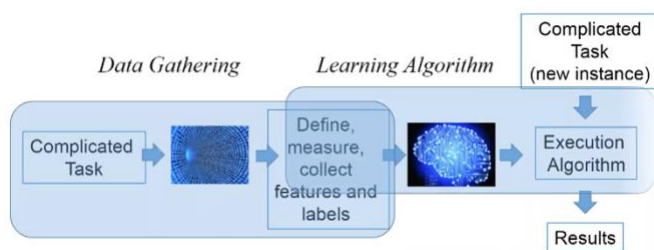


מסגרת משימות למידה: הרצאה 1 – Linear Regression



מושגים בלמידה חישובית ממידע:

- אוסף של דאטה, דגימות, טובות ורעות (חיוביות ושליליות) שיקראו samples או instances
- איסוף מידע אודות הדגימות הנ"ל, שיקרא features תכונות של הדגימות (חלקן רלוונטיות וחלקן אינן, לא בהכרח נדע מראש מה רלוונטי עבורנו ומה לא רלוונטי עבורנו לדעת)
- הפעלת אלגוריתם הלמידה שייצור אלגוריתם מבצע = שידע לקבל דגימות חדשות ולחשב עבורנו

סוגי למידה:

- רגרסיה: בהינתן דגימה $\{x_i, y_i\}$ נמצא פונקציה f כך ש $y=f(x)$
- קלסיפיקציה: בהינתן דגימה $\{x_i, y_i\}$ כאשר y_i הינו 0 או 1 בדאטה אימון, נקבע, עבור כל x חדש אם x שייך ל- c_0 או ל- c_1 .
- Density Estimation
- Clustering

רגרסיה ליניארית

דוגמה: בהינתן גודל של בית נרצה לקבוע מחיר עבורו. להלן דאטה אימון שלנו שמכיל 10 דגימות ו feature יחיד. מתוכו נרצה ללמוד פונקציה מכלילה. ההנחה היא שאנחנו רוצים מודל לינארי – ולכן נעשה רגרסיה לינארית – הדבר מגביל את מרחב ההיפותזות שלנו למרחב הפונקציות הלינאריות בלבד.

כיצד מייצגים את ההיפותזה h במרחב הפונקציות הלינאריות? על ידי משוואה לינארית:

$$y = \theta_0 + \theta_1 x$$

בדוגמה שלנו יש פיז'ר בודד, מה עושים כאשר יש מספר פיז'רים?

כאשר $X_i =$ הוא וקטור הפיז'רים שלנו, עבור כל דגימה i .

$$\mathbf{x}^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)})$$

למען הנוחות, נוסיף פיז'ר "קבוע" עבור כל i , $X_0(i) = 1$

עבור θ פיז'רים, ההיפותזה הלינארית שלנו תהיה מיוצגת על ידי וקטור הפרמטרים טטה:

$$\theta = (\theta_0, \theta_1, \theta_2, \dots, \theta_n)$$

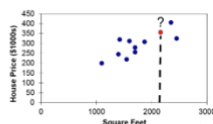
• We say that $\theta = (\theta_0, \theta_1, \theta_2, \dots, \theta_n)$

is the vector of parameters that defines our function (or our model)

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

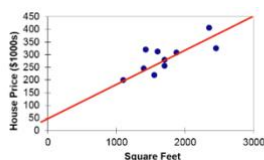
כך שמתקיים:

Scatter plot of House Price (y) vs House Size (x)



Prediction:

Given house of size x , what would be its price $y = f(x)$?



- We assume/hypothesize that the relationship between the observed and the independent/explained variable is linear and thereby conduct our search
- This is our **Hypotheses Space** – all linear functions

וקטור הפרמטרים טטה למעשה מגדיר את הפונקציה שלנו (המודל שלנו), וכך למעשה בממד גבוה תיראה הפונקציה שנרצה ללמוד, לכן עלינו ללמוד את הטעות. נרצה למצוא את הטעות הטובה ביותר – שניב את המודל הטוב ביותר.

המכפלה הפנימית של טטה באיקס החדש (i) תניב לנו את מחיר הבית.

הדאטה אימון training data תעזור לנו להגיע לטטה הטובה ביותר. בעולם "מושלם" היינו רוצים למצוא טטה שמניבה טעות 0, כלומר ממש למצוא מודל שעובר בכל הנקודות, אלגוריתם למידה כזה נקרא consistent learner. זהו לא המקרה בעולם האמיתי וכן לא בדוגמה שלנו. אבל עדיין נרצה את הטעות הקטנה ביותר. אז איך נמדוד את השגיאה שלנו? הטעות בכל דגימה (instance) הינה:

$$(\theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \dots + \theta_n x_n^{(i)} - y^{(i)}) = \theta \cdot \mathbf{x}^{(i)} - y^{(i)}$$

כעת נבצע ממוצע על כל דגימות האימון במטרה לקבל את פונקציית העלות שלנו Mean square error:

cost function:

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

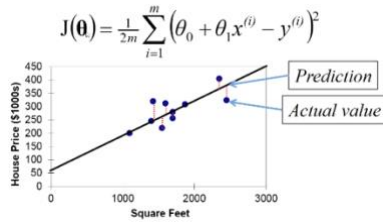
אנו משתמשים בשגיאה הריבועית כדי שטעות בכיוונים שונים לא תתבטל, וגם זו פונקציה smoother מערך מוחלט (שאינה גזירה ב0).

כך נוכל לחשב את הטעות עבור כל טטה וכך להעריך את איכות הטטה, על ידי חישוב השגיאה. ומכיוון שנרצה להגיע לטטה בעלת השגיאה המינימלית, נרצה למצוא את הטטה שמניבה ערך מינימלי של הפונקציה J – פונקציית העלות. הטעות הן המודל.

- Hence, our best hypothesis θ^* would be the one that minimizes the cost function:

$$\theta^* = \arg \min_{\theta} [J(\theta)] = \arg \min_{\theta} \left[\frac{1}{2m} \sum_{i=1}^m (\theta \cdot \mathbf{x}^{(i)} - y^{(i)})^2 \right]$$

- How can we find it?



• The directional derivative in the direction of the vector $u = (u_1, u_2)$ (a scalar!) can also be written as:

$$D_u f(x_1, x_2) = \nabla f \cdot \vec{u} = |\nabla f| |\vec{u}| \cos \beta = |\nabla f| \cos \beta$$

• Where β is the angle between u and ∇f

• However, $\cos \beta \leq 1$.

• Therefore:

1. The greatest increase in the function happens in the direction of the gradient (i.e. $\beta=0$)
2. The greatest decrease is in the direction $-\nabla f$ (i.e. $\beta=180^\circ$)

✖

השאלה, איך מביאים למינימום את פונקציית המחיר / העלות, שהיא בעלת מספר משתנים ?

- הדרך הראשונה למצוא מינימום של פונקציה היא לגזור אותה ולהשוות אותה ל-0.

- הדרך השנייה: גרדיאנט דיסנט.

שתי השיטות יכולות להוביל אותנו למינימום לוקאלי, ולא גלובלי כפי שנרצה!

נזכיר שנגזרת חלקית מחושבת על ידי גזירה על פי משתנה אחד כך שכל השאר נותרים קבועים.

הנגזרת הכיוונית בכיוון u מוגדרת:

$$D_u f(x_1, x_2) = \lim_{s \rightarrow 0} \frac{f(x_1 + su_1, x_2 + su_2) - f(x_1, x_2)}{s} = \left(\frac{df}{ds} \right)_u$$

Define the
GRADIENT of f :

$$\nabla f = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2} \right)$$

הגרדיאנט של פונקציה f מוגדר להיות:

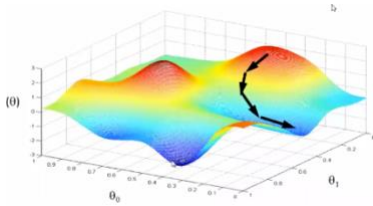
$$D_u f(x) = \nabla f(x) \cdot u$$

משפט: עבור כל כיוון $u = (u_1, u_2)$ וכל נקודה $x = (x_1, x_2)$ מתקיים:

כלומר הנגזרת הכיוונית בכיוון u בנקודה x שווה למכפלה הפנימית של הגרדיאנט של הפונקציה f בנקודה x עם וקטור הכיוון u .

במטרה להגיע למינימום של פונקציה נתחיל בנקודה מסוימת ונלך נגד הגרדיאנט. תהליך זה נקרא **גרדיאנט דיסנט**.

גרדיאנט דיסנט לא מבטיח מינימום לוקאלי! ומתאים גם לאלגוריתמים שאינם מניבים פונקציה שהיא בהכרח לינארית.



אלגוריתם גרדיאנט דיסנט Gradient Descent

• נתחיל עם ערך רנדומלי כלשהו של טטה $\theta = (\theta_0, \theta_1, \theta_2, \dots, \theta_n)$

• עד שנגיע למינימום נבצע את הצעד הבא: (הליכה נגד הגרדיאנט)

For all j ,

$$\text{Update } \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

עבור כל j , נעדכן את טטה j להיות:

כלומר נעדכן את טטה j להיות טטה j עם "צעד קטן" נגד הגרדיאנט.

