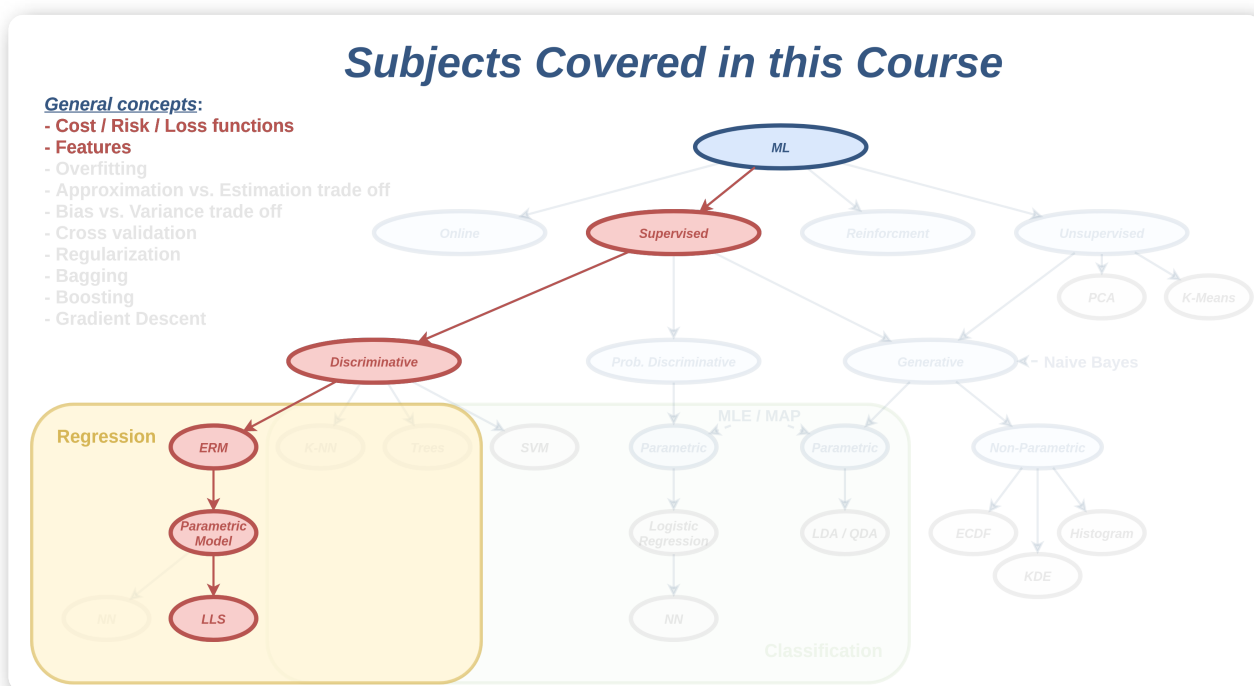


# הרצאה 2 - רגרסיה לינארית

Slides PDF Code

## מה נלמד היום



## Supervised learning (למידה מונחית)

בהרצאה הקודמת הכרנו באופן כללי את סוגי הבעיות בהם עוסק התחום של מערכות לומדות והתחלנו לדבר על בעיות מסוג supervised learning (למידה מונחית). בעיות אלו הם הבסיסיות ביותר בתחום והבנה טובה של בעיות אלו היא הבסיס להבנה של כל שאר הבעיות במערכות לומדות. אנו נעסוק בבעיות מסוג זה לאורך רוב רובו של הסימסטר.

על מנת להבין ולהגדיר מה הם בעיות supervised learning עלינו ראשית לחזור על הנושא של בעיות חיזוי.

### בעיית החיזוי

בתרגול 2 מופיעה תזכורת מורחבת יותר של התחום, אך לצורך הרצאה זו נסתפק בתיאור הקצר המופיע פה. בבעיית החיזוי אנו מנסים לחזות את ערכו של משתנה אקראי לא ידוע, לרוב על סמך משתנים אקראיים ידועים. בעיות חיזוי הם **מאד** נפוצות ומופיעות במגוון רחב של תחומים בהנדסה ומדע. בהנדסת חשמל בעיות אלו מופיעות בתחומים כגון עיבוד אותות, תקשורת ספרתית ובקרה. יתרה מזאת, בעיות חיזוי מלוות אותנו כמעט בכל פעולה יום יומית, לדוגמא:

כאשר אנחנו מתלבטים האם לקחת איתנו מטריה ביציאה מהבית, או למעשה מנסים לבצע חיזוי של האם ירד גשם או לא על סמך פרטי מידע שיש בידינו, כגון התחזית ששמענו, כמות העננים בשמים, צבע העננים וכו'.

ביום יום אנחנו אולי לא מנסים לפתור את בעיות החיזוי באופן מתמטי, אך אנו כן מחזיקים בראש איזה שהוא מודל של הקשרים הסטטיסטיים בין המשתנים השונים, ואנו מנסים לבצע את החיזוי על סמך אותו מודל באופן איכותי.

## הקשר ל supervised learning

בבעיות חיזוי קלאסיות, אנו מניחים שהפילוג של כל המשתנים האקראיים ידוע, וכי האתגר הוא מציאת החזאי האופטימאלי של סמך הפילוג. לעומת זאת, בבעיות supervised learning (ובבעיות במערכות לומדות באופן כללי) אנו מניחים כי הפילוג אינו ידוע ובמקומו נתון לנו מדגם של דגימות מתוך אותו פילוג. את החזאי נאלץ כעת לבנות על סמך המדגם במקום על סמך הפילוג. במהלך הקורס אנו נדון בכיר שיטות שונות לבנות חזאים באופן זה ונדון בבעיות הקיימות בשיטות אלו.

## סימונים ושמות

בקורס זה אנו נשתמש בסימונים והשמות הבאים:

- **Labels** (תיות / תגיות):  $y$  (או במקרה הוקטורי  $\mathbf{y}$ ) - יהיה המשתנה / הוקטור האקראי שאותו אנו מנסים לחזות. בקורס זה ה  $y$ , labels, יהיו כמעט תמיד סקלריים.
- **Observations \ measurements** (תצפיות או מדידות):  $x$  (או במקרה הוקטורי  $\mathbf{x}$ ) - יהיה הוקטור האקראי אשר מכיל את המשתנים שלפיהם נרצה לבצע את החיזוי. במקרים מסויימים החיזוי יהיה על פי משתנה יחיד ואז  $x$  יהיה סקלר.
- $\hat{y}$  - תוצאת חיזוי כל שהיא.
- מרחב החזאים \ השערות  $\mathcal{H}$ . במרחב זה נמצאים כל החזאים האפשריים.
- $\hat{y} = h(\mathbf{x})$  - פונקציית החיזוי. נשים לב כי  $h \in \mathcal{H}$ .
- אנו נשתדל להשתמש ב  $D$  לסימון האורך של הוקטור  $\mathbf{x}$

## The dataset (המדגם)

כפי שצינו, את הבניה של החזאי אנו נעשה על פי מדגם מתוך הפילוג הלא ידוע. לרוב המדגם יהיה מורכב מזוגות של  $\mathbf{x}$  ו  $y$  אשר יוצרו מתוך  $N$  דגימות בלתי תלויות:

$$D = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$$

נשתדל להשתמש תמיד ב  $N$  לסימון מספר הדגימות שבמדגם.

## הנחת ה i.i.d

במערכות לומדות אנו תמיד נניח כי הדגימות במדגם נוצרו כולם מאותו הפילוג באופן בלתי תלוי אחת בשניה. זאת אומרת שזוג המשתנים  $\{\mathbf{x}^{(i)}, y^{(i)}\}$  הינו בלתי תלוי סטטיסטית בזוג המשתנים  $\{\mathbf{x}^{(j)}, y^{(j)}\}$  כאשר  $i \neq j$ .

## מהו החזאי האופטימאלי

באופן כללי, כל פונקציה אשר ממפה מהמרחב של  $\mathbf{x}$  למרחב של  $y$  היא פונקציית חיזוי חוקית. נשאלת אם כן השאלה מהי פונקציית החיזוי המוצלחת ביותר? באופן כללי היינו מעוניינים למצוא חזאי אשר לעולם לא טועה. בפועל, מכיוון שע  $y$  הינו משתנה אקראי לא נוכל אף פעם לחזותו במדוייק (מלבד במקרים מיוחדים בהם  $y$  נקבע באופן חד ערכי על ידי  $\mathbf{x}$ ).

מכיוון שבעבור כל חיזוי שנבחר אנו מצפים לשגיאה כל שהיא, אנו צריכים להגדיר דרך להשוות בין הטעויות שאותם מבצעים החזאים שונים. אנו צריכים להחליט לדוגמא איך נבחר בין חזאי שעושה כל הזמן שגיאות בינוניות לבין חזאי אשר רוב הזמן עושה שגיאות ממש קטנות אך פעם בכמה זמן עושה שגיאה מאד גדולה. הנושא הראשון שנעסוק בו בהרצאה יהיה הדרך שבה נרצה להשוות בין הביצועים של חזאים שונים.

## Regression vs. Classification

מוקבל לחלק את הבעיות ב supervised learning לשני תתי תחומים:

- **בעיות regression (רגרסיה)** - בעיות בהם  $y$  הוא משתנה רציף.
- **בעיות classification (סיווג)** - בעיות בהם  $y$  הוא משתנה בדיד אשר יכול לקבל ערכים מתוך סט ערכים סופי (ולרוב קטן).

דוגמאות:

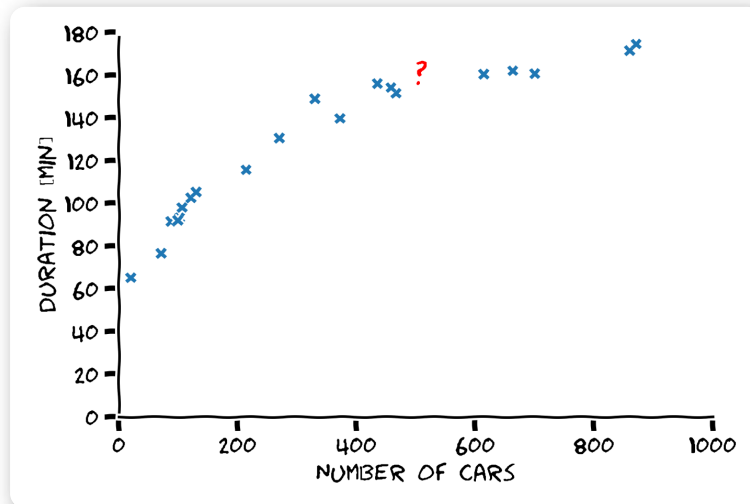
- **רגרסיה:** חיזוי זמן נסיעה בכביש החוף, חיזוי מרחקים לאובייקטים בתמונה, חיזוי מחירים של דירות וכו'.
- **סיווג:** חיזוי של המחלה בה חולה אדם מסוים על פי הסימפטומים שלו, חיזוי של האם דואר מסוים הוא spam או לא, חיזוי של האם עסקת אשראי מסוימים היא לגיטימית או הונאה וכו'.

(בעיקרון יכולים להיות גם בעיות בהם  $y$  בדיד ולא סופי. בבעיות מסוג זה לרוב פשוט מניחים ש $y$  רציף והופכים את הבעיה לבעיית רגרסיה ולבסוף מעגלים את התוצאה).

כפי שנראה בהמשך הקורס, אבחנה זו חשובה מכיוון שהאופי של  $y$  ישפיע על הדרך שבה ננסה לפתור את הבעיה.

## בעיית רגרסיה לדוגמא

נסתכל על בעיית החיזוי של זמן הנסיעה בכביש החוף על סמך מספר המכוניות בכביש:



ננסח את הבעיה בצורה פורמאלית:

- $y$  - Labels - המשתנה האקראי של זמן הנסיעה.
- $x$  - Measurements - המשתנה האקראי של מספר המכוניות על הכביש.
- $h$  - פונקציית החיזוי אשר מקבלת את מספר המכוניות על הכביש ומוציאה חיזוי של זמן הנסיעה.
- $\mathcal{D}$  - מדגם הנתון של הזוגות של (מספר מכוניות, זמן נסיעה).

המטרה שלנו הינה להשתמש ב  $d$  על מנת למצוא חזאי  $\hat{y} = h(x)$  אשר יהיה כמה שיותר מוצלח תחת קריטריון שאותו נצטרך להגדיר.

## השלבים הכלליים לפתרון הבעיה

הרעיון מאחורי כמעט כל השיטות במערכות לומדות הוא זהה:

1. נגדיר קריטריון מתמטי אשר מודד עד כמה מודל מסוים מצליח לבצע את המשימה
2. נבחר משפחה רחבה של מודלים בתקווה שלפחות אחד מהם יהיה מוצלח מספיק.
3. נחפש מבין כל המודלים במשפחה את המודל המוצלח ביותר.

## The cost function (פונקציית המחיר)

נתחיל מהשלב הראשון של הגדרת הקריטריון של שלפיו נרצה לבחור את החזאי שלנו. הדרך המקובלת לעשות זאת הינה על ידי הגדרת פונקציית מחיר.

פונקציית המחיר  $C(h)$  היא פונקציה אשר מעניקה לכל חזאי ציון. לרוב נהוג להגדיר את הפונקציה כך שכל שהציון נמוך יותר החזאי טוב יותר. אנו נגדיר את החזאי האופטימאלי כחזאי בעל הציון הנמוך ביותר מבין כל החזאים האפשריים. נסמן את החזאי האופטימאלי ב  $h^*$ :

$$h^* = \arg \min_{h \in \mathcal{H}} C(h)$$

ישנם דרכים רבות להגדיר את פונקציית המחיר, ואין דרך "נכונה" לעשות זאת. באופן כללי, הבחירה של פונקציית המחיר צריכה להתאים לבעיה שאותה רוצים לפתור בעזרת החזאי. באופן כללי פונקציית המחיר אמורה לשקף את המחיר אותו "נשלם" על שימוש בחזאי נתון כל שהוא, בפועל, במרבית המקרים משתמשים באחת מכמה פונקציות מחיר נפוצות.

נציג את אחת הדרכים הפופולריות להגדיר פונקציית מחיר, אשר עושה זאת על ידי שימוש בפונקציה הנקראת פונקציית הפסד (loss function).

## Risk and loss functions (פונקציות סיכון והפסד)

פונקציית הסיכון הינה מקרה פרטי של פונקציית המחיר והיא מוגדרת באופן הבא:

בעוד שפונקציית המחיר מנסה לתת ציון ליכולת החיזוי הכללית של החזאי, פונקציית הפסד  $l$  נותנת ציון לחיזוי בודד מסויים. שבהינתן חזאי  $h$  ודגימה אקראית עם  $\mathbf{x}$  ו  $y$  מסויים, ההפסד המשווין לדגימה זו יהיה:

$$l(h(\mathbf{x}), y) = l(\hat{y}, y)$$

את פונקציית המחיר הכוללת ניתן כעת להגדיר כתוחלת של פונקציית ה loss על פני הפילוג של  $\mathbf{x}$  ו  $y$

$$C(h) = \mathbb{E}[l(h(\mathbf{x}), y)]$$

במקרים כאלה, מוקבל לכנות את פונקציית המחיר, פונקציית **risk** (סיכון), ולסמנה באות  $R$ :

$$R(h) = \mathbb{E}[l(h(\mathbf{x}), y)]$$

### פונקציות loss נפוצות

- פונקציית loss נפוצה לבעיות classification היא פונקציית ה **Zero-One loss** אשר מוגדרת באופן הבא:

$$l(\hat{y}, y) = I\{\hat{y} \neq y\}$$

לפונקציית ה risk אשר משתמשת ב loss הזה קוראים: **misclassification rate**.

- פונקציית loss נפוצה לבעיות regression היא פונקציית ה **loss**  $l_2$  אשר מוגדרת באופן הבא:

$$l(\hat{y}, y) = (\hat{y} - y)^2$$

לפונקציית ה risk אשר משתמשת ב loss הזה קוראים: **(mean squared errors (MSE)**.

- במקרים רבים נהוג להשתמש דווקא בשורש של ה MSE בכדי שפונקציית המחיר תחזיר ערכים באותם יחידות כמו  $y$ . במקרה זה קוראים לפונקציית המחיר **(root mean squared errors (RMSE)**. התוספת של השורש לא משנה את בעיית האופטימיזציה (משום שהיא פונקציה מונוטונית עולה) ולכן הוא לא משפיע על החזאי המתקבל.

- פונקציית loss נפוצה נוספת לבעיות regression היא פונקציית ה **loss**  $l_1$  אשר מוגדרת באופן הבא:

$$l(\hat{y}, y) = |\hat{y} - y|$$

לפונקציית ה risk אשר משתמשת ב loss הזה קוראים: **(mean absolute errors (MAE)**.

## בעיה: הפילוג של המשתנים האקראיים לא ידוע

הבעיה עם האופן שבו הגדרנו את פונקציית הסיכון הינה העובדה שהיא מוגדרת על ידי תוחלת על פני הפילוג של המשתנים האקראיים בבעיה שהוא כאמור לא ידוע. בעיה זו קיימת לא רק בפונקציות מחיר מסוג סיכון אלא גם בסוגים שונים של פונקציות מחיר אשר כמעט תמיד תלויות בפילוג של המשתנים האקראיים בבעיה.

## שיערוך אמפירי של פונקציית המחיר / סיכון

בכדי לנסות ולהתמודד עם בעיה זו נוכל במקום לנסות ולחשב את ערכה של פונקציית המחיר באופן אנליטי, לנסות ולשערך את ערכה של פונקציית הסיכון על סמך אוסף של דוגמאות מתוך הפילוג (מדגם). שיעור על סמך דוגמאות מכונה שיערוך אמפירי.

## Empirical risk (סיכון אמפירי)

הסיכון האמפירי מוגדר על ידי החלפת התוחלת בפונקציית הסיכון בגרסא האמפירית שלה. אנו נשתמש בסימון  $\hat{\mathbb{E}}_{\mathcal{D}}$  על מנת לסמן את תחולת האמפירית המבוססת על המדגם נתון  $\mathcal{D}$ .

$$\mathbb{E}[f(\mathbf{x})] \approx \hat{\mathbb{E}}_{\mathcal{D}}[f(\mathbf{x})] = \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}^{(i)})$$

ניתן להראות כי כאשר מספר הדגימות  $N$  הולך לאין סוף התוחלת האמפירית מתכנסת לתוחלת האמיתית במובן הסתברותי. הסיכון האמפירי המקבל מהחלפה זו הינו:

$$R(h) = \mathbb{E}[l(h(\mathbf{x}), y)] \approx \hat{R}(h) = \frac{1}{N} \sum_{i=1}^N [l(h(\mathbf{x}^{(i)}), y^{(i)})]$$

### התאמת יתר

השימוש בגרסא האמפירית של פונקציית המחיר היא במקרים רבים בעייתית והיא גורמת בין היתר לתופעה המוכנה overfitting (התאמת יתר). בשלב זה אנו נתעלם מבעיה זו ואנו נעסוק בה בהרחבה בהרצאה הבאה.

## גישות לפתרון בעיות supervised learning

לפני שנעבור לשלב הבניה של החזאי אנו נציג שתי גישות שונות לבהם ניתן לגשת לבעיה. באופן כללי ניתן לחלק את השיטות לפתרון בעיות supervised learning לשתי הגישות הבאות:

**גישה גנרטיבית (generative)** - בגישה זו אנו ננסה להשתמש במדגם על מנת לנסות וללמוד את הפילוג הלא ידוע. שיטה זו נקראת גנרטיבית משום שהיא לומדת את הפילוג שמתוכו נוצרו (generated) הדגימות.

**גישה דיסקרימינטיבית (discriminative)** - בשיטה זו ננסה לבנות חזאי אופטימאלי על סמך הגרסא האמפירית של פונקציית המחיר בתקווה שהוא יקבל ציון טוב גם בגרסא הלא אמפירית של פונקציית המחיר (זאת אומרת, שהוא ידע להכיל בצורה טובה).

לכל אחד מהגישות יש את היתרונות והחסרונות שלה. במהלך הקורס נכיר אלגוריתמים משני הגישות ונעמוד על ההבדלים ביניהם.

## (Empirical risk minimization (ERM

בעבור מקרים בהם פונקציית המחיר מוגדרת כפונקציית סיכון, הגישה הדיסקרימינטיבית הבסיסית ביותר הינה לנסות לחפש באופן ישיר חזאי אשר ממזער את הסיכון האמפירי. הגישה הזו מכונה ERM (Empirical Risk Minimization) והיא מוגדרת על ידי בעיית האופטימיזציה הבאה:

$$h_{\mathcal{D}}^* = \arg \min_{h \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^N [l(h(\mathbf{x}^{(i)}), y^{(i)})]$$

שימו לב כי הוספנו את הכיתוב  $\mathcal{D}$  מתחת לחזאי האופטימאלי של בעיית האופטימיזציה זו. עשינו זאת משתי סיבות:

1. על מנת להדגיש את התלות של החזאי במדגם (לכל מדגם יהיה חזאי אופטימאלי אחר).
2. בכדי להבדיל את החזאי המתקבל משיטת ה ERM מהחזאי האופטימאלי של הבעיה המקורית אשר באופן עקרוני יהיה שונה.

## מודלים פרמטריים

לרוב אנו נרצה להגביל את החזאי שלנו למשפחה מצומצמת של פונקציות. לרוב אנו נרצה לעשות זאת על ידי בחירה של משפחה של פונקציות אשר מוגדרות על ידי מודל פרמטרי. ישנם שני סיבות עיקריות לכך:

1. כפי שנראה בהרצאה הבאה הגבלה זו חשובה בכדי לשפר את יכול ההכללה של החזאי ולסייע בהקטנת בעיית overfitting.
2. יותר פרקטי לנסות לחפש פרמטרים של מודל מאשר חיפוש כללי של פונקציה במרחב הפונקציות.

מודל פרמטרי מגדיר את המבנה הכללי של הפונקציות במשפחה עד כדי מספר סופי של פרמטרים אשר חופשיים להשתנות. את הפרמטרים של המודל נסמן בעזרת הוקטור  $\theta$ . אנו נשתמש ב  $h(\mathbf{x}; \theta)$  בכדי לתאר חזאי מהמשפחה הפרמטרית עם פרמטרים  $\theta$ . אין כל מגבלה על הצורה הכללית של המודל הפרמטרי, המודל הפרמטרי יכול להיות לדוגמא:

$$h(\mathbf{x}; \theta) = \frac{\theta_1^3 x_1 + x_4^{\theta_2}}{\log(\theta_3 x_2)}$$

דוגמאות נוספות למודלים פרמטריים:

1. פונקציות ליניאריות:  $h(\mathbf{x}; \theta) = \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$
2. פולינומים:  $h(\mathbf{x}; \theta) = \theta_1 + \theta_2 x_1 + \theta_3 x_1^2 + \theta_4 x_1^3$
3. טור פוריה סופי:  $h(\mathbf{x}; \theta) = \theta_1 \sin(\pi x) + \theta_2 \cos(\pi x) + \theta_3 \sin(2\pi x) + \theta_4 \cos(2\pi x)$
4. רשתות נוירונים

מודל פרמטרי למעשה ממפה כל פונקציה מהמשפחה הפרמטרית לוקטור. היתרון בעבודה עם וקטורים הינו שיש לנו סט עשיר של כלים בהם אנו יכולים להשתמש. לדוגמא, מכיוון שניתן לגזור לפי וקטורים, נוכל להשתמש ב gradient decent על מנת לחפש את המודל האופטימאלי במרחב המודלים. מיכיון שכל וקטור כעת מגדיר מודל מסויים (ולהיפך) ניתן לרשום את בעיית האופטימיזציה של מציאת המודל האופטימאלי כבעיית אופטימיזציה על וקטור הפרמטרים (במקום על  $h$ ):

$$\theta^* = \arg \min_{\theta} C(h(\cdot; \theta))$$

או במקרה של ERM:

$$\theta^* = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N [l(h(\mathbf{x}_i; \theta), y_i)]$$

## מודל ליניארי

המודל הפרמטרי הפשוט ביותר הינו המודל הליניארי. המודל הליניארי הוא בעל המבנה הבא::

$$h(\mathbf{x}; \theta) = \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_D x_D$$

דרך נוחה יותר לכתוב את המודל הזה היא בצורה וקטורית:

$$h(\mathbf{x}; \theta) = \mathbf{x}^\top \theta$$

## איבר היסט (bias)

ניתן להוסיף למודל גם איבר bias על מנת לקבל מודל מהצורה הבאה:

$$h(\mathbf{x}; \theta) = \theta_1 + \mathbf{x}^\top [\theta_2, \theta_3, \dots, \theta_{D+1}]^\top$$

לשם הנוחות בכדי לשמור על הכתיב הוקטורי של המודל נפריד לרוב את איבר ה bias משאר הפרמטרים. לרוב נסמן אותו בעזרת  $b$  או  $\theta_0$ :

$$h(\mathbf{x}; \theta, \theta_0) = \theta_0 + \mathbf{x}^\top \theta$$

אנו נראה מיד דרך נוחה יותר להוספת איבר ההיסט בעזרת שינוי של הוקטור  $\mathbf{x}$  כך שהביטוי  $\mathbf{x}^\top \theta$  יכיל גם את איבר ההיסט.

## Linear Least Squares (LLS)

מקרה מיוחד של בעיית ERM עם מודל ליניארי, הוא המקרה שבו משתמשים בפונקציית MSE (פונקציית risk עם loss ריבועי  $l_2$ ):

$$l(\hat{y}, y) = (\hat{y} - y)^2$$

השימוש במודלים ליניאריים וב MSE נפוץ מאד ולכן בעיית LLS מופיעה בתחומים רבים. בעיית האופטימיזציה המקבלת בעבור LLS הינה:

$$\theta_{\mathcal{D}}^* = \arg \min_{\theta} \frac{1}{N} \sum_{i=0}^N l(h(\mathbf{x}^{(i)}; \theta), y^{(i)}) = \arg \min_{\theta} \frac{1}{N} \sum_{i=0}^N (\mathbf{x}^{(i)\top} \theta - y^{(i)})^2$$

## כתיב מטריצי

את בעיה זו ניתן לרשום גם בצורה קומפקטית על ידי הגדרת הוקטור והמטריצה הבאים:

- נגדיר את וקטור התגיות  $\mathbf{y}$  כוקטור של כל התגיות במדגם:

$$\mathbf{y} = [y^{(1)}, y^{(2)}, \dots, y^{(n)}]^\top$$

- נגדיר את המטריצת המדידות  $X$  כמטריצה של כל ה  $\mathbf{x}$ -ים במדגם:

$$X = \begin{bmatrix} - & \mathbf{x}^{(1)} & - \\ - & \mathbf{x}^{(2)} & - \\ & \vdots & \\ - & \mathbf{x}^{(N)} & - \end{bmatrix}$$

בעזרת הגדרות אלו ניתן לרשום את בעיית האופטימיזציה של LLS באופן הבא:

$$\theta_{\mathcal{D}}^* = \arg \min_{\theta} \frac{1}{N} \|X\theta - \mathbf{y}\|_2^2$$

(ניתן להראות זאת על ידי רישום הורמה כסכום ושימוש בעובדה ש  $(X\theta = [\mathbf{x}^{(1)\top} \theta, \mathbf{x}^{(2)\top} \theta, \dots, \mathbf{x}^{(N)\top} \theta]^\top$ )

## פתרון סגור

מה שמויחד בבעיית האופטימיזציה של LLS הינה העובדה שניתן להגיע לפתרון סגור לפרמטרים האופטימליים על ידי גזירה והשוואה ל-0. הפתרון המתקבל הינו: (את החישוב עצמו אתם תראו בתרגול 3)

$$\begin{aligned} \nabla_{\theta} \left( \frac{1}{N} \|X\theta - \mathbf{y}\|_2^2 \right) &= 0 \\ \Rightarrow \theta &= (X^\top X)^{-1} X^\top \mathbf{y} \end{aligned}$$

פתרון זה נכון רק כאשר המטריצה  $X^\top X$  הפיכה. בתרגול אנו נדון במשמעות של תנאי זה.

## דוגמא

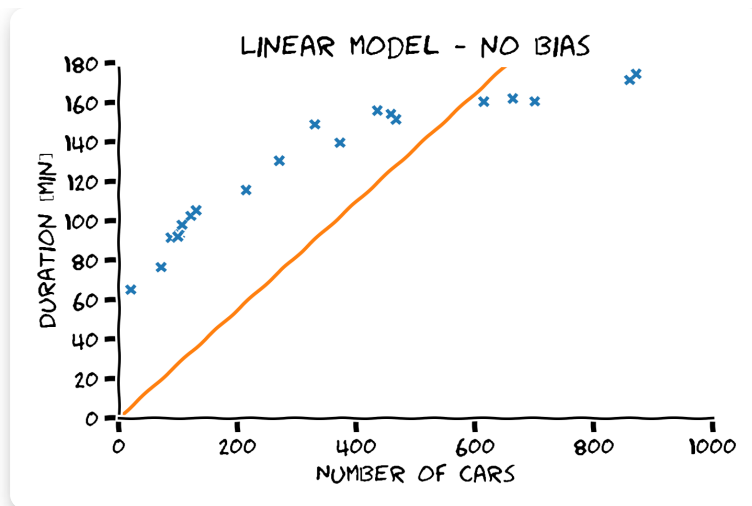
נשתמש במודל ליניארי ובפתרון של בעיית ה LLS בכדי למצוא חזאי לבעיית שיערוך זמן הנסיעה. נתחיל במודל ללא היסט:

$$h(x; \theta) = \theta x$$

את הפרמטר  $\theta$  האופטימאלי נוכל לחשב על ידי הצבה של הנקודות במדגם לתוך:

$$\theta_{\mathcal{D}}^* = (X^\top X)^{-1} X^\top \mathbf{y}$$

כאשר  $X = [x^{(1)}, x^{(2)}, \dots, x^{(N)}]^\top$ . התוצאה המקבלת הינה



## הוספת איבר היסט

בכדי להוסיף איבר היסט, עלינו להשתמש במודל מהצורה של:

$$h(x; \theta) = \theta_1 + \theta_2 x$$

הבעיה אם צורה זו הינה, שהפתרון הסגור שמצאו מתייחס למקרה שבו אין איבר היסט, לשם כך ננסה לנסח מחדש את הבעיה כך שיתקבל מודל ללא איבר היסט.

## עיבוד מקדים

נשים לב לעובדה שבבואנו לבצע משימת חיזוי או לא חייבים להשתמש בנתונים בצורת הגולמית ומותר לנו לבצע עיבוד מקדים של הנתונים לפני שאנו מזינים אותם לחזאי. נניח לדוגמא שאנו מעריכים שיהיה נוח יותר לבצע את החיזוי של זמן הנסיעה על פי הריבוע של כמות המכוניות על הכביש. במקרה שכזה נוכל פשוט להעלות המספר המכוניות בריבוע לפני שאנו מזינים אותם לחזאי. העיבוד המקדים יכול למעשה לפעול על כל וקטור המדידות  $\mathbf{x}$ , לעבד אותו ולייצר ממנו וקטור חדש. אנו נסמן ב  $\Phi$  את הפונקציה אשר מקבלת את המידע הגולמי  $\mathbf{x}$  מייצרת ממנו את  $\mathbf{x}_{\text{new}}$ :

$$\mathbf{x}_{\text{new}} = \Phi(\mathbf{x})$$

פעולת החיזוי במקרים אלו תהיה:

$$\hat{y} = h(\Phi(\mathbf{x}); \theta)$$

את קלט החדש  $\mathbf{x}_{\text{new}}$  מקובל לכנות וקטור **המאפיינים (features)**. השימוש במאפיינים מאפשר לנו מספרים דברים:

- הרחבת מודלים פשוטים, כגון המודל הליניארי למודלים מורכבים יותר (כפי שנראה כאן ובתרגול).
- שינוי האופן בו מיוצג המידע כך שיהיה לחזאי קל יותר לבצע את בעיית החיזוי. לדוגמא:
  - החלפת היחידות שבהם השתמשו לתיאור מדידה מסוימת
  - הפיכת תמונת פנים לוקטור של מאפיינים של פנים כגון: המרחק בין העיניים, גוון העור, עד כמה הפנים אליפטיות וכו'
  - ניקוי רעשים להקלטות audio.
- לקראת סוף הקורס נראה גם כיצד ניתן להשתמש במאפיינים על מנת להתמודד עם בעיית ה overfitting בעזרת שיטה המוכנה הורדת מימד.

במקרים רבים נרצה להתייחס לפונקציה אשר מייצרת איבר ספציפי ב  $\mathbf{x}_{\text{new}}$ , לשם כך נוח להתייחס לפונקציה  $\Phi$  כוקטור של פונקציות  $\varphi_i(\mathbf{x})$  אשר מייצרות כל אחת איבר אחד בוקטור  $\mathbf{x}_{\text{new}}$ :

$$x_{\text{new},i} = \varphi_i(\mathbf{x})$$

סימון מתמטי מקובל בו נשתמש הינו הסימון הבא:

$$\Phi(\mathbf{x}) = [\varphi_1(\mathbf{x}), \varphi_2(\mathbf{x}), \dots, \varphi_M(\mathbf{x})]^T$$



כאן  $\Phi$  מוצגת כוקטור של פונקציות, כאשר הפעלה וקטור שכזה על  $x$  מייצרת את הוקטור של הפלטים של הפונקציות  $\varphi_i$

## דוגמא: הוספה של איבר ההיסט בעזרת מאפיינים

על ידי שילוב של מודל לינארי עם מאפיינים נוכל לקבל חזאים מהצורה:

$$\hat{y} = h(x; \theta) = h_{\text{linear}}(\Phi(x); \theta) = \Phi(x)^\top \theta = \theta_1 \varphi_1(x) + \theta_2 \varphi_2(x) + \dots + \theta_M \varphi_M(x)$$

נחזור כעת לדוגמא של שיערוך זמן הנסיעה. נראה כעת כיצד ניתן להוסיף את איבר ההיסט על ידי שימוש במאפיינים. העבור הבחירה של במאפיינים הבאים:

$$\varphi_1(x) = 1, \quad \varphi_2(x) = x$$

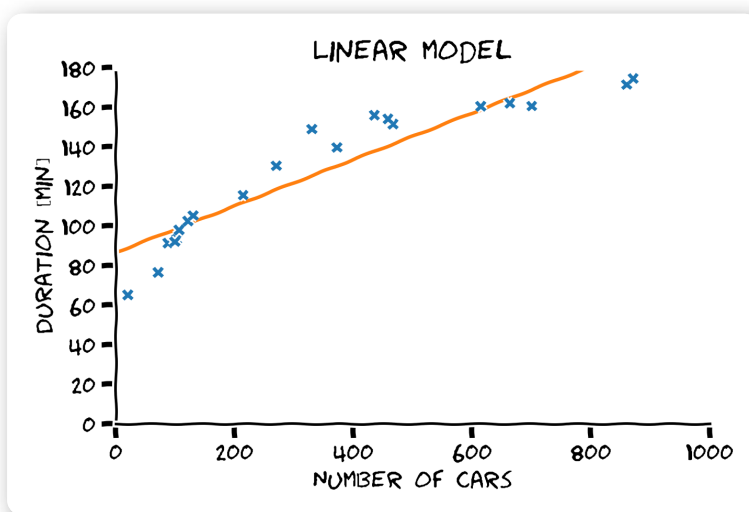
כל דגימה  $x$  תהפוך לוקטור  $x_{\text{new}} = [1, x]^\top$  ומודל החיזוי שלנו יהיה:

$$h(x; \theta) = \theta_1 + \theta_2 x$$

המטריצת המדידות  $X$  תהיה כעת:

$$X = \begin{bmatrix} 1 & x^{(1)} \\ 1 & x^{(2)} \\ \vdots & \vdots \\ 1 & x^{(N)} \end{bmatrix}$$

הצבה של מטריצה זו בנוסחה ל  $\theta_D^*$  נותנת את המודל הלינארי הבא:



באותו אופן ניתן כמובן להשתמש במודל הלינארי בכדי לייצג מגוון רחב של פונקציות כגון פולינומים או קומבינציה של גאוסיאנים, כפי שיודגם בתרגול 3.

