

מבוא לבינה מלאכותית - תרגול 1.

rgb יחזקאל אימרה.

4 בנובמבר 2025

- יש הרבה סוגים למידה בפועל, הנה חלוקה:
1. למידה מפוקחת: לכל דוגמה יש את התווית שלה (label).
 - (א) סיווג: חיזוי קטגוריות.
 - (ב) רגסיה: חיזוי ערך רציף.
 2. למידה לא מפוקחת: יש דוגמאות, אין תוויות.
 - (א) אשכול: חלוקת דוגמאות לקבוצות.
 - (ב) הורדת ממדדים: הפתוחה כמוות הפיצרים של המודל.
 - (ג) זיהוי אנומליות: גילוי דוגמאות חריגיות מבין סך כל הדוגמאות.
 3. למידה באמצעות חיזוקים: המודל לומד דרך ניסוי וטעייה. המודל מבצע פעולות ומתקבל חיזוקים (חיוביים או שליליים) ומשפר את המדייניות שלו.
 4. למידה מונחית עצמית: המערכת מייצרת לעצמה תוויות מתוך הנתונים כדי ללמידה "צוגים יותר טובים".
 5. למידה גנרטיבית: המערכת מייצרת נתונים חדשים מהתיקה מציעה קורס בלמידה לא מפוקחת, והמחלקה למדעי המחשב מציעה קורסים בשער סוגים הלמידה.

חלק I

מודלי רגסיה.

1 רגסיה לינארית.

לאלו שלא זוכרים, בלמידה מפוקחת יש שתי נישות מרכזיות: קלסיפיקציה ורגסיה.

הגדרה 1.1. מודל רגסיה הוא מודל סטטיסטי הנועד להעריך קשר בין משתנים.

דוגמה 2.1. דוגמאות למודל רגסיה:

1. חיזוי מחיר של בית לפי גודל ומיקום.

2. חיזוי ציון של סטודנט לפי מספר שעות למידה.

3. חיזוי הטמפרטורה של מחר על סמך נתוני מזג האוויר של הימים האחרונים.

איך עובד מודל רגסיה?

(1) הנתות המודל: נניח יש לנו אוסף של N דוגמאות והתיוגים שלהם, $\{(x_i, y_i)\}_{i=1}^N$ כאשר $x_i \in \mathbb{R}^d$ וכן $y \in \mathbb{R}$. רגסיה לינארית, כשמה כן היא, מניחה כי הקשר בין x ו- y הוא לינארי, כלומר שקיים d סקלרים b ו- w_i כך שמתקיים $y_i = \sum_{n=1}^d x_{ij}w_j + b$.

השאלה, איך אנחנו מגלים מיהם הסקלרים b ו- w_i הנ"ל?

(2) פונקציית הפסד: לכל דוגמה (x_i, y_i) המודל חוצה ערך כלשהו \hat{y}_i . נגידו את השגיאה של המודל להיות $y_i - \hat{y}_i$. אנחנו מעוניינים במספר אחד שייהי תלוי בכל השגיאות שלנו ושיגיד לנו כמה רע המודל שלנו. יש לכך כמה בחירות טבעיות:

• סכום השגיאות: $(\hat{y}_i - y_i)^2$. הבעה היא שפה אם יש לי שנייה חיובים ושגיאה שלילית הן תבטלו, מה שיכל להוביל לדיווחים שגויים.

• סכום ערכים מוחלטים: $|\hat{y}_i - y_i|$. יותר טוב, אבל פה הבעה היא שזה לא גיר - הרי אם אנחנו רוצים לעזור את השגיאה הכללת שלנו נצטרך לגזר אותה.

• סכום השגיאות בריבוע: $(\hat{y}_i - y_i)^2$. מעולה בשbillנו. זה גיר ונוטן יותר משקל לכל שהשגאה גדולה.

הגדרה 1.3. פונקציית הפסד (loss function) היא פונקציה המappa את התחזית של המודל ואת הערך האמתי למספר ממשי, המציג את מידת הטעות או העלות של התחזית.

באופן טבעי, ככל שפונקציית ההפסד מניבה ערך נמוך יותר עברו תחזיות המודל, כך ביצועי המודל טובים יותר. מטרתנו היא לעזור את ערך פונקציית הפסד.

הגדרה 4.1. נגידר את ה-*loss function* טעوت הריבועית המומוצעת (MSE) להיות

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

כאשר θ הם הפרמטרים של המודל שלנו שאנו מראנו מעתה ללמידה (למוד) = למאער את ה-*loss* שלנו, ו- \hat{y} זה הערך שהמודל פולט לדגימה x_i עם פרמטרים θ . נרצה למצאו $\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta)$.

זהו פונקציית הפסד סטנדרטית בעיות רגרסיה בכללית לאו דווקא ברגRESSED LINEARITY.

(3) תהליך הלמידה: נציין לפרוטוקול כי קיים פתרון סגור לבעה זו: כתוב את הדאטה שלנו בכתיב מטריציוני:

$$X = \begin{pmatrix} 1 & - & x_1^T & - \\ 1 & - & x_2^T & - \\ \vdots & & \vdots & \\ 1 & - & x_N^T & - \end{pmatrix} \text{ המוגדר להיות } X \in \mathbb{R}^{N \times (d+1)} \bullet$$

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \text{ כל הערכים שלנו בצורה וקטור } y \in \mathbb{R}^N \bullet$$

$$\beta = \begin{pmatrix} b \\ w_1 \\ \vdots \\ w_d \end{pmatrix} \text{ כל הפרמטרים שאנו רוצים ללמידה } \beta \in \mathbb{R}^{d+1} \bullet$$

כלומר המודל שלנו הופך להיות

$$\hat{y} = X\beta$$

ופונקציית ההפסד שלנו היא

$$\mathcal{L}(\beta) = \frac{1}{N} \|y - X\beta\|_2^2$$

כדי למאער, נזכיר ביחס ל- β ונשווה לאפס:

$$\nabla_{\beta} \mathcal{L} = -\frac{2}{N} X^T (y - X\beta)$$

כלומר קיבל

$$X^T y = X^T X \beta$$

ולפי גאוס הפתרון יהיה

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

הבעיה: לרוב, המטריצה $X^T X$ לא הפיכה, לכן נצטרך להיות יותר חכמים בחיפוש שלנו אחריו $\hat{\beta}$.

הפתרון: ניעזר באלגוריתם מورد הגרדיאנט:

אלגוריתם 5.1. אלגוריתם מורד הגרדיאנט:

• צחزو נקודות התחלת θ_0 .

• חשבו את $\nabla_{\theta} \mathcal{L}(\theta_0)$.

• צחزو גודל צעד $\eta > 0$.

• חשבו $\theta_{k+1} = \theta_k - \eta \nabla_{\theta} \mathcal{L}(\theta_k)$.

בצורת כתיבה יותר כללית כתוב

$$\theta^+ = \theta - \eta \nabla_{\theta} \mathcal{L}(\theta)$$

מוטיבציה: הגרדיאנט של פונקציה הוא כיוון העליה הכי מהיר שלה, לכן לחסוך מהמיקום הנוכחי כפולה שלילית של הגרדיאנט תוריד אותנו אל עבר מינימום מקומי. לפיכך הוספה של כפולה שלילית של הגרדיאנט של פונקציית ההפסד θ אמורה להוביל לאורך זמן לירידה בערך MSE.

חשוב! עד שbow באלגוריתם מورد הגרדיאנט הוא בחירה נconaה של גודל הצעד: גודל צעד גדול מדי יגרום לכך שבхиים לא נתכנס, וגודל צעד קטן מדי יגרום לכך שנתקנסן מקצת מאוד איטי (לציריך).

בדרכ' בוחרים $\eta = 0.05$, ואם צריך משנים את הערך במהלך האימון.

הגדירה 6.1. **היפרפרמטר** הוא פרמטר שקובעים על מנת להגדיר חלקו במודל.

במקרה שלנו קצב הלמידה הוא היפר פרמטר, כי הוא קבוע כמו מהר מורד מגרדיינט מתכנס וכאן הוא קבוע האם יהיה למודל איטוי או מהיר.

הערה 7.1. בעיה מרכזית של מורד המדריאנט היא שאנו יכולים ליתקע במינימום מקומי (לציר!) לא משנה כמה פעמים נרים מורד המדריאנט, שכן ישבו אנשים חכמים והמציאו אלגוריתמים אחרים הידועים כטור optimizers (נרחיב עליהם בהמשך הקורב אولي), כגון *SGD*.

2. *Momentum* שזה ככל שאנו מתקדים בצריך אחד יותר זמן, נתקדם בו יותר מהר.

3. *Nesterov* שזה כמו *Momentum* אבל קצת יותר חכם.

4. *AdaGrad* שמוריד את הה-*learning rate* ככל שימושים לאמן.

5. *RMS* שמוריד גם הוא את ה-*lr* אבל מונע מה לקטון יותר מדי.

6. *Adam* שמשלב *RMS* ו-*Momentum*.

7. ולאחרונה גם *Moun*.

מסקנה: ככל יותר התהום עדין פועל ואין באמת שיטה הכח טובה למצוא את הפרמטרים הכי טובים אם אין פתרון סגור.

שאלה 8.1. מתי עוזרים את האימון?

פתרון: יש כמה אפשרויות:

1. כאשר $\epsilon < \|\nabla \mathcal{L}\|$.

2. נגמר התקיכוב נניח מסיבה מסוימת יש לו רק שעה על המחשב להריץ את האלגוריתם מה שיצא אליו מרווח.

3. אם חשוב להגיע לערך כלשהו של המודל (פרטיו עוד רגע) ומצילוחים, עוזרים.

4. כוששיות מספר עצדים שהגדורי מראש.

2 מטריקות הערכת המודל.

עכשו אחרי שאימנו את המודל שלנו, איך אנחנו יודעים בפועל כמה טוב הוא? רק בגלל שדחפנו את הדעתה שלנו למודל ואמרנו תלמד לא אומר שהוא למד. אנחנו צריכים מدد כלשהו שיאמר לנו כמה טוב המודל בכלל מליל את הלמידה שלו, לומר לנו לפחות דרך כלשהו להעריך את המודל שלנו.

הרעינו הוא מאד פשוט: במקומות לדוחף את כל הדעתה שלנו למודל נחלק אותה לשתי קבוצות זרות: **קבוצת אימון (train set)** ו**קבוצת מבחון (test set)**. כך נוכל לבדוק את השגיאה של המודל על דעתה שהוא עוד לא ראה ושאנו יודעים מה הערך האמיתי שלו.

נזכיר כי

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \quad (1)$$

מודד את המרחק בריבוע ממוצע של החיזויים (פרדיקציות) שלנו מהערכים האמיתיים, אך אם נקטין את ה-*MSE* שלנו נקבל מודל טוב יותר. הבעיה היחידה היא שהיחסות מידת *MSE* הן היחידות מידת *RMSE*。

$$RMSE = \sqrt{MSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2} \quad (2)$$

כיוון שה-*RMSE* באותו יחידות מידת כמו של *u*, יהיה יותר קל לפרש אותו ולהחליט לפיו האם המודל טוב או לא טוב. עוד אפשרות היא להגיד

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|^2 \quad (3)$$

היתרון של *MAE* זה שהוא גיש לאוטליירים (נקודות שונות משאר הדעתה) ולפיכך הוא בחירה טובה אם לדעתה של *i* יש לעיתים ערכים קיצוניים.

עוד אפשרות נקראת *R²* והוא הדרך הכח שכייה לסכום כמה טוב המודל שלו:

$$R^2 = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2} \quad (4)$$

כאשר $\bar{y} = \frac{1}{N} \sum_i y_i$ זה הממוצע של y .

האינטואיציה היא שהמודנה מודד עד כמה רוחקות התצפיות מהאמות, בעוד שהמכנה מודד עד כמה הנתונים רוחקים מניבוי הממוצע בלבד. אם המודל מושלם אז $R^2 = 1$, אם המודל מוצע (תרתי משמע חוצה כל הזמן את הממוצע) אז $R^2 = 0$, ואם $R^2 < 0$ אז המודל ממש גרוע. דרך נוספת למדוד טובות המודל שלנו היא על ידי לשרטט את y מול \hat{y} : ככל שהמודל יותר טוב נצפה מהזוג (y_i, \hat{y}_i) להיות קרוב יותר לישר $x = y$. לזה קוראים scatter plot.

ניתן גם לשרטט את השארית $Residual_i = y_i - \hat{y}_i$ וככל שהמודל יותר טוב נצפה שהשאריות יהיו קרובות לאפס. הערה 1.2. חולקה טובות של הדאטה ל- $train \setminus test$ היא $80\% \setminus 20\%$ בהתאם, ונחיב על כך בהמשך.

תרגיל 2.2. פתרו את כלל האימון של אלגוריתם מורד הגרדיינט לרגרסיה ליניארית.

פתרון. נזכיר כי כלל העדכון הוא $\theta = \theta - \eta \nabla \mathcal{L}(\theta)$. אם נציב b נקבל

$$\mathcal{L}(w, b) = \frac{1}{N} \sum_{i=1}^N (y_i - (w^T x_i + b))^2$$

נגזר ביחס לכל אחד מהפרמטרים:

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial}{\partial w} \frac{1}{N} \sum_{i=1}^N (y_i - (w^T x_i + b))^2 = \frac{2}{N} \sum_{i=1}^N -x_i(y_i - (w^T x_i + b)) = \frac{2}{N} \sum_{i=1}^N -x_i(y_i - \hat{y}_i) = \frac{2}{N} \sum_{i=1}^N x_i(\hat{y}_i - y_i)$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial}{\partial b} \frac{1}{N} \sum_{i=1}^N (y_i - (w^T x_i + b))^2 = -\frac{2}{N} \sum_{i=1}^N (y_i - (w^T x_i + b)) = -\frac{2}{N} \sum_{i=1}^N (y_i - \hat{y}_i) = \frac{2}{N} \sum_{i=1}^N (\hat{y}_i - y_i)$$

$$\sum_{i=1}^N x_i(\hat{y}_i - y_i) = X^T(\hat{y} - y) \text{ וכאן לפי כפל מטריצות } \hat{y} = \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_N \end{pmatrix} \in \mathbb{R}^N, y = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} \in \mathbb{R}^N \text{ וכן } X = \begin{pmatrix} -x_1^T - \\ \vdots \\ -x_N^T - \end{pmatrix} \in \mathbb{R}^{N \times d}$$

כלומר

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{2}{N} X^T(\hat{y} - y)$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{2}{N} \sum_{i=1}^N (\hat{y}_i - y_i)$$

ולכן נקבל

$$w^+ = w - \eta \frac{2}{N} X^T(\hat{y} - y)$$

$$b^+ = b - \eta \frac{2}{N} \sum_{i=1}^N (\hat{y}_i - y_i)$$

הערה 3.2. אם נסמן במקומות $X = \begin{pmatrix} 1 & - & x_1^T & - \\ 1 & - & x_2^T & - \\ \vdots & & \vdots & \\ 1 & - & x_N^T & - \end{pmatrix}, y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}, \beta = \begin{pmatrix} b \\ w_1 \\ \vdots \\ w_d \end{pmatrix}$ נוכל לכתוב

$$\beta^+ = \beta - \eta \frac{2}{N} X^T(\hat{y} - y)$$