• אלפא היא פרמטר של האלגוריתם שנקרא learning rate גודל צעד העידכון שמוגדר.

כאשר חישוב הגרדיאנט של פונקציית העלות הינו:

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1, \dots, \theta_n) = \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{2m} \sum_{i=1}^m \frac{\partial}{\partial \theta_j} (\theta_0 + \theta_1 x_1^{(i)} + \dots + \theta_j x_j^{(i)} + \dots + \theta_n x_n^{(i)} - y^{(i)})^2$$

$$= \frac{1}{2m} \sum_{i=1}^m 2(\theta_0 + \theta_1 x_1^{(i)} + \dots + \theta_j x_j^{(i)} + \dots + \theta_n x_n^{(i)} - y^{(i)}) \cdot x_j^{(i)}$$

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

האלגוריתם הסופי להפעלת גרדיאנט דיסנט על פונקציית העלות (השיגאה של טטה):

- מתחילים מערך התחלתי של טטה

- חוזרים על הצעד עד תנאי העצירה: מעדכנים טטה j חדש להיות טטה j

הנוכחי פחות learning rate (אלפא) כפול הנגזרת החלקית של פונקציית

העלות ביחס לטטה j , מחושב על הטטה הנוכחי.

- נזכור ש- $x_0(i) = 1$

• Start with some value for $\theta = (\theta_0, \theta_1, \theta_2, \dots, \theta_n)$

• Repeat until you reach a minimum: ✖

For all j ,

$$\text{Update } \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

• $\alpha > 0$ is a parameter of the algorithm called the **learning rate**

• Updates are simultaneous (in all $n + 1$ directions)

• In the general case this process can still be trapped in local minima!

• Initialize $\theta = (\theta_0, \theta_1, \theta_2, \dots, \theta_n)$

• Repeat until you reach a minimum (or stop cdn):

○ For all $0 \leq j \leq n$,

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

• In words: set the new θ_j to the current θ_j minus the learning rate (α) times the partial derivative of the error function with respect to θ_j , computed at the current θ .

• Also remember that $x_0^{(i)} = 1$

$$X \cdot \theta = y$$

$$\begin{pmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_n^{(1)} \\ x_0^{(2)} & x_1^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & & \vdots \\ x_0^{(m)} & x_1^{(m)} & \dots & x_n^{(m)} \end{pmatrix} \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}$$

- If X is square and non singular we can write $\vec{\theta} = X^{-1}y$
- Most of the time X is overdetermined ($m \gg n$).

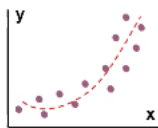
- $X^T X$ can still be singular (or not full rank) and not have an inverse. This can be resolved with some more algebra.
- This pinv technique doesn't work for all error functions J . Gradient descent is more general.
- Gradient descent allows parallelization.

Polynomial Regression

- We can expand our feature space by using functions of the original features.
- For example, if we want to use a cubic function feature space we can define:

$$x_0 = 1, x_1 = x, x_2 = x^2, x_3 = x^3$$

then use regular regression and in essence we are learning the function

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$


אם המטריצה X הייתה ריבועית, היינו משתמשים באינברס שלה.

סודו אינברס:

- Setting the gradient to zero yields the necessary condition for minimum (see notes):

$$X^T X \cdot \vec{\theta} = X^T \vec{Y}$$

- Now, $X^T X$ is square and often nonsingular and so we can solve for $\vec{\theta}$ uniquely as:

$$\vec{\theta} = \text{pinv}(X) \vec{Y} \quad \text{where} \quad \text{pinv}(X) = (X^T X)^{-1} X^T$$

- The $n \times m$ matrix $\text{pinv}(X) = (X^T X)^{-1} X^T$ is called the **pseudo inverse** of X (which is $m \times n$)
- If X is square it is just its inverse.

הערות לגבי סודו אינברס:

- יתכן מצב בו למכפלה לא יהיה אינברס, ניתן לפתור זאת עם פעולות אלגבריות אשר לא נתעמק בהן בקורס.
- הפינב לא יעבוד לכל פונקציית עלות j , גרדיאנט דסנט הוא כללי יותר.
- גרדיאנט דסנט מאפשר פרלליזציה.

נשתמש בטכניקה של רגרסיה גם עבור פונקציות פולינומאליות: