

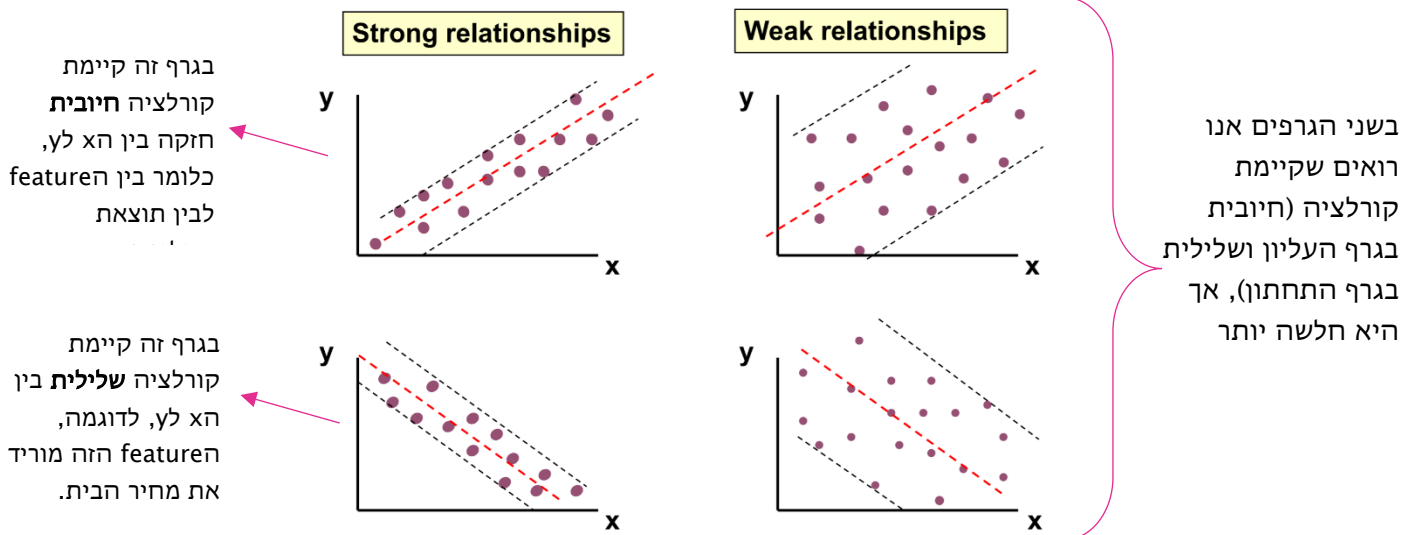
למידה חישובית - שיעור 2

Correlation

אנו מעוניינים למדוד כמה feature כלשהו משפיע על התוצאה של האלגוריתם שלנו. הקורלציה היא היחס בין שני המשתנים (x ו- y בשרטוט שמולנו). המשתנים יכולים להיות קשורים אחד לשני, או לא. בקונטקסט של רגרסיה לינארית- אנו מדברים על הקשר בין feature לבין התוצאה של האלגוריתם הלומד.

בציר ה- x נשים את המשתנה ה"מסביר", לדוגמה גודל מ"ר של בית.

בציר ה- y נשים את המשתנה ה"מוסבר", לדוגמה מחיר הבית.



הערה: חשוב לשים לב שלא בהכרח יהיה קשר בין משתנים, או שיהיה קשר שאינו לינארית. אם הקשר בין המשתנים אינו לינארי, אז הקורלציה לא תהיה מושלמת.

Pearson Correlation

פונקציה המודדת מתי היחס בין 2 משתנים הוא לינארי:

$$r = \frac{\sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^m (x_i - \bar{x})^2 \cdot \sum_{i=1}^m (y_i - \bar{y})^2}} = \frac{\partial_{xy}}{\partial_x \partial_y}$$

r = Sample Pearson correlation coefficient
 m = Number of samples
 x = Value of an explaining variable
 y = Value of the dependent variable

r נקרא מקדם הקורלציה, יכול לקבל ערכים בין מינוס 1 ל-1.

נשים לב שאם ערכי x וגם ערכי y מעל הממוצע או מתחת לממוצע (מעל או מתחת ל \bar{x} , \bar{y}), המונה של הביטוי יהיה חיובי ובגלל שהמכנה תמיד חיובי נקבל r חיובי.

ככל ש r יתקרב ל-0, נקבל יחס לינארי חלש יותר.

ככל ש r יתקרב למינוס 1, נקבל יחס לינארי חזק שלילי.

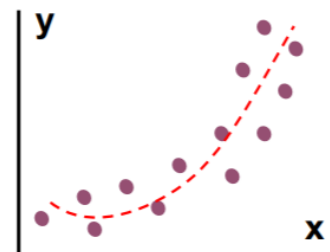
ככל ש r יתקרב ל-0, נקבל יחס לינארי חזק חיובי.

רגרסיה פולינומיאלית

לא נוכל להסביר כל data בצורה לינארית, ופעמים רבות נוכל להסביר את data בצורה פולינומיאלית.

בדוגמה זו נוכל להעביר קו לינארי, אבל אם נעשה זאת לא נסביר את data בצורה המיטבית.

על מנת להסביר את data בצורה המיטבית נשתמש ברגרסיה לינארית (במקרה זה בפונקציה ממעלה שניה).



נוכל להשתמש במנגנון שלמדנו ברגרסיה לינארית, על מנת לייצג מידע ממעלה שניה. אם ברגרסיה לינארית יש feature עמודה אחת שמייצגת מעלה אחת, נוסיף עמודה שתייצר את הfeature בריבוע ונפתור את הבעיה כפי שפתרנו רגרסיה לינארית ובכך יצרנו רגרסיה לפולינום ממעלה 2.

נוכל לבצע את התהליך עבור כל מעלה שהיא.

הערה: נלמד בהמשך הקורס שכל שנעלה את מעלת הפולינום, הסיכוי שנמצא את מה שאנו מחפשים גבוה יותר.

נציג כעת טכניקה נוספת לפתרון בעיית הרגרסיה הלינארית, ללא Gradient Decent. בשונה מ Gradient Decent, הטכניקה הזו לא בהכרח תעבור על כל פונקציה.

$$\mathbf{X} \cdot \vec{\theta} = \vec{Y}$$

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

$$\bar{x}\theta - y : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

הערה: x ו θ קבועים ו θ הוא המשתנה שלנו. כאשר \bar{x} ו θ מכפילים את \bar{x} ו θ בנוסף מעבירים את θ ממימד n למימד m (גדול יותר).

$$Error(\theta) = J(\theta) = \|\bar{x}\theta - y\|^2 : \mathbb{R}^n \rightarrow \mathbb{R}^+$$

בד"כ לא נוכל להביא את הטעות להיות בדיוק 0 כי \bar{x} ו θ מעבירים מרחב קטן למרחב גדול. בנוסף לרוב יש רעש בניסוי או גם כשהתופעה לינארית, כשנמדוד את התופעה לא נקבל לינאריות טהורה, ולכן נחפש קירוב. לכן נמדוד את הקירוב ע"י נורמה בריבוע.

$$\|v\|^2 = \sum_{i=1}^n v_i^2$$

אנו רוצים להביא את הטעות $J(\theta)$ (שהיא הנורמה בריבוע) למינימום. בפונקציה עם מספר משתנים נביא למינימום על ידי השוואת הגרדיאנט ל-0:

$$\nabla J(\theta) = \left(\frac{\partial J}{\partial \theta_j}(\theta) \right)_{j=1}^n$$

אנו רוצים לדעת את הטעות כפונקציה של כל קבוצת משקולות שאני שמה עליהם. נחשב כל אחת מהנגזרות החלקיות

$$\frac{\partial J}{\partial \theta_j}(\theta) = \frac{\partial}{\partial \theta_j} \left(\sum_{i=1}^m \left(\underbrace{x(i, :) \cdot \theta}_{\text{מכפלה של השורה } i \text{ עם הוקטור } \theta} - y^{(i)} \right)^2 \right)$$

$$\frac{\partial}{\partial \theta_j} \sum_{i=1}^m \left(\sum_{k=1}^n x(i, k) \theta_k - y^{(i)} \right)^2 \xrightarrow{\text{נגזור לפי } \theta_j} \sum_{i=1}^m 2 \cdot \left[\sum_{k=1}^n (x(i, k) \theta_k - y^{(i)}) \right] x(i, j)$$

$$2 \cdot \sum_{i=1}^m \underbrace{x(i, j)}_{\text{עמודה}} \left(\underbrace{x(i, :) \cdot \theta - y^{(i)}}_{\text{וקטור ההפרש}} \right) \xrightarrow{\text{מכפלה פנימית}} 2 \cdot x(:, j) \cdot (x\theta - y)$$

קיבלנו את הנגזרת החלקית של J (הטעות) לפי θ_j . כעת נשווה ל-0 כדי לקבל מינימום:

$$2 \cdot x(:, j) \cdot (x\theta - y) = 0 \rightarrow x^T \cdot (x\theta - y) = 0 \rightarrow x^T \cdot x\theta - x^T y = 0$$

$$x^T x \theta = x^T y \quad \text{פונקציה} \quad \xrightarrow{\cdot (x^T x)^{-1}} \quad \underbrace{(x^T x)^{-1} (x^T x)}_I \theta = (x^T x)^{-1} x^T y \rightarrow \theta = (x^T x)^{-1} x^T y$$

הפונקציה: $(x^T x)^{-1} x^T y$ שהתקבלה נקראת $pinv$ (קיצור של pseudo inverse), שנותנת את הקירוב הטוב ביותר ל- x ההפיכה.

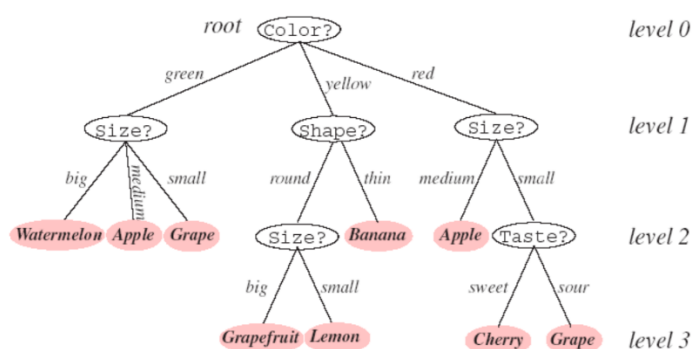
הראנו את pseudo inversen שפותר את בעיית הרגרסיה הלינארית, אבל רק אותה. אם ניתן לפונקציית $pinv$ רגרסיה פולינומיאלית (בהנחה ולא נעשה לבעיה הפולינומיאלית התאמות כלשהן) ולכן אנו מסיקים שהיתרון הגדול של Gradient Decent הוא שהוא עובד בכל שאלת אופטימיזציה.

עצי החלטה

עד כה דיברנו על קיום מפריד לינארי, אך לא דיברנו על הדרך ללמוד מפריד הלינארי.

עצי החלטה היא דרך למיין דברים לפי feature שלהם. זהו עץ עם שאלות כך שכל instance שעובר בעץ יישאל שאלה לגבי feature מסוים של הdata. תהליך הסיווג מתבצע כך שכל instance מתחיל משורש העץ ומתפזר לאורך העץ בהתאם לתשובות לשאלות. בכל צומת בעץ מתבצעת החלטה (ומכאן שמו) בנוגע לסיווג הinstance הנוכחי.

תהליך הסיווג נפסק כאשר הinstance הפך להיות עלה. נסיק מכך שבמקרה הגרוע כל עלה בעץ הוא instance, לכן מספר העלים שווה למספרים הinstance בdata. אנו כמובן יכולים להגיע למצב בו יש מספר instance בכל עלה ולכן זה חסם עליון על מספר העלים בעץ. ברוב המוחלט של המקרים לא נגיע אליו.

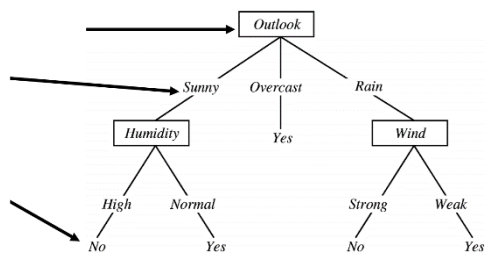


בדוגמה שמולנו, instance של לימון יתחיל בראשית העץ. feature הראשון עליו הוא ישאל יהיה ה-color. בהתאם לתשובתו הוא ירד במורד העץ (ימינה או שמאלה). בהמשך ישאל על feature של ה-shape ולבסוף על feature של ה-size.

בתום תשאול כל הinstance נקבל עץ בינארי.

דוגמה נוספת ופופולרית של עצי החלטה:

נרצה לדעת האם אדם ישחק טניס או לא.



נבדוק מה האפשרות למשחק טניס בכל מזג אוויר. נתקדם בעץ עד אשר נגיע לעלה. בשלב זה נדע האם התשובה היא yes/no.

הערה: ניתן להפוך כל עץ החלטה שאינו בינארי לעץ החלטה בינארי.

האלגוריתם שיבנה את עץ ההחלטה הוא **אלגוריתם הלמידה**. העץ עצמו הוא אלגוריתם הביצוע.

נשים ♥ שאם התשובה לשאלתנו היא מספרים רציפים, נבחר חתכים לערכים אלה ונעבוד לפיהם.

לדוגמה אם נרצה לחשב הסתברות להתקף לב, ואחד מהפיצ'רים שמעניינים אותנו הוא גיל, נמיר את המספר הרציף שמייצג את הגיל לערך רציף כלשהו שיכול להוות תשובה בינארית, לדוגמה: במקום גיל 55 או 65 נשאלה האם גדול מ65.

על ידי תהליך הקטגוריזציה אנו מאבדים את הדיוק ואת המדרג. דוגמה טובה לאיבוד המדרג הינה סיווג צבעים לפי מדרג כלשהו. ירוק הוא כמובן לא יותר גדול או קטן מאדום ולכן אנו מאבדים משהו מהמדרג של המשתנה.

אנו כמובן יכולים להפוך משתנה קטגוריאלי למשתנה רציף, אבל בפעולה זו אנו מייצרים מדרג שאינו קיים.

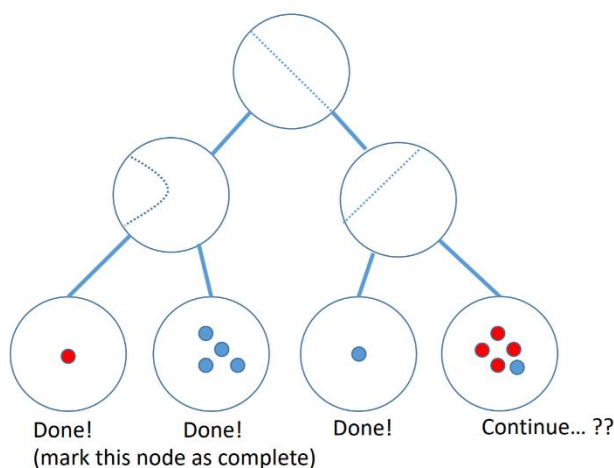
מתי נשתמש בעצי החלטה:

1. כאשר הנתונים שלנו training data הם כאלה שעונים על attribute מסוימים
2. כאשר פונקציית המטרה שלנו אינה רציפה
3. כאשר נרצה להיפחית את המבנה של עץ
4. כאשר data שיש ברשותנו שאנו רוצים ללמוד אינו מסודר

אלגוריתם ליצירת עץ החלטה - מימוש באמצעות תור

- (1) ניצור שורש שמכיל את כל ה-training data שאותו אנו מעוניינים ללמוד.
- (2) נכניס את השורש לתור
- (3) נוציא את ה-noden הבא מהתור.
 - a. אם ה-noden מסווג עם סיווג 1 (perfectly classified), סמן את ה-noden שסיימנו את עבודתנו עמו וחזור לשלב (3), כלומר נוציא את ה-noden הבא מהתור.
 - b. אם ה-noden לא מסווג 1, נשים ב-A את ה-attribute שמחלק את המידע ב-nodes הנוכחי בצורה הטובה ביותר, ונגדיר שה-node הנוכחי מחולק על ידי A.
 - c. ניצור בנים ל-n ככמות החלוקות ש-A משרה ונעביר את כל ה-instances לבנים
 - d. נכניס את כל הבנים לתור.
- (4) אם התור לא ריק, חזור לשלב (3)

איך נדע מתי להפסיק?



ייתכן שנגיע למצב בו ה-attribute שלנו לא מספיקים על מנת לסווג את העץ בצורה מושלמת, או שאנו לא מעוניינים בסיווג מושלם.

הבחירה האם ומתי להפסיק לסווג את ה-instances בעץ תלויה ב-data שלנו ובהחלטה שלנו לגבי רמת הדיוק. לא נתעמק בנושא בקורס.

עולה השאלה איך נבחר *attribute* (תכונה).

אנו נרצה אלגוריתם שיחליט עבורנו בכל שלב (כלומר באופן לוקלי) באיזה *attribute* להשתמש על מנת לפצל את הענף הנוכחי, וכן היכן לשים את הסף לפיצול. אנו כמובן מעוניינים לבחור ב-*attribute* שישירו לנו חלוקה טובה ככל שניתן.

נשתמש במדד אי וודאות כאשר 1 הוא אי וודאות גבוה ביותר, ו0 הוא ודאות מוחלטת. זוהי תכונה של צומת בעץ. אנו מעוניינים למצוא *attribute* כך שניתן יהיה להוריד את אי הודאות ככל שניתן.

Goodness of Split

אנו מעוניינים לדעת האם ה-*attribute* מחלק טוב את ה-*data*. נעשה זאת על ידי שימוש במדד אי הודאות. נרצה לדעת כמה 1 יש לעומת 0.

Impurity Criteria

פונקציה $\varphi: [0,1]^k \rightarrow \mathbb{R}$ המקיימת את התנאים הבאים להסתברויות הבאות $P = (p_1, \dots, p_k) \in [0,1]^k$

$$\varphi(P) \geq 0 \quad (1)$$

$$\exists i \text{ s.t. } p_i = 1 \text{ כאשר } \varphi(P) = 0 \quad (2)$$

$$\varphi(P) = \frac{1}{k} \text{ כאשר } p_i = \frac{1}{k} \text{ } 1 \leq i \leq k \quad (3)$$

$$\varphi(P) \text{ סימטרית ביחס למרכיבים של } P \quad (4)$$

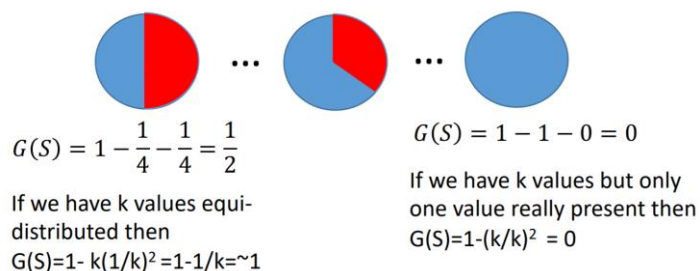
$$\varphi(P) \text{ גזירה אינסוף פעמים בתחום.} \quad (5)$$

Gini Impurity

מדד ג'יני מודד כמה פעמים *instance* שנבחר באופן אקראי מה-*data* שלנו, יתוייג בצורה שגויה, אם תויג באופן אקראי.

$$G(S) \equiv 1 - \sum_{i=1}^c (p_i)^2 = 1 - \sum_{i=1}^c \left(\frac{|S_i|}{|S|} \right)^2$$

תחום הפונקציה הינו $[0,1]$. אם יתקבלו ערכים נמוכים, נסיק כי יש מעט פיזור. אם יתקבלו ערכים גבוהים נסיק כי יש יותר פיזור.



תהליך השימוש בGINI:

1. עבור כל node שעבר פיצול באמצעות attribute כלשהו, נמדוד את הGINI שהקבל טרם הפיצול - $G(S)$.
 2. נמדוד את הGINI שהתקבל לאחר הפיצול - עבור כל הבנים של הnode שפיצלנו.
 3. נבצע ממוצע משוקלל לערכים שקיבלנו ב- (2) - $\sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} G(S_v)$.
 4. נבדוק האם חל שיפור בimpurity בעקבות הפיצול באמצעות הattribute הנוכחי.
 5. נחזור על התהליך עם כל הattributes.
- בתום התהליך המתואר, נבחר את הattribute שהניב לנו את השינוי הגדול ביותר בimpurity.

**הוכחה במודל לכך שבפונקציית GINI יש תמיד שיפור.

כלומר הגדרנו קריטריון לפיצול:

$$\text{GiniGain}(S, A) = \Delta G(S, A) = G(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} G(S_v)$$

- A זה ה-attribute שבאמצעותו אנו מחלקים.
- v זה ערך אפשרי של ה-attribute.
- S זאת קבוצת ה-Instances בnode הנוכחי.

לא דיברנו על כך בשיעור, אבל נכיר את הGini Index:

מחשב פיזור. אם נקבל Gini Index נמוך, נדע שהפיזור נמוך גם כן, ולהפך:

$$G(S) = 1 - \sum_{i=1}^c p_i^2 = 1 - \sum_{i=1}^c \left(\frac{|S_i|}{|S|} \right)^2$$