

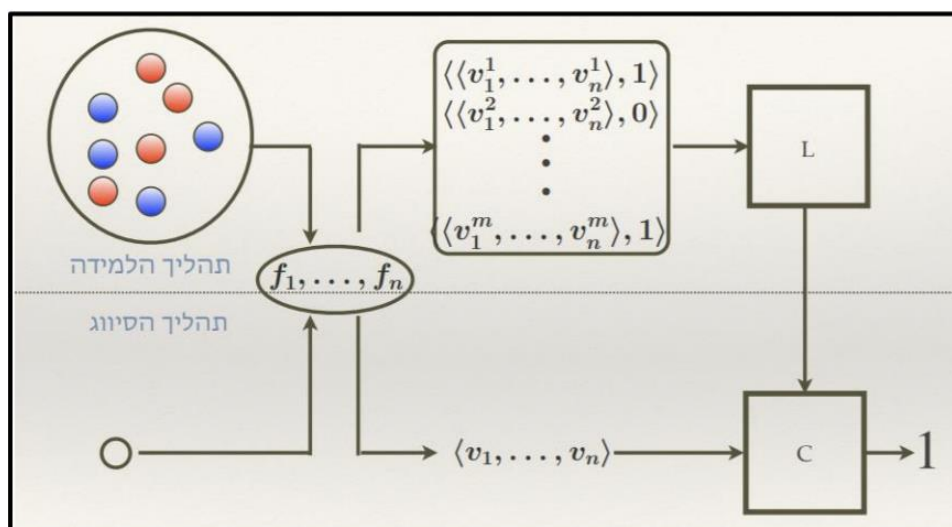
אלגוריתם לסיווג דוגמאות תחת הגבלת תקציב

לרכישת ערכי תכונות

פרויקט בבינה מלאכותית

236502

הפקולטה למדעי המחשב, טכניון



מגישים:

sivan.y@campus.technion.ac.il	318981586	סיון יוסבשוילי
razle@campus.technion.ac.il	316579275	רז לוי

תוכן עניינים

מבוא	3
אלגוריתמים נאיביים	5
PCA	9
אלגוריתם ביניים	11
פונקציית score	13
אלגוריתם חיפוש בגרף	16
אלגוריתם חיפוש מקומי	21
אלגוריתם חיפוש בהשראת אלגוריתמים גנטיים	24
מתודולוגיית ניסויים	28
ולידציה	30
ניסוי: פונקציית המחיר הטובה ביותר	31
ניסוי: מסווג הבסיס הטוב ביותר	34
ניסוי: אלגוריתם הלמידה הטוב ביותר	37
ניסוי: זמן ריצה כתלות בכמות התכונות	42
ניסוי: אלגוריתם חיפוש מקומי הטוב ביותר	45
סיכום	48

מבוא

הרעיון שלנו לפרויקט עלה מתוך הצורך לענות על בעיה רפואית מציאותית מחיי היום-יום.

נציג כעת דוגמא שתלווה אותנו לכל אורך הפרויקט ותדגים את הצורך הכללי שאנו רוצים להציע לו פתרון - ישראל ישראלי, רוצה לגשת לרופא המשפחה שלו, ולבקש ממנו שיאמר לו מה הסיכויים שבעתיד יחלה במחלה כלשהי.

הרופא יכול לשלוח את ישראל לסדרת בדיקות ארוכות (ויקרות) שבסופן יוכל לענות לישראל על השאלה עם וודאות די גבוהה.

פעמים רבות, בדיקות מיוחדות אינן נמצאות בסל הבריאות, ולכן ישראל יאלץ לשלם עליהן.

ישראל אינו יכול להרשות לעצמו להוציא סכומי עתק, ולכן הוא מאתגר את הרופא –

"בהינתן שאני מוכן לשלם X שקלים על בדיקות רפואיות, מה הן הבדיקות שעליי לעבור כדי שתוכל לומר לי

בוודאות גבוהה מה הסיכויים שבעתיד אחלה במחלה כלשהי?"

באופן כללי, פעמים רבות כשאנו נתקלים בבעיית סיווג, הדוגמא אותה אנו מנסים לסווג לא מכילה את כל המידע (יש בה חלקים חסרים).

