

החלטה - 2 Decision Trees

עצי החלטה הם תת-קבוצה של אלגוריתמים ממיינים - קלסיפיקציה (Classifiers).

נקבל דאטה אימון = המון וקטורים של תכונות שלגביהם נדע האם הם פלוס או מינוס, נכניס אותם לאלגוריתם הלמידה שיניב עבורנו היפותזה, מודל, ממין (קלאסיפיקציה).

Classification

- Given $\{x_i, y_i\}$ where $y_i \in \{0,1\}$ as training data, determine for a new x if $x \in C_0$ or $x \in C_1$

מושגים פורמליים (עבור מרחב בדיד, דיסקרטי):

- X הוא מרחב הדגימות (x שייך ל- X).
- קונספט c הוא תת קבוצה של מרחב הדגימות, כלומר c שייך לקבוצת החזקה של X . $H = X$. (נקרא גם מרחב ההיפותזה)
- דאטה אימון D training set היא קבוצה של זוגות $\langle x, c(x) \rangle$ כך ש- x היא דגימה ממרחב הדגימות X ו- $c(x) \in \{+1/-1\}$.
- אנחנו מחפשים אחר היפותזה (או מודל) h ששייך ל- H (תת קבוצה של דגימות), כך ש- $c(x) = h(x)$ עבור כל x המיוצג ב- D (או עבור רוב x -ים בקבוצת האימון D).

- Converting from continuous to discrete can be done by quantization (binning):



- Converting from discrete to continuous can be done by encoding:

- Blue = 1, Green = 2, Brown = 3 etc.

יכולים להיות אטריביוטים/פיצורים נומריים = רציפים, וניתר להמיר בין משתנים רציפים לבדידים ולהפך. אבל, כדאי שנעשה זאת בזהירות רבה, מפני שאם ננקוט בשיטה בה ההבדל בין 2 ל-3 גדול מההבדל שבין 3 ל-1 כנראה שהימננו באופן לא נכון.

עצי החלטה

להלן דוגמה לגבי מר סמית' שמחליט האם ללכת לשחק טניס על פי מצב הרוח, הלחות, הטמפרטורה והתחזית (האטריביוטים). להלן דוגמה לדגימות של מר סמית' ועץ ההחלטה שנבנה מהדאטה.

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

מתי משתמשים בעצי החלטה?

- כאשר הדגימות שלנו מתוארות על ידי זוגות של תכונות וערכים.
- כאשר פונקציית המטרה שלנו היא בדידה (לא בהכרח, DT regression).
- כאשר ההיפותזה הנדרשת היא בעלת מבנה לוגי ואינטרפרטציה היא חשובה.
- כאשר יש לנו דאטה אימון שהיא possibly noisy, even inconsistent.

השאלה שלנו השיעור: איך נבנה את העץ הזה באופן הנכון/הטוב ביותר? איך יודעים איזו שאלה שואלים קודם?

יצירת העץ

- ניקח את כל הדגימות וניצור מהם את השורש של העץ
- נאתחל תור של הצמתים בעץ
- כל עוד יש צמתים לא שלמים בתור נבצע:

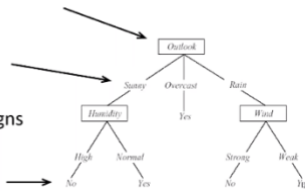
- ניקח את הצומת הבא n
- אם הדגימות ב- n טהורות/מונוכרומטיות (כולם בצבע כחול למשל), נסמן את הצומת כשלם ונמשיך לצומת הבא בתור.
- אחרת, נכניס ל- A את האטריביוט (וערך החסימה של ערכיו) ה"טוב ביותר" עבור הקבוצה ב- n .

- נשים את A כאטריביוט ההחלטה עבור n
- לכל אינטרוול של ערכים של A , ניצור צומת-בן לצומת n
- נציג את הדגימות לילדים של n
- נכניס את כל הילדים הלא ריקים לתור

איך נדע מיהו אטריביוט ההחלטה הטוב ביותר?

בדוגמה משמאל יש לנו שני אטריביוטים $A1, A2$. איך נדע איזה עדיף? נשים לב שהילדים של האטריביוט $A2$ מניב ילדים כמעט-מונוכרומטים, מה שמקרב אותנו להתפלגות טהורה, מה שנותן לנו להרוויח יותר אינפורמציה. לכן, נרצה להפחית את רמת האי-טהורות (או הבילבול, או אי הוודאות) בצעד הבא שלנו.

- Each internal node tests an attribute
- Each branch corresponds to an attribute value
- Each leaf node assigns a classification

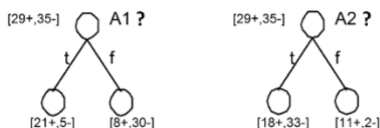


Create the root node with all samples

Insert the node to initialize a queue, Q

While there are more incomplete nodes in Q do:

- Get next node n
- If training examples in n are perfectly classified Then mark node as complete and continue to next node in Q
- Else $A \leftarrow$ the "best" decision attribute (and boundary values) for the set in n
 - Assign A as the decision attribute for n
 - For each interval of values of A , create a new descendant of n
 - Distribute training examples to descendant nodes
 - Insert all (non empty) descendant nodes to Q



פונקציית אי וודאות היא פונקציה המקבלת התפלגות דיסקרטית ומחזירה מספר ממשי, המקיימת את התנאים הבאים עבור וקטור ההסתברויות:

for probability distributions $P = (p_1, \dots, p_k) \in [0,1]^k$:

- $\varphi(P) \geq 0$
- The minimal value is attained when $\exists i$ s.t. $p_i = 1$.
- The maximal value is attained when $1 \leq \forall i \leq k, p_i = 1/k$.
- It is symmetric with respect to the components of P
- It is smooth (infinitely differentiable) in the relevant range



- פונקציה זו צריכה להיות אי שלילי (הכול ודאי 0)
- הערך המינימלי = אי וודאות מינימלי = מקסימום וודאות, מתקיים כאשר יש מאורע שקורה תמיד, כלומר קיים i כך ש- $p_i = 1$. העיגול הכחול המושלם
- אי וודאות מקסימלית מתקיימת כאשר ההתפלגות היא יוניפורמית = ההתפלגות של כל המאורעות היא יוניפורמית עם אותה ההסתברות. העיגול החצוי העליון
- פונקציה זו תהיה סימטרית בהתייחס לרכיבים של הוקטור P (וקטור ההסתברויות)
- היא חלקה ככל האפשר (נגזרות רציפות למען נוחות מתמטית).

פונקציית Gini Impurity

- מודד חוסר שוויון בכלכלה
- מניב טווח ערכים 0-1, כך שערכים נמוכים משמעותם אי וודאות פחותה (יותר וודאות), ערכים גבוהים משמעותם אי וודאות גבוהה.



$$G(S) = 1 - \frac{1}{4} - \frac{1}{4} = \frac{1}{2}$$

If we have k values equi-distributed then $G(S) = 1 - k(1/k)^2 = 1 - 1/k = \sim 1$

$$G(S) = 1 - 1 - 0 = 0$$

If we have k values but only one value really present then $G(S) = 1 - (k/k)^2 = 0$

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

(נזכיר שפונקציה זו מקבלת S משום שהיא מקבלת תת קבוצה של דגימות שנמצאים באותו הצומת) להלן דוגמה לחישוב אי וודאות על ידי פונקציית Gini: כאשר יש לנו שני צבעים $c=2$

בחזרה לנושא השני: אנו מחפשים אחר אטריביוט A שיניב את ממוצע משוקלל ה-Gini הטוב ביותר לאחר פיצול, ולכן נגדיר את פונקציית הרווח, שמחשבת לנו את הוודאות שהרווחנו מהפיצול על ידי אטריביוט ספציפי A , על תת הקבוצה של הדגימות S של הצומת הנוכחי:

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

↑ Change in Impurity
↑ Impurity Before Split
↑ Weighted Average of Impurity of All Groups After Splitting

