathbf{x})$ אשר אינו תלוי ב θ . בעיית האופטימיזציה שנרצה לפתור הינה:

$$\begin{aligned}
 \theta^* &= \arg \min_{\theta} - \sum_{i=1}^N \log \left(p_{\mathbf{x},y}(\mathbf{x}^{(i)}, y^{(i)}; \theta) \right) \\
 &= \arg \min_{\theta} - \sum_{i=1}^N \log \left(p_{y|x}(y^{(i)}|\mathbf{x}^{(i)}; \theta)p_{\mathbf{x}}(\mathbf{x}^{(i)}) \right) \\
 &= \arg \min_{\theta} - \sum_{i=1}^N \log \left(p_{y|x}(y^{(i)}|\mathbf{x}^{(i)}; \theta) \right) - \sum_{i=1}^N \log \left(p_{\mathbf{x}}(\mathbf{x}^{(i)}) \right) \\
 &= \arg \min_{\theta} - \sum_{i=1}^N \log \left(p_{y|x}(y^{(i)}|\mathbf{x}^{(i)}; \theta) \right)
 \end{aligned}$$

המעבר האחרון נובע מהעובדה שהסכום השני שבשורה הלפני אחרונה לא תלוי ב θ ולכן ניתן להתעלם ממנו. המשמעות של תוצאה זו היא שאנו יכולים למצוא את הפרמטרים של המודל של הפילוג המותנה ללא צורך לדעת או לשערך את הפילוג של \mathbf{x} .

ניתן כמובן להגיע לאותה תוצאה גם עבור משערך MAP, שם פשוט יופיע איבר ה prior על θ בנוסף לסכום המופיע בבעיית ה MLE.

היתרון של הגישה הדיסקרימינטיבית הסתברותית על הגישה הגנרטיבית

הבעיה בגנרטיבית הייתה הקושי למצוא מודלים שייצגו פילוגים משותפים חוקיים, ספציפית מודלים פרמטריים שמייצגים את $p_{\mathbf{x},y}(\mathbf{x}, y)$ צריכים לקיים את התנאים הבאים:

$$\begin{aligned}
 p_{\mathbf{x},y}(\mathbf{x}, y; \theta) &\geq 0 && \forall \mathbf{x}, y, \theta &.1 \\
 \int \int p_{\mathbf{x},y}(\mathbf{x}, y; \theta) d\mathbf{x} dy &= 1 && \forall \theta &.2
 \end{aligned}$$

את התנאי הראשון קל לקיים. אך התנאי השני הוא בעייתי.

בגישה הדיסקרימינטיבית הסתברותית הפונקציה שאותה נרצה למדל הינה $p_{y|x}(y|\mathbf{x})$. מסתבר שלמצוא מודלים פרמטריים חוקיים לפונקציה זו היא משימה פשוטה בהרבה, במיוחד עבור בעיות סיווג, בהם y הוא דיסקרטי וסופי. עבור בעיות סיווג התנאים על המודל יהיו:

$$p_{y|x}(y|\mathbf{x}; \theta) \geq 0 \quad \forall \mathbf{x}, y, \theta \quad .1$$

$$\sum_{y=1}^C p_{y|x}(y|\mathbf{x}; \theta) = 1 \quad \forall \mathbf{x}, \theta \quad .2$$

התנאי הראשון לא השתנה הרבה, אך בתנאי השני האינטגרל על כל הערכים של \mathbf{x} ו y התחלף בסכום על מספר סופי של איברים. זהו שינוי משמעותי. נראה כעת כיצד ניתן לבנות מודלים המקיימים תנאים אלו.

סיווג בינארי

נתחיל במקרה הפשוט של סיווג בינארי (בו $C = 2$). נסמן את שתי המחלקות ב 0 ו 1. במקרה זה התנאי השני יהיה:

$$p_{y|x}(0|\mathbf{x}; \theta) + p_{y|x}(1|\mathbf{x}; \theta) = 1 \quad \forall \mathbf{x}, \theta$$

דרך פשוטה לקיים תנאי זה הינה למצוא פונקציה פרמטרית כל שהיא $f(\mathbf{x}; \theta)$ אשר מחזירה ערכים בין 0 ל 1 ולהגדיר את המודל באופן הבא:

$$p_{y|x}(1|\mathbf{x}; \theta) = f(\mathbf{x}; \theta)$$

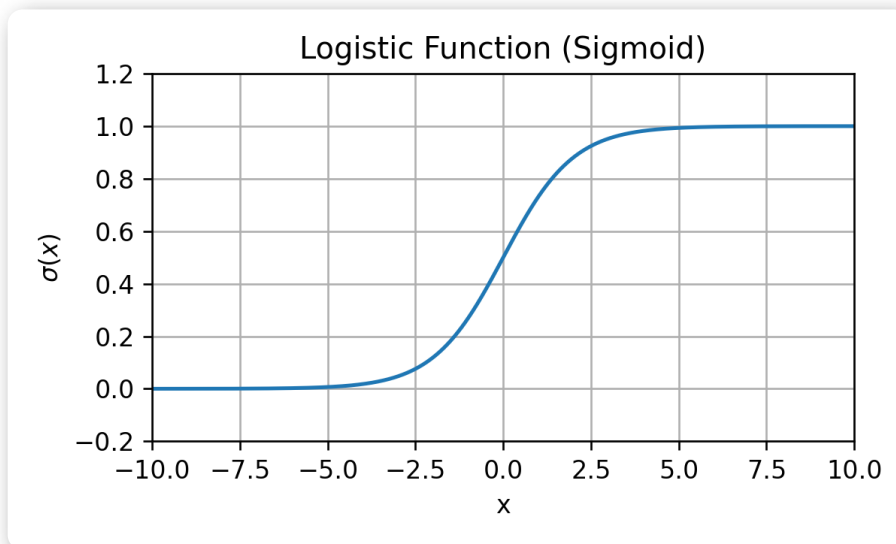
$$p_{y|x}(0|\mathbf{x}; \theta) = 1 - f(\mathbf{x}; \theta)$$

הפונקציה הלוגיסטית

הדרך הנפוצה ביותר לייצר פונקציה אשר מחזירה ערכים בין 0 ל 1 היא על ידי שימוש בפונקציה הנקראת הפונקציה הלוגיסטית. פונקציה זו מסומנת לרוב באות σ ומגדרת באופן הבא:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

הפונקציה נראית כך:



הערה: בתחום של מערכות לומדות מקובל לכנות את הפונקציה הזו **סיגמואיד (sigmoid)**.

בעזרת הפונקציה הלוגיסטית נוכל להפוך כל פונקציה פרמטרית (ללא שום מגבלה על הפונקציה) $f(\mathbf{x}; \theta)$ לפונקציה המחזירה ערכים בין 0 ל 1 על ידי הרכבה שלה עם הפונקציה הלוגיסטית: $\sigma(f(\mathbf{x}; \theta))$.

זאת אומרת, שכל מודל פרמטרי מהצורה:

$$p_{y|x}(1|\mathbf{x}; \boldsymbol{\theta}) = \sigma(f(\mathbf{x}; \boldsymbol{\theta}))$$

$$p_{y|x}(0|\mathbf{x}; \boldsymbol{\theta}) = 1 - \sigma(f(\mathbf{x}; \boldsymbol{\theta}))$$

יהיה מודל פרמטרי חוקי עבור f שמקבלת ערכים חיוביים ושלייים.

תכונות

חלק מהתכונות אשר הופכות את הפונקציה הלוגיסטית לבחירה נוחה כפונקציה אשר ממפה את התחום של $[-\infty, \infty]$ לתחום של $[0, 1]$ היא העובדה שהיא עושה זאת באופן רציף ומונוטוני עולה. בנוסף יש לה שתי תכונות מתמטיות נוספות שמקלות על העבודה איתה:

- היא מקיימת את הקשר הבא: $1 - \sigma(z) = \sigma(-z)$ אשר מסייע במקרים רבים לפשט ביטויים.
- לנגזרת של הלוג של הפונקציה יש צורה פשוטה: $\frac{d}{dz} \log(\sigma(z)) = \sigma(z)(1 - \sigma(z))$. הדבר לרוונטי מאד כאשר מחשבים את הנגזרת של ה log-likelihood.

Binary Logistic Regression

ב Binary Logistic Regression נשתמש במודל שהצגנו קודם:

$$p_{y|x}(1|\mathbf{x}; \boldsymbol{\theta}) = \sigma(f(\mathbf{x}; \boldsymbol{\theta}))$$

$$p_{y|x}(0|\mathbf{x}; \boldsymbol{\theta}) = 1 - \sigma(f(\mathbf{x}; \boldsymbol{\theta}))$$

(עם פונקציות פרמטריות כל שהם) ונמצא את הפרמטרים של המודל בעזרת MLE (או בעזרת MAP אשר מוסיף איבר גולריזציה). בעיית האופטימיזציה שיש לפתור במקרה זה הינה:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} - \sum_{i=1}^N \log(p_{y|x}(y^{(i)}|\mathbf{x}^{(i)}; \boldsymbol{\theta}))$$

$$= \arg \min_{\boldsymbol{\theta}} - \sum_{i=1}^N I\{y^{(i)} = 1\} \log(\sigma(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}))) + I\{y^{(i)} = 0\} \log(1 - \sigma(f(\mathbf{x}^{(i)}; \boldsymbol{\theta})))$$

דרך נפוצה נוספת לרשום את בעיית האופטימיזציה הזו עושה שימוש בעובדה שבמקרה של משתנה בינארי מתקיים:

$$I\{y = 1\} = y$$

$$I\{y = 0\} = (1 - y)$$

מכאן שניתן לרשום את בעיית האופטימיזציה באופן הבא:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} - \sum_{i=1}^N y^{(i)} \log(\sigma(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}))) + (1 - y^{(i)}) \log(1 - \sigma(f(\mathbf{x}^{(i)}; \boldsymbol{\theta})))$$

(בתרגול נראה דרך נוספת לרשום את בעיית האופטימיזציה הזו).

במרבית המקרים לא ניתן יהיה לפתור את בעיית האופטימיזציה הזו באופן אנליטי (על ידי גזירה והשוואה ל-0) ואנו נחפש את הפתרון בשיטות נומריות כגון אלגוריתם ה gradient descent עליו נרחיב בהמשך ההרצאה.

עבור פונקציית מחיר של misclassification rate החזאי האופטימלי יהיה:

$$h(\mathbf{x}) = \arg \max_y p_{y|x}(y|\mathbf{x}; \boldsymbol{\theta}) = \begin{cases} 1 & \sigma(f(\mathbf{x}; \boldsymbol{\theta})) > 0.5 \\ 0 & \text{else} \end{cases} = \begin{cases} 1 & f(\mathbf{x}; \boldsymbol{\theta}) > 0 \\ 0 & \text{else} \end{cases}$$

סיווג לא בינארי

ניתן להרחיב את השיטה לבניית מודלים בעזרת הפונקציה הלוגיסטית גם למקרה שבו y אינו משתנה אקראי בינארי. הדרך לעשות זאת הינה באמצעות פונקציית ה softmax.

פונקציית ה Softmax

פונקציית ה softmax היא מעיין הרחבה של הפונקציה הלוגיסטית מהמקרה של משתנה אקראי בינארי למשתנה אקראי דיסקרטי סופי אשר מקבל אחד מ C ערכים. נסמן את הערכים אלו ב $\{1, 2, \dots, C\}$. פונקציית ה softmax לוקחת וקטור כלשהו \mathbf{z} באורך C ומייצרת ממנו וקטור חדש אשר יכול לייצג פילוג דיסקרטי חוקי. הפונקציה מוגדרת באופן הבא:

$$\text{softmax}(\mathbf{z}) = \frac{1}{\sum_{c=1}^C e^{z_c}} [e^{z_1}, e^{z_2}, \dots, e^{z_C}]^T$$

או לחילופין, הערך של האיבר ה i של הפונקציה הינו:

$$\text{softmax}(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{c=1}^C e^{z_c}}$$

ניתן להסתכל על הפונקציה הזו כעל פונקציה המבצעת את שתי השלבים הבאים:

1. על מנת להפוך את כל רכיבי הוקטור לחיוביים, כל איבר בוקטור z_i מוחלף באקספוננט שלו e^{z_i} .
2. בכדי שסכום הערכים של הוקטור יהיה אחד מנרמלים את הוקטור על ידי חלוקת איברי הוקטור בסכום האיברים: $\sum_c e^{z_c}$.

תכונות

- אינווריאנטיות לתוספת של קבוע (לכל אברי הוקטור): $\forall i: \text{softmax}(\mathbf{z} + \mathbf{a})_i = \text{softmax}(\mathbf{z})_i$.
 - $\frac{\partial}{\partial z_j} \log(\text{softmax}(\mathbf{z}))_i = \delta_{i,j} - \text{softmax}(\mathbf{z})_j$.
- ($\delta_{i,j} = I\{i = j\}$ קרונוקר של הדלתא)

הפונקציה הלוגיסטית כמקרה פרטי

עבור וקטור באורך 2: $\mathbf{z} = [a, b]$, נקבל:

$$\begin{aligned} \text{softmax}(\mathbf{z})_1 &= \frac{e^a}{e^a + e^b} = \frac{1}{1 + e^{b-a}} = \sigma(a - b) \\ \text{softmax}(\mathbf{z})_2 &= \frac{e^b}{e^a + e^b} = 1 - \sigma(a - b) \end{aligned}$$

Non-Binary) Logistic Regression)

בדומה לפונקציה הלוגיסטית נוכל להשתמש בפונקציית ה softmax על מנת לבנות פילוגים חוקיים של משתנים דיסקרטיים סופיים. עבור C פונקציות פרמטריות כלשהן, $f_c(\mathbf{x}; \theta_c)$, ניתן לבנות מודל פרמטרי חוקי לפילוג המותנה באופן הבא:

$$p_{y|\mathbf{x}}(y|\mathbf{x}; \theta) = \frac{e^{f_y(\mathbf{x}; \theta_y)}}{\sum_{c=1}^C e^{f_c(\mathbf{x}; \theta_c)}}$$

לשם נוחות נסמן:

- את הוקטור θ כוקטור אשר כולל את כל C וקטורי הפרמטרים: $\theta = [\theta_1^T, \theta_2^T, \dots, \theta_C^T]^T$.
- את הפונקציה \mathbf{f} כפונקציה המאגדת את כל C הפונקציות הפרמטריות: $\mathbf{f}(\mathbf{x}; \theta) = [f_1(\mathbf{x}; \theta_1), f_2(\mathbf{x}; \theta_2), \dots, f_C(\mathbf{x}; \theta_C)]^T$

בעזרת סימונים אלו נוכל לרשום את המודל הפרמטרי באופן הבא:

$$p_{y|\mathbf{x}}(y|\mathbf{x}; \theta) = \text{softmax}(\mathbf{f}(\mathbf{x}; \theta))_y$$

משערה ה MLE של מודל זה יהיה נתון על ידי:

$$\begin{aligned} \theta^* &= \arg \min_{\theta} - \sum_{i=1}^N \log(p_{y|\mathbf{x}}(y^{(i)}|\mathbf{x}^{(i)}; \theta)) \\ &= \arg \min_{\theta} - \sum_{i=1}^N \log(\text{softmax}(\mathbf{f}(\mathbf{x}^{(i)}; \theta))_{y^{(i)}}) \end{aligned}$$

היתירות בייצוג של מודל ה logistic regression

בדומה למקרה הבינארי שבו לא היינו צריכים להגדיר 2 פונקציות פרמטריות, אחת ל $y = 0$ ואחת ל $y = 1$, גם במקרה הכללי מספיק להגדיר $C - 1$ פונקציות פרמטריות. הדבר נובע מהעובדה שאם מגדירים את ההסתברות של $C - 1$ מחלקות, המחלקה הנותרת תקבע באופן מוחלט כך שהיא תשלים את ההסתברות ל-1.

באופן דומה, אם נוסיף את אותו הערך לכל הפונקציות f_c לא נשנה את הפילוג המתקבל. במילים אחרות, כל שינוי מהצורה של $f_c(\mathbf{x}; \theta_c) \rightarrow f_c(\mathbf{x}; \theta_c) + g(\mathbf{x})$ לא ישנה את הפילוג המותנה $p_{y|x}(y|\mathbf{x}; \theta)$.

במקרים מסויימים נרצה לבטל יתירות זו. ניתן לעשות זאת על ידי קיבוע של אחת הפונקציות הפרמטריות, לרוב הראשונה $c = 1$, כך שהיא תהיה שווה זהותית ל 0: $f_1(\mathbf{x}; \theta_1) = 0$. שינוי שכזה כאמור לא יפגע ביכולת הייצוג של המודל ויבטל את היתירות שיש בייצוג של כל פילוג. בחירה כזו גם תקטין את מספר הפרמטרים שיש ללמוד.

Linear Logistic Regression

הגרסא הלינארית של הרגרסיה הלוגיסטית היא המקרה שבו בוחרים את הפונקציות הפרמטריות להיות פונקציות לינאריות:

$$f_c(\mathbf{x}; \theta_c) = \theta_c^T \mathbf{x}$$

במקרה זה פונקציית ה objective היא קמורה (convex) ולכן מובטח ש gradient descent, במידה והוא מתכנס, יתכנס למינימום גלובלי.

Gradient descent (שיטת הגרדיאנט)

בבעיות אופטימיזציה בהם לא ניתן להגיע לפתרון סגור על ידי גזירה של פונקציית המטרה והשוואה ל-0 נאלץ להשתמש בשיטות נומריות. אחת השיטות הנפוצות ביותר במערכות לומדות (בעיקר בעבודה עם רשתות נוירונים, אך לא רק) לפתרון בעיות אופטימיזציה הינו אלגוריתם ה gradient descent אותו הצגנו בקצרה בתרגול 1. נציג אותו כעת בצורה יותר מפורטת

הרעיון מאחרי Gradient descent הוא פשוט. האלגוריתם מנסה למצוא מינימום מקומי של פונקציית המטרה על ידי התקדמות בצעדים קטנים בכיוון שבו הפונקציה יורדת הכי מהר. אילוסטרציה:



לאלגוריתמים איטרטיביים מסוג זה, אשר מנסים בכל איטרציה לשפר את מצבם לעומת המצב הנוכחי (מבלי התייחס צורה הגלובאלית של הפונקציה) אנו קוראים אלגוריתמים חמדנים (greedy). כפי שצינו קודם אלגוריתמים כאלה לא מתיימרים להתכנס לאופטימום הגלובאלי, אלא רק ינסו להשתפר כל הזמן עד אשר יגיעו לאופטימום מקומי. הדרישה היחידה על הבעיה לשם השימוש באלגוריתם הינה היכולת לחשב את הנגזרת של פונקציית המטרה.

עבור בעיית המינימיזציה:

$$\arg \min_{\theta} g(\theta)$$

האלגוריתם פועל באופן הבא:

- מאתחלים את $\theta^{(0)}$ לנקודה אקראית כל שהיא.

- חוזרים על צעד העדכון הבא עד שמתקיים תנאי עצירה כל שהוא:

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_{\theta} g(\theta^{(t)})$$

את הפרמטר η יש לקבוע מראש, והוא יקבע את גודל הצעדים שהאלגוריתם יעשה.

תנאי עצירה

ישנם מספר דרכים להגדיר תנאי עצירה לאגוריתם:

- הגעה למספר צעדי עדכון שנקבע מראש: $t > \text{max-iter}$.
- כאשר הנורמה של הגרדיאנט קטנה מערך סף מסויים שנקבע מראש: $\|\nabla_{\theta} g(\theta)\|_2 < \epsilon$.
- כאשר השיפור בפונקציית המטרה קטן מערך סף מסויים שנקבע מראש: $g(\theta^{(t-1)}) - g(\theta^{(t)}) < \epsilon$.
- שימוש בעצירה מוקדמת על מנת להתמודד עם התאמת יתר (נרחיב על כך בהרצאה הבאה)

הבעיות של האלגוריתם

התכנסות למינימום מקומי ותלות באיתחול

כפי שצינו האלגוריתם הוא אלגוריתם חמדן אשר מתכנס למינימום מקומי אשר תלוי באיתחול של האלגוריתם.

לא ניתן לקבוע בוודאות האם האלגוריתם התכנס

בעיה נוספת של האלגוריתם הינה שלרוב לא נוכל לדעת האם האלגוריתם הגיע לנקודת מינימום כל שהיא או שהוא הגיע לאיזור שבו השיפוע קטן והוא עדיין מתקדם לאיטו לכיוון המינימום.

בעיית הבחירה של גודל הצעד

שני הבעיות הקודמות אומנם מונעות מאיתנו להגיע לנקודה האופטימלית, אך הם עדיין לא מפריעות לאלגוריתם להניב תוצאות טובות. הבעיה העיקרית של האלגוריתם הינה הבעיה של בחירת גודל צעד אשר עלולה למנוע מהאלגוריתם להניב תוצאות רלוונטיות תוך מספר סביר של צעדים.

על מנת להדגים את הבעיה נצא לטיול בנחל יהודיה. בדומה לאפיק של נחל יהודיה גם בבעיות אופטימיזציה לרוב יהיו במרחב כיוונים שבהם השיפועים / גרדיאנטים יהיו גדולים כמו הכיוונים הניצבים לכיוון זרימת הנחל (הדפנות של האפיק) וכיוונים שבהם השיפוע יהיה קטן, כמו הכיוון שבו הנחל זורם.

נניח ואנו מתחילים מאחת השפות של הנחל ואנו רוצים להתקדם במורד הנחל (ולהגיע לכינרת). נתחיל עם בחירה של גודל צעד קטן אשר יורד לאט ומגיע לאפיק הנחל:



בתחילת הדרך, כאשר אנו על הדפנות של הנחל, השיפוע / הגרדיאנט יהיה גדול ויצביע לכיוון האפיק (המקום שבו המים זורמים). בשלב זה לא נרצה שהפרמטר η יהיה גדול מידי בכדי שלא נקח צעדים גדולים מידי יגרמו לנו לעבור לצידו השני של האפיק. במקרה זה הבעיה תתחיל כאשר נגיע לאפיק עצמו שם השיפוע / גרדיאנט יהיה קטן מה שיגרום לכך שנתקדם בצעדים מאד קטנים ותהליך ההתכנסות יהיה מאד איטי. תופעה זו תרגום לרוב לכך שהאלגוריתם לא ידע לתוצאה משמעותית במספר צעדים סביר.

נרצה אם כך להגדיל את גודל הצעד אך כם גודל צעד גדול הוא בעייתי. נראה מה קורה כאשר אנו מנסים לעשות זאת:



כאשר נגדיל את גודל הצעד אנחנו נהיה בבעיה דווקא בשלב הראשון של האלגוריתם שבו אנו נמצאים עדיין נמצאים על דפנות הנחל שם הגרדיאנטים גדולים. הגדלה של η יגרום לכך שנתקשה להגיע לאפיק עצמו ואנו נעשה צעדים גדולים מידי בכיוון הניצב לכיוון שבו זורם הנחל. גם מצב זה יגרום לכך שהאלגוריתם לא יצליח להתכנס לתוצאה סבירה במספר צעדים סביר.

ננצל את הדוגמא הזו להראות עוד מקרה אחד שבו אנו מגדילים את η אף יותר:



בחירה של η גדול מאד לא רק שתמנע מאיתנו להגיע לאפיק עצמו, אלא עלולה אף להרחיק אותנו ממנו יותר בכל צעד ולגרום לאלגוריתם להתבדר.

עקב בעיה זו אלגוריתם ה gradient descent בצורתו הפשוטה כפי שהוצגה כאן אינו מאד שימושי שכן הוא לרוב לא יצליח להניב תוצאות טובות במספר סביר של צעדים. למזלנו ישנם מספר שיפורים קלים שניתן לעשות לאלגוריתם על מנת להתמודד בצורה טובה עם בעיה זו. לצערנו בקורס זה לא נספיק לכסות שיפורים אלו. אנו נציין השניים מהשיפורים הנפוצים ביותר בתחום::

1. הוספה של רכיב תנע (מומנטום) לאלגוריתם
2. שימוש בגודל צעד אשר משתנה במהלך הריצה

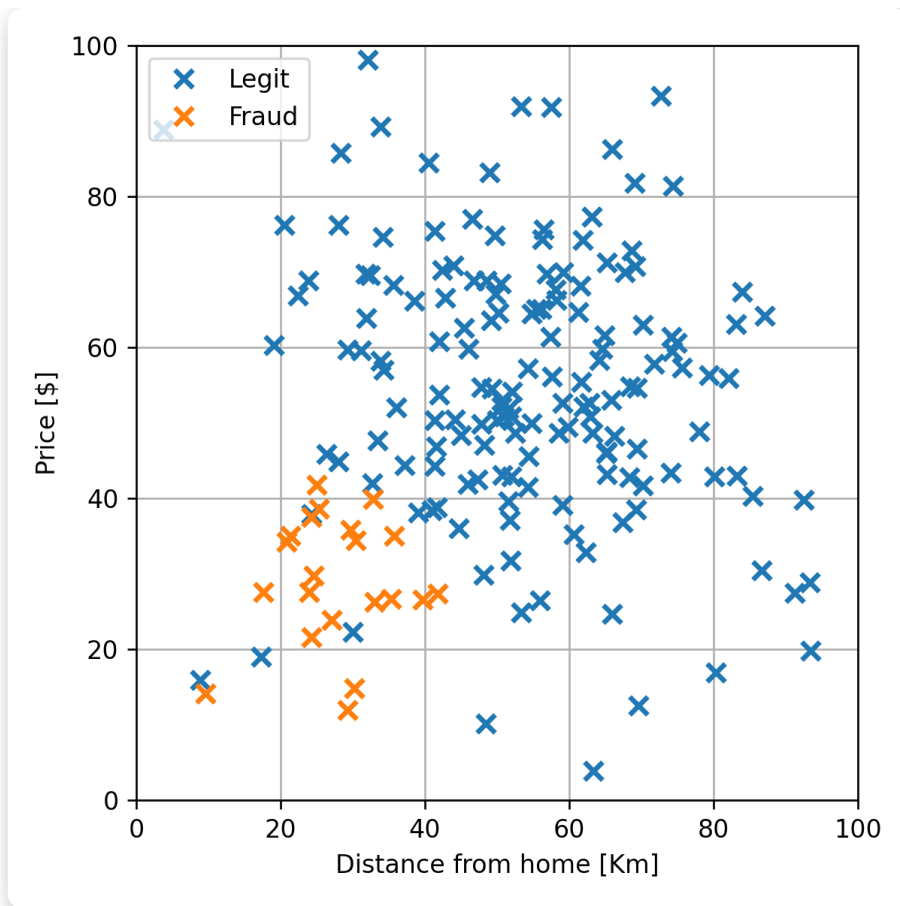
שני מקורות מצויינים לקריאה על נושא זה הם שתי הכתבות הבאות:

1. [An overview of gradient descent optimization algorithms](#)
2. [Why Momentum Really Works](#) (כתבה זו מכילה דוגמאות אינטרקטיביות מצויינות אשר עוזרות להבין את הבעיה והפתרון)

בתרגיל הבית הרטוב אתם תשתמשו במימוש קיים של גרסא משופרת נפוצה של אלגוריתם ה gradient descent בשם ADAM. שיטה זו עושה שימוש מתוחכם בתנע בשביל להתמודד עם הבעיה של גודל הצעד.

דוגמא: Linear Logistic Regression

נחזור לבעיה של חיזוי עסקאות החשודות כהונאות אשראי.



נשתמש במודל של linear logistic regression על מנת למדל את הפילוג המתונה:

$$p_{y|x}(y|\mathbf{x}; \boldsymbol{\theta}) = \begin{cases} \sigma(\mathbf{x}^\top \boldsymbol{\theta}) & y = 1 \\ 1 - \sigma(\mathbf{x}^\top \boldsymbol{\theta}) & y = 0 \end{cases}$$

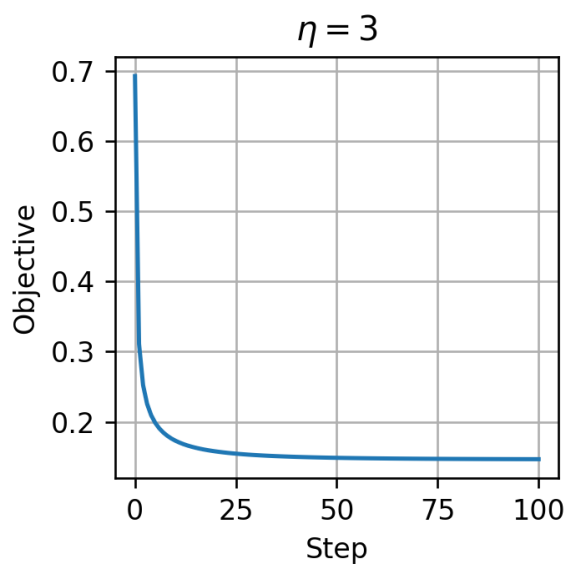
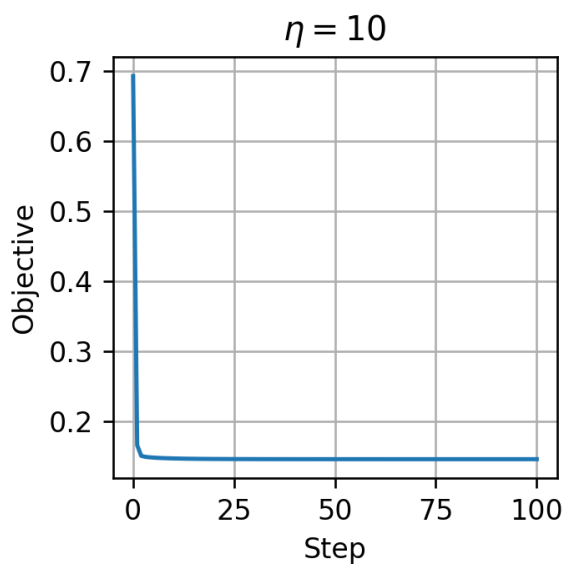
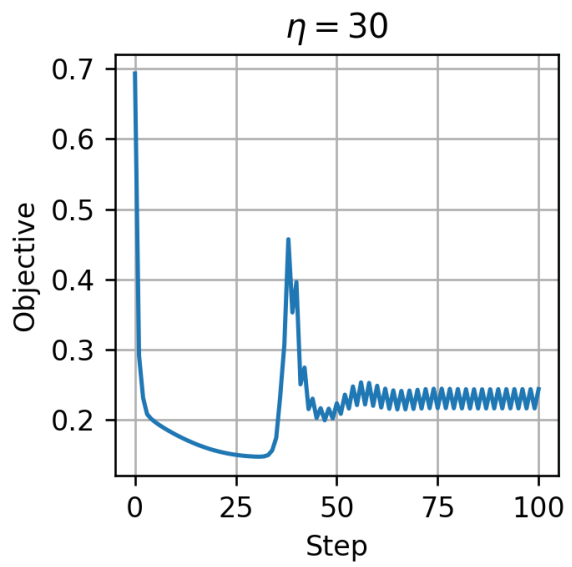
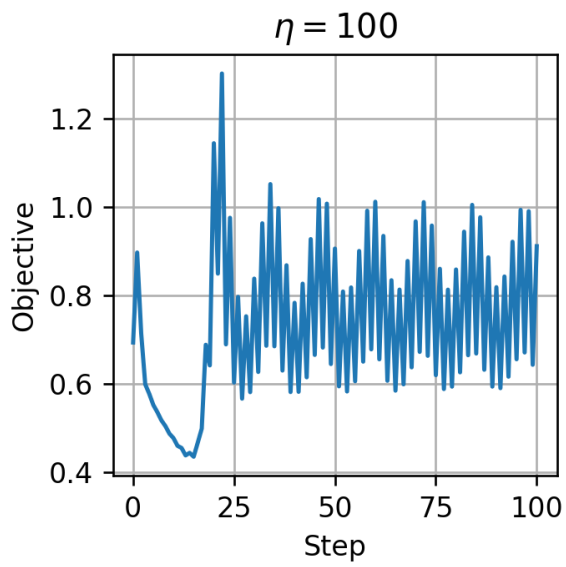
נמצא את הפרמטרים של המודל בעזרת MLE:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} - \sum_{i=1}^N y^{(i)} \log(\sigma(\mathbf{x}^{(i)\top} \boldsymbol{\theta})) + (1 - y^{(i)}) \log(1 - \sigma(\mathbf{x}^{(i)\top} \boldsymbol{\theta}))$$

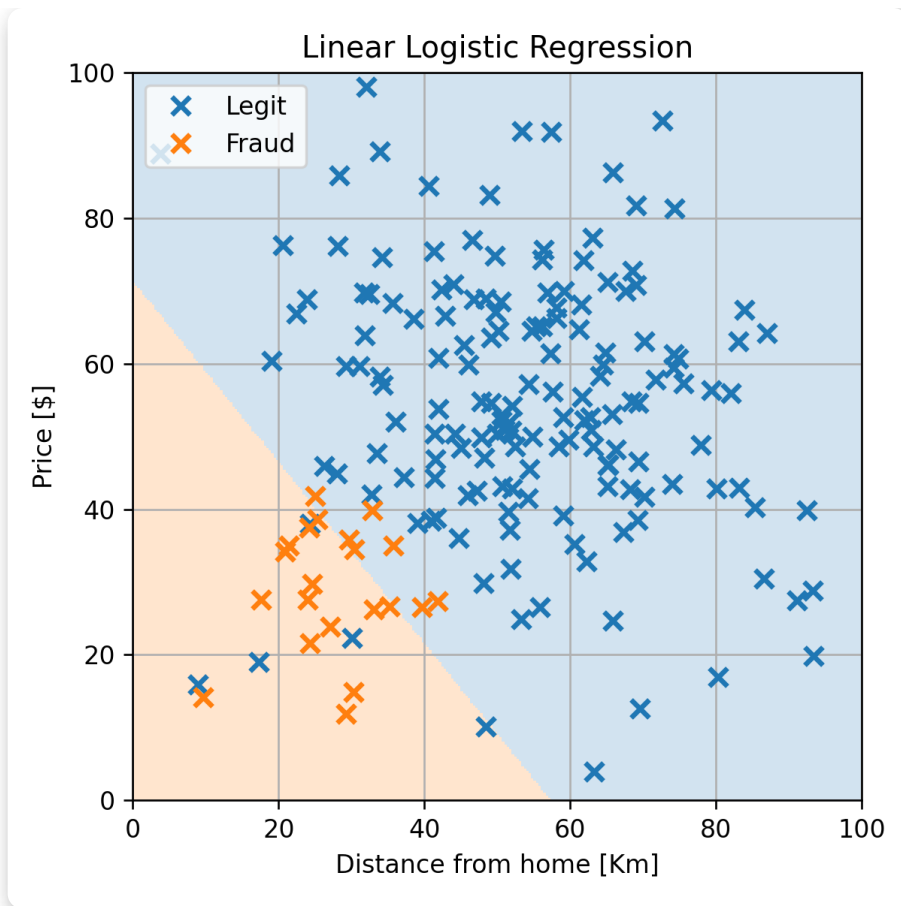
כלל העדכון של האלגוריתם יהיה:

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \eta \sum_{i=1}^N \left(y^{(i)} (1 - \sigma(\mathbf{x}^{(i)\top} \boldsymbol{\theta})) - (1 - y^{(i)}) \sigma(\mathbf{x}^{(i)\top} \boldsymbol{\theta}) \right) \mathbf{x}^{(i)}$$

בכדי לבחור את η נריץ את האלגוריתם מספר קטן של צעדים ונסתכל על הערך של ה objective כפונקציה של מספר הצעדים עבור ערכי η שונים:



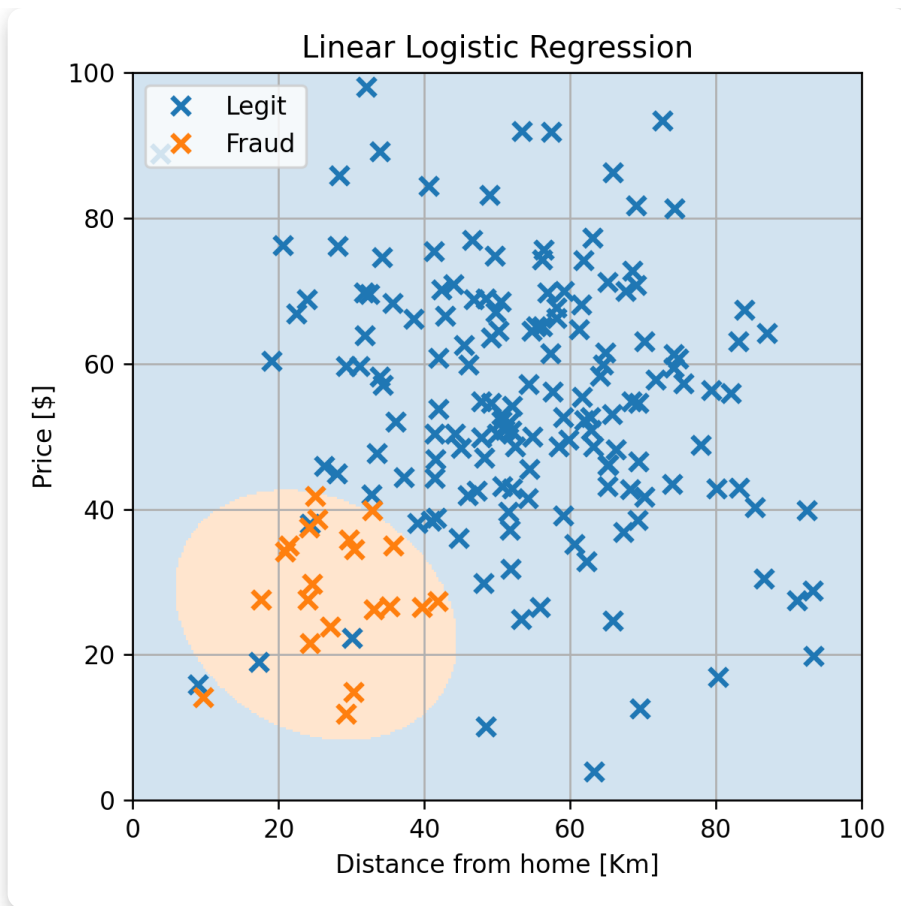
הערכים של $\eta = 30$ ו $\eta = 100$ מתאימים בדיוק למקרה של η גדול מידי כפי שראינו קודם ולכן נבחר את $\eta = 10$.
 נריץ את האלגוריתם עם בחירה זו של η ונקבל את החזאי הבא:



אשר מניב תוצאה של misclassification rate של 0.02 על ה test set.

שימוש במודל מסדר גבוה יותר

כפי שציינו קודם, היתרון של הגישה הדיסקרימינטיבית ההסתברותית היא שנוכל להשתמש בכל מודל שנרצה. לדוגמה נוכל להחליף את $f(\mathbf{x}; \theta)$ להיות פולינום מסדר שני ונקבל את התוצאה הבאה:



עם misclassification rate של 0 על ה test set.

בהרצאה הבאה נראה כיצד ניתן להשתמש ברשתות נוירונים כפונקציה הפרמטרית $f(\mathbf{x}; \boldsymbol{\theta})$.