נרצה תחת אילוץ התקציב, לרכוש את ערכי התכונות החסרות על מנת שבסופו של דבר, יהיה לנו מספיק מידע כדי לסווג את הדוגמה כהלכה.

ניסוח הבעיה באופן פורמלי-

בהינתן:

סט $\{X, Y\}$, אשר מורכב מקבוצת דוגמאות $X = \{\bar{x}_1, \dots, \bar{x}_m\}$, $x_i \in R^D$, כאשר D הוא מספר התכונות וקבוצת תיוגים בינאריים $Y = \{y_1, \dots, y_m\}$, כאשר y_i הוא התיוג לדוגמה x_i .

$C = \{c_1, \dots, c_D\}$ קבוצת המחירים עבור כל תכונה, כאשר c_i היא עלות התכונה ה- i (המחיר הוא אי שלילי).

נרצה לייצר מודל, כך שבהינתן דוגמא חדשה $\hat{x} \in R^k$, כאשר $k < D$ (בה חסרים $(D - k)$ תכונות) ותקציב מקסימלי \hat{C} , מחזיר בכל שלב, איזו תכונה הוא מעוניין לרכוש על מנת לסווג את הדוגמה.

בסיום התהליך, יפלוט המודל את הסיווג של הדוגמה.

לאורך הדו"ח נציג מספר גישות לפתרון הבעיה.
 באופן כללי, כל האלגוריתמים שנציע יהיו בעלי מבנה של אימון וסיווג כאשר מטרת האימון היא לאסוף מידע על הדוגמאות ולחקור את הקשרים בין התכונות.
 מטרת הסיווג היא להציע סיווג לדוגמה, כאשר במהלך הסיווג, האלגוריתם יבחר לרכוש תכונות (תחת מגבלת התקציב), כאשר הבחירה לרכוש תכונה בכל שלב תתבסס על המידע שצבר האלגוריתם עד לאותו שלב.

פסאודו קוד לאלגוריתם הכללי (ממומש בקובץ abstract_algorithm.py)

אימון:

• קלט:

(1) קבוצה $\{X, Y\}$ כאשר $X = \{\bar{x}_1, \dots, \bar{x}_m \mid \bar{x}_i \in \mathbb{R}^D\}$ קבוצת דוגמאות (D מספר התכונות) ו- $Y = \{\bar{y}_1, \dots, \bar{y}_m \mid \bar{y}_i \in \{0,1\}\}$

(2) קבוצת מחירים לתכונות $C = \{\bar{c}_1, \dots, \bar{c}_D \mid \bar{c}_i \in \mathbb{Z}\}$ בהתאמה.

...

סיווג:

• קלט:

(1) דוגמה $\hat{x} \in \mathbb{R}^k$ כאשר $k < D$ מספר התכונות האפשרי לכל דוגמה.

(2) תקציב מקסימלי \tilde{C} .

(3) קבוצת אינדקסים K של התכונות הנתונות.

• פלט:

(1) סיווג הדוגמה.

...

האלגוריתמים "הנאיביים"

תחילה, חשבנו על האלגוריתם "הנאיבי" שיאפשר את פתרון הבעיה. החשיבה על אלגוריתם "נאיבי" תעזור לנו למפות את הבעיות שעלולות לעלות במהלך הפתרון, למצוא את "צוואר הבקבוק" של הבעיה ובכך להתקדם אל עבר מציאת אלגוריתם טוב יותר. במהלך החשיבה, עלו שתי גישות אפשריות לאלגוריתם "נאיבי" שכזה.

גישה א' - האלגוריתם הריק

כאשר מטופל מגיע לרופא בבקשה לקבלת אבחנה רפואית ובהינתן שידועים לרופא $k < D$ תוצאות בדיקות על המטופל (D המספר המקסימלי לתוצאות הבדיקות), ייתכן כי הרופא יבחר לנסות ולאבחן את המטופל בעזרת k התכונות הללו בלבד.

נציג כעת את האלגוריתם בצורה פורמלית. לצורך כך נבחר אלגוריתם למידה כלשהו (כדוגמת KNN או ID3).

אימון:

• קלט:

(1) קבוצה $\{X, Y\}$ כאשר $X = \{\bar{x}_1, \dots, \bar{x}_m \mid \bar{x}_i \in \mathbb{R}^D\}$ קבוצת דוגמאות (D מספר התכונות)

ו- $Y = \{\bar{y}_1, \dots, \bar{y}_m \mid \bar{y}_i \in \{0,1\}\}$

(2) קבוצת מחירים לתכונות $C = \{\bar{c}_1, \dots, \bar{c}_D \mid \bar{c}_i \in \mathbb{Z}\}$ בהתאמה.

1. שמור את קבוצת הדוגמאות והתיוגים $\{X, Y\}$.

סיווג:

• קלט:

(1) דוגמה $\hat{x} \in \mathbb{R}^k$ כאשר $k < D$ ו- D מספר התכונות האפשרי לכל דוגמה.

(2) תקציב מקסימלי \tilde{C} .

(3) קבוצת אינדקסים K של התכונות הנתונות.

• פלט:

(1) סיווג הדוגמה.

1. הרץ את שלב האימון של אלגוריתם הלמידה שבחרת עבור קבוצת הדוגמאות והתיוגים ששמרת בשלב האימון.
2. הרץ את שלב הסיווג של אלגוריתם הלמידה שבחרת.
3. החזר את הפלט שהתקבל.

מימשנו את "האלגוריתם הריק" בקובץ `naive_algorithm.py`. ניתן להבחין שהבעיה המרכזית באלגוריתם היא הוויתור על האופציה לקבלת מידע נוסף. על אף שברשותנו תקציב לרכוש תכונות נוספות ובכך להקטין את ה-loss (שגיאת הסיווג) ולהגיע לאחוזי דיוק (accuracy) טובים יותר- אנו "מסתפקים" בתכונות הנתונות לנו ולכן ה-loss גבוה וה-accuracy נמוך. כדאי לציין שהאלגוריתם הנ"ל הוא המהיר ביותר האפשרי, שכן אנו לא מבצעים פעולות כלל מלבד הרצת אלגוריתם סיווג מוכר מהקורס.

גישה ב'- האלגוריתם הרנדומלי

כאשר מטופל מגיע לרופא בבקשה לקבלת אבחנה רפואית ובהינתן שידועים לרופא $k < D$ תוצאות בדיקות על המטופל (D המספר המקסימלי לתוצאות הבדיקות), ייתכן כי הרופא יבקש מהמטופל לבצע בדיקות רנדומליות כל עוד ברשות המטופל תקציב מספק.

נציג את האלגוריתם בצורה פורמלית. לצורך כך נבחר אלגוריתם למידה כלשהו (כדוגמת KNN או ID3).

אימון:

• קלט:

(1) קבוצה $\{X, Y\}$ כאשר $X = \{\bar{x}_1, \dots, \bar{x}_m \mid \bar{x}_i \in \mathbb{R}^D\}$ קבוצת דוגמאות (D מספר התכונות) ו- $Y = \{\bar{y}_1, \dots, \bar{y}_m \mid \bar{y}_i \in \{0,1\}\}$.

(2) קבוצת מחירים לתכונות $C = \{\bar{c}_1, \dots, \bar{c}_D \mid \bar{c}_i \in \mathbb{Z}\}$ בהתאמה.

1. שמור את קבוצת הדוגמאות והתיוגים $\{X, Y\}$.
2. שמור את קבוצת המחירים לתכונות C .

סיווג:

• קלט:

(1) דוגמה $\hat{x} \in \mathbb{R}^k$ כאשר $k < D$ - מספר התכונות האפשרי לכל דוגמה.

(2) תקציב מקסימלי \tilde{C} .

(3) קבוצת אינדקסים K של התכונות הנתונות.

• פלט:

(1) סיווג הדוגמה.

1. חשב את קבוצת התכונות שטרם נבחרו.

2. כל עוד יש עדיין תכונה שלא נבחרה ותקציב הרכישה גדול מאפס:

2.1. בחר תכונה רנדומלית משלב אימון.

2.2. הסר אותה מקבוצת התכונות שטרם נבחרו.

2.3. אם מחיר התכונה שבחרת קטן או שווה לתקציב רכישת התכונות:

2.3.1. החסר את מחיר התכונה מתקציב רכישת התכונות.

2.3.2. הוסף את התכונה לקבוצת התכונות שנבחרו ("רכוש את התכונה").

3. לאחר שחישבת את קבוצת התכונות הנבחרות, הרץ את שלב האימון של אלגוריתם הלמידה שבחרת

(כאשר התכונות הנתונות לכל דוגמה הן התכונות בקבוצת התכונות הנתונות שחישבנו בשלב 2) עבור

קבוצת הדוגמאות והתייגים ששמרת בשלב האימון.

4. הרץ את שלב הסיווג של אלגוריתם הלמידה שבחרת.

5. החזר את הפלט שהתקבל.

מימשנו את "האלגוריתם הרנדומלי" בקובץ `naive_algorithm.py`.

בהשוואה "לאלגוריתם הריק", כעת אנו אכן מנצלים את התקציב שברשותנו לקבלת מידע נוסף, אלא שאנו

עושים זאת בצורה רנדומלית ולכן אנו למעשה עלולים לרכוש תכונה שלא בהכרח מוסיפה מידע, או שמוסיפה

מעט מידע ביחס לתכונות נוספות שיכולנו לרכוש.

בנוסף, ניתן להבחין שאחת הבעיות באלגוריתם היא חוסר האופטימליות - ייתכן שנבחר קבוצת תכונות קטנה

שגודלה קטן מגודל מהקבוצה האופטימלית (למשל אם קיימת תכונה שמחירה כל התקציב ומספר תכונות

שסכום מחיריהן הוא כל התקציב - אם באופן רנדומלי נבחרה התכונה היקרה, לא נוכל להמשיך לרכוש תכונות

בעוד שאם היינו בוחרים רנדומלית בכל פעם את התכונות הזולות, היינו יכולים לרכוש מספר רב יותר של

תכונות). כלומר, באלגוריתם הנ"ל ייתכן שנוסיף מעט מידע. יש לציין שהמצב הזה לא בהכרח יפגע בסיכוינו לסווג כראוי, שכן ייתכן שתכונה אחת תהיה משמעותית יותר ממספר תכונות.

- נוכל באופן דומה לחשוב על אלגוריתם חמדן שבשלב האימון יבחר בכל פעם את התכונה הזולה ביותר ובכך נקבל קבוצת תכונות בגודל אופטימלי. האלגוריתם האופטימלי הנ"ל ממומש גם הוא בקובץ `naive_algorithm.py` בשם "OptimalAlgorithm". יתרונות וחסרונות האלגוריתם יהיו זהים "לאלגוריתם הרנדומלי" לרבות הדיון על היחס בין הגודל האופטימלי לתרומת התכונות.

מהדיון על האלגוריתם "הנאיבי" אנו מסיקים שעלינו לנצל מידע נוסף (כלומר לרכוש תכונות).

כמו כן, נסיק מספר נקודות חשובות שעלינו לחקור במטרה להגיע לאלגוריתם מספק:

1. מציאת דרך לבחור את התכונה שתוסיף לנו את המידע הרב ביותר. על מנת לעשות זאת, עלינו להגדיר "פונקציית מחיר" לכל תכונה שתיתן ביטוי לתרומת התכונה לסיווג נכון.
2. עם זאת, ראינו שישנו Trade-off בין אופטימליות מספר התכונות לתרומתן. עלינו למצוא דרך להתמודד עם ה-Trade-off הנ"ל, שכן למשל אם התכונה שמוסיפה את המידע הרב ביותר היא תכונה A , ייתכן שנעדיף דווקא לרכוש שתי תכונות זולות יותר תכונה A , אך שתיהן יחד מוסיפות מידע גדול יותר מאשר תכונה A .
3. עלינו למצוא דרך להתמודד עם בעיות פרקטיות כגון "ניבוי המחיר הנתון", שכן בזמן האימון איננו מודעים לתקציב שיינתן רק בשלב הסיווג, לכן עלינו למצוא דרך "לנבוא" את גודל התקציב שיינתן בשלב הסיווג.

PCA

במהלך סיעור המוחות שקיימנו, עלה רעיון להשתמש ב-PCA.

נסביר בקצרה את הקשר של הנושא לבעיה שאנו דנים בה ומדוע הוא אינו כיוון פתרון נכון.

ראשית, נחدد שתי הגדרות חשובות-

בואנו להשתמש ב- $k < d$ תכונות, קיימות שתי גישות:

- Feature selection – בחירת תת קבוצה $F' \subset F$ של תכונות מתוך התכונות הקיימות.
- Dimensionality Reduction – יצירת תכונות חדשות המשתמשות בתכונות הקיימות

$$F' = \{\phi_1(F), \dots, \phi_k(F)\}$$

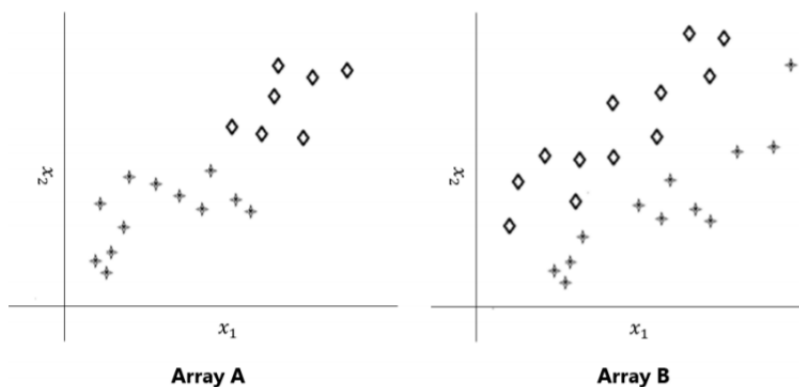
שיטת ה-PCA (Principal Components Analysis, בעברית: ניתוח גורמים ראשוניים) הינה שיטה שמטרתה הורדת מימד של dataset נתון.

שיטה זו מתמירה באופן לינארי את ה-dataset למערכת קואורדינטות חדשה, בה המידע בקואורדינטות החדשות הוא אורתוגונלי עם שונות הולכת וקטנה.

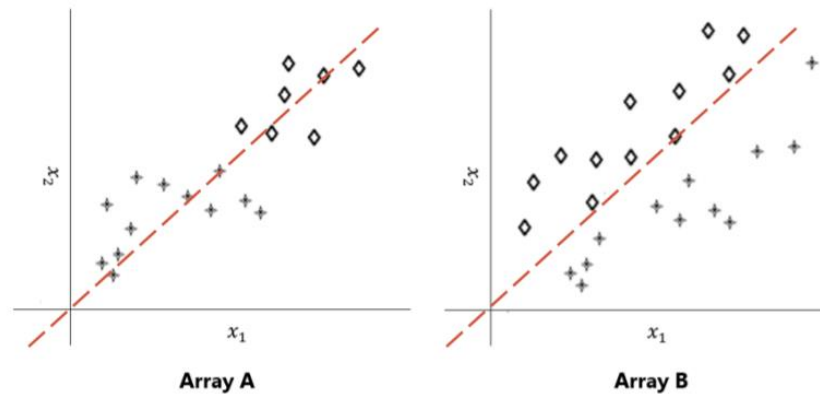
שיטה זו נפוצה בעיקר להורדת מימד בסטטיסטיקה ולמידת מכונה (הומצאה ב-1901 על ידי קרל פיטרסון).

חסרון בולט של השיטה הוא התעלמות מתגיות. כך לדוגמה, אם נבחר לפרק לגורם ראשוני יחיד, באיור המצורף שני התרשימים יהיו בעלי אותו גורם ראשוני, למרות שרק עבור array A הגורם הראשוני ישמור על האפשרות להפריד את הדאטה לשני התיוגים השונים ואילו עבור Array B הוא יפגעו ביכולת זו (שברור שקיימת – ה-dataset פריד לינארית).

ה-dataset המקורי:



לאחר פירוק לגורמים ראשוניים :



הסיבה הראשונה שלא ניעזר ב-PCA היא עקב אכילס זה של השיטה.

יתר על כן, נבחין ש-PCA לא מבצע בחירת תכונות (feature selection) כפי שאנו מנסים לעשות בבעיה שלנו, אלא יוצר תכונות חדשות מתוך התכונות הקיימות (שהינן קומבינציה של התכונות הקיימות). זו סיבה נוספת לכך שבחרנו לא להיעזר ב-PCA.