נוסחה זו מחשבת כמה אי וודאות הייתה לי לפני הפיצול ומחזירה כמה אי וודאות משוקללת הרווחתי אחרי הפיצול על ידי האטריביוט. אנחנו נרצה למקסם את פונקציית ה-Gain. $(S_v$ כל מי שבתוך S ועונה לערך של v)

פונקציית האנטרופיה Entropy

נשתמש בנוסחה של אנטרופי כפונקציית אי וודאות עבור משתנה רנדומי X שלוקח n ערכים שונים עם הסתברויות p_i :

$$H(X) = - \sum_{i=1}^n p_i \log(p_i)$$

(לוג בסיס 2, כאשר ידוע לנו כי לוג על ערכים בין 0-1 הוא שלילי ולכן המינוס)

- מודד את המידע הממוצע שמתקשר לתוצאה של משתנה מקרי.
- יהיו קבוצה S ודאטה עם c (יתכן יותר מ-2) קלאסים. P_i היא החלק היחסי של קלאס i בקבוצה S . האנטרופיה של S :

$$Entropy(S) = - \sum_{i=1}^c p_i \log_2 p_i = - \sum_{i=1}^c \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}$$

נמיר את האנטרופיה לנוסחת רווח $Gain(S, A) =$ פונקציית הרווח הצפוי של האנטרופיה לאחר פיצול על פי אטריביוט A .

$$\Delta \phi(S, A) \equiv \phi(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \phi(S_v)$$

תכונות נוספות של אנטרופיה:

- הערך המקסימלי שלה (בה יש מקסימום אי וודאות), הינה כאשר ההתפלגות

$$H\left(P = \left(\frac{1}{k}, \dots, \frac{1}{k}\right)\right) = \log k.$$

יוניפורמית/אחידה כיאה לפונקציית אי וודאות:

- **נרצה להוכיח איפה אנטרופיה מקסימלית:**

נחמש בלמה 1: פונקציית לוג היא פונקציית "עצובה" קעורה כלפי מטה (נגזרת שניה של \log היא שלילית).

ובלמה 2: אי שוויון ינסן (מופיע משמאל <<) אשר הוא נכון לכל פונקציה עצובה.

Lemma 2 (Jensen's Inequality):

The following holds for any sad function f :

$$\forall x_1, \dots, x_k \text{ and } \forall \lambda_1, \dots, \lambda_k \in [0,1] \text{ s.t. } \sum \lambda_i = 1.$$

$$f\left(\sum_{i=1}^k \lambda_i x_i\right) \geq \sum_{i=1}^k \lambda_i f(x_i).$$

$$H(p_1, \dots, p_k) = - \sum_{i=1}^k p_i \log(p_i)$$

$$= \sum_{i=1}^k p_i \log\left(\frac{1}{p_i}\right)$$

$$\leq \log\left(\sum_{i=1}^k p_i \cdot \frac{1}{p_i}\right)$$

$$= \log k$$

$$f\left(\sum_{i=1}^k A_i x_i\right) \geq \sum_{i=1}^k A_i f(x_i)$$

הצעד האחרון בהוכחה:

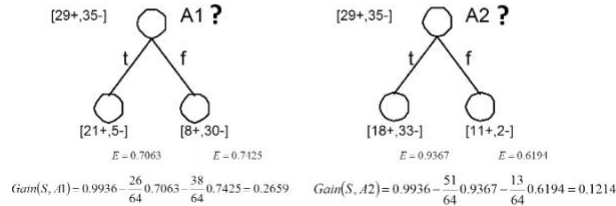
- השתמשנו בשתי הלמות לעיל במטרה להראות היכן אנטרופיה מקסימלית.

בחזרה לדוגמה שלנו: נחשב את הרווח האנטרופי עבור שני האטריביוטים A1, A2

כך נדע איזה מבין האטריביוטים עדיף, נשים לב שהרווח גבוהה יותר עבור A1 כלומר הוא מניב עבורנו אי וודאות פחות מאטריביוט A2.

Which attribute is best?

$$Entropy([29+, 35-]) = -\frac{29}{64} \log_2\left(\frac{29}{64}\right) - \frac{35}{64} \log_2\left(\frac{35}{64}\right) = 0.9936$$



The information gain values for the 4 attributes are:

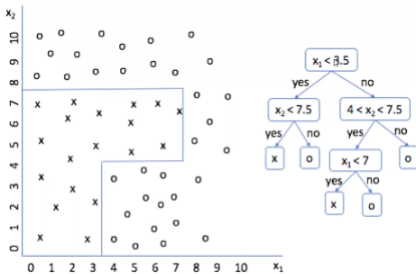
- Gain(S, Outlook) = 0.247
- Gain(S, Humidity) = 0.151
- Gain(S, Wind) = 0.048
- Gain(S, Temperature) = 0.029

where S denotes the collection of all training examples

נחזור חזרה לדוגמה ההתחלתית עם מר סמית': יש לנו 14 אינסטנסים (בעולם האמיתי יהיו הרבה יותר), ונרצה להחליט על מה שואלים: תחזית, טמפרטורה, רוח וכו'. נחשב עבור הדגימות שלנו ועבור האטריביוטים את הרווח האנטרופי במטרה להחליט על השאלה שתישאל ראשונה. התקבל **השתחזית (Outlook) הוא האטריביוט שמניב את הרווח הגבוהה ביותר** ולכן הוא ישאל ראשון.

נוכל לעבוד גם עם דאטה שאיננו קטגורי: להלן דוגמה (DTs in continuous space)

DTs in continuous space

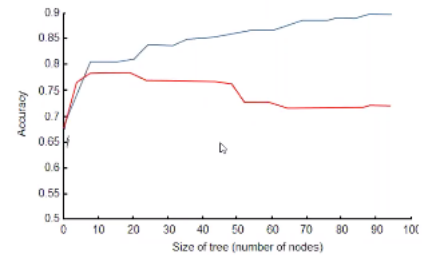


עובדות נוספות

- גרדיי אינו אופטימלי – הגישה שהצגנו היא גישה חמדנית, אנו מחפשים בכל שלב אחר הרווח המקסימלי עבור אותו שלב בעץ.
- הרווח תמיד אי שלילי – עבור כל פיצול שנעשה, נרוויח וודאות.

תמיד נרצה להימנע מ-Overfitting! לכן לא נרצה להרחיב את העץ יותר מידי כדי שלא יתאים לדאטה הנתון יותר מידי ויוכל להגיע להחלטה עם מצב חדש. כדי לעשות זאת נשתמש בדאטה סט נוסף שיקרא "validation" set. נלמד את הדאטה על ידי הדאטה אימון אבל נמדוד את הטעות על ה-validation set. (להלן גרף המתאר את הדיוק, לכאורה מדויק יותר ככל שיש לנו יותר אטריביוטים) מכאן שנרצה להפסיק להרחיב את העץ לפני שנגיע למונוכרומטיות מושלמת (לפני שגיא 0), או להרחיב את העץ לשגיאה 0 ואז לבצע קיצוץ post-prune, או לשלב בין שתי השיטות.

גישת Chi Square ששואלת – האם פיצול לפי אטריביוט מועמד נותן לנו התפלגות לפי קלאס שיש לה יותר כוח מהנוכחית? האם הבנים מאוד דומים לאבא? הם הרווח הוא סטטיסטית משמעותי? (הרחבה בתרגול 2)



בעיות נוספות

- התמודדות עם **ערכים חסרים** של חלק מהאטריביוטים : השלמת דאטה או לקיחת הממוצע
- חיפוש אחר ערך חותך, נקודת הסף, **באטריביוטים רציפים**
- פיצול אינפרמציה ו- **Gain Ratio** עבור אטריביוטים עם ערכים רבים - מבחינת החישוב יש יתרון עבור אטריביוטים שיש להם מספר ערכים גבוהה יותר ולכן נגזיר Gain Ratio (הרחבה בתירגול)
- הכללת העלות עבור שמירת האטריביוטים, ופונקציית שגיאה ממושקלת
- גבולות מורכבים