אלגוריתם ביניים

בדרך אל פיתוח אלגוריתמים מתוחכמים ואופטימליים במובנים מסוימים, חשבנו על אלגוריתם מתוחכם באופן פשוט שיכול להוות פתרון ביניים להשוואה בהמשך.

אלגוריתם שונות מקסימלית

כאשר מטופל מגיע לרופא בבקשה לקבלת אבחנה רפואית ובהינתן שידועים לרופא $k < D$ תוצאות בדיקות על המטופל יכול הרופא להציע למטופל לעבור את הבדיקות ה"אקספרסיביות" ביותר. לדוגמה אם יודעים שבדרך כלל כשעוברים את בדיקה x מקבלים המון מידע, ימליץ הרופא למטופל לעבור את בדיקה x . נציג כעת את האלגוריתם בצורה פורמלית. לשם כך נבחר אלגוריתם למידה.

אימון:

• קלט:

(1) קבוצה $\{X, Y\}$ כאשר $X = \{\bar{x}_1, \dots, \bar{x}_m \mid \bar{x}_i \in \mathbb{R}^D\}$ קבוצת דוגמאות (D מספר התכונות)

ו- $Y = \{\bar{y}_1, \dots, \bar{y}_m \mid \bar{y}_i \in \{0,1\}\}$

(2) קבוצת מחירים לתכונות $C = \{\bar{c}_1, \dots, \bar{c}_D \mid \bar{c}_i \in \mathbb{Z}\}$ בהתאמה.

1. חשב לכל תכונה את השונות שלה. שמור רשימה של התכונות ממיונות לפי השונות מהגבוהה לנמוכה.

סיווג:

• קלט:

(1) דוגמה $\hat{x} \in \mathbb{R}^k$ כאשר $k < D$ ו- D מספר התכונות האפשרי לכל דוגמה.

(2) תקציב מקסימלי \tilde{C} .

(3) קבוצת אינדקסים K של התכונות הנתונות.

• פלט:

(1) סיווג הדוגמה.

1. כל עוד יש עדיין תכונה שלא נבחרה ותקציב הרכישה גדול מאפס:
 - 1.1. מצא את התכונה שאינה נבחרה עם השונות המקסימלית.
 - 1.2. הסר אותה מקבוצת התכונות שטרם נבחרו.
 - 1.3. אם מחיר התכונה שבחרת קטן או שווה לתקציב רכישת התכונות:
 - 1.3.1. החסר את מחיר התכונה מתקציב רכישת התכונות.
 - 1.3.2. הוסף את התכונה לקבוצת התכונות שנבחרו.
2. לאחר שחישבת את קבוצת התכונות הנבחרות, הרץ את שלב האימון של אלגוריתם הלמידה שבחרת עבור קבוצת הדוגמאות והתיוגים ששמרת בשלב האימון.
3. הרץ את שלב הסיווג של אלגוריתם הלמידה שבחרת.
4. החזר את הפלט שהתקבל.

מימשנו את "אלגוריתם השונות המקסימלית" בקובץ `mid_algorithm.py`. נשים לב כי אלגוריתם זה מתוחכם יותר מאלגוריתמים קודמים, שכן הוא מבצע בחירה של תכונות, והיא נעשית תוך לוגיקה בסיסית שאמורה לקדם אותנו אל עבר המטרה. אלגוריתם זה הינו אלגוריתם מהיר, והלוגיקה שהוספנו לבחירת התכונות "זולה" ביחס לתועלת שתביא בבחירה חכמה של תכונות (שיפור במיקום על עקומת `accuracy effort tradeoff`). חיסרון של גישה זו היא להסתכל על תכונות טוב אך ורק מנקודת מבט אחת (תכונה טובה היא תכונה שהשונות שלה גבוהה). אמירה זו אולי נכונה, אך כמובן שאפשר לפתח אותה באופן מעמיק ומתוחכם יותר, ונצפה שפיתוחים אלו יביאו לשיפור בתוצאות אל מול אלגוריתם בסיס זה. לכן, נרצה להמשיך ולחקור מה היא תכונה "טובה" בבעיה שלנו.

פונקציית score

נגדיר פונקציה שתעזור לנו להעריך כמה טוב יהיה להוסיף תכונה ספציפית כאשר יש לנו סט תכונות מסוים.

לצורך הפרויקט, נגדיר שתי פונקציות score ובבחן את התנהגויות האלגוריתמים עם שתי הפונקציות האלו. בשתי הפונקציות כמובן, ניקח בחשבון את המחיר של התכונה הנוספת לסט התכונות הנתון.

- פונקציה A:

מוטיבציה – נרצה לראות כי הוספת תכונה תעזור לנו בחיזוי הסיווג. לשם כך נרצה לראות כי היא אכן קורולטיבית לתווי המטרה (לדוגמה, אם אנו מתעסקים בבעיית חיזוי מצטייני סמסטר, כנראה שלא נרצה לקנות תכונה כמו "מספר עוקבים ברשתות החברתיות").

נרצה להבחין כי תכונה מוסיפה לנו מידע חדש שלא קיים לנו בסט התכונות הקיימות (לדוגמה, אם יש לנו תכונה שהיא גיל, לא נרצה לקנות תכונה שהוא פונקציה לינארית של הגיל).

כמובן שכיאה לבעיה כמו שלנו, ככל שהתכונה יקרה יותר, כך בהיבט כולל פחות משתלם לנו לקנות אותה. לסיכום, נסתכל על שלושה פרמטרים-

(a) תכונה "טובה" היא תכונה שקורולטיבית לתגית.

(b) תכונה "טובה" היא תכונה שאינה קורולטיבית לתכונות אחרות.

(c) תכונה "טובה" היא תכונה זולה.

לפיכך, נגדיר את הפונקציה הבאה:

def ScoreFunctionA(features, new_feature, γ):

return
$$\frac{a(\text{new_feature}) + \gamma \cdot b(\text{new_feature}, \text{features})}{c(\text{new_feature})}$$

בהמשך נבצע ולדיציה בה נכוון את ההיפרפרמטר γ .

- פונקציה B:

מוטיבציה – לבעיה כמו שלנו, ככל שהתכונה יקרה יותר, כף בהיבט כולל פחות משתלם לנו לקנות אותה. נרצה להיות בטוחים ככל הניתן כי הסיווג שלנו נכון. לשם כך נגדיר את המושג הבא:

ודאות הסיווג

באלגוריתמים שנציע אנו נשתמש באלגוריתם למידה בסיסי בקבלת ההחלטות בכל שלב של האלגוריתם בזמן הסיווג.

הסיווג הסופי שלנו מתבצע על תת קבוצה $F' \subset F$ ולכן בכל שלב בדרך נוכל להשתמש בקבוצת התכונות \bar{F} הקיימת לנו בשלב זה.

נגדיר מדד הנקרא ודאות לסיווג בינארי במצב זה המנסה לחשב את ודאות סט הפיצ'רים בסיווג:

def GetCertainty (sample, features, new_feature, learning_algo):

```
learning_algo.predict(sample [features + new_feature])
yes_prob, no_prob = learning_algo.predict_proba(sample)
certainty = entropy(yes_prob, no_prob)
return 1-certainty
```

כאשר המטרה שלנו היא למקסם את פלט הפונקציה.

נסתכל על שני פרמטרים –

(c) תכונה "טובה" היא תכונה זולה.

(d) תכונה "טובה" היא תכונה המשפרת את ודאות הסיווג.

נגדיר את הפונקציה הבאה:

def ScoreFunctionB(features, new_feature):

```
return  $\frac{d(\text{new\_feature}, \text{features})}{c(\text{new\_feature})}$ 
```

בשלב הניסויים, נבחן איזו פונקציה מטיבה עם האלגוריתמים שפיתחנו.

- נרצה להכליל את הגדרת ה-score עבור קבוצת תכונות נתונה:

def ScoreFunction (features):

return $\sum_{f \in \text{features}}$ ScoreFunction(features \ {f}, f)

כאשר ScoreFunction היא פונקציית score כלשהי.

במילים אחרות, על מנת לדעת את גודל ה-score עבור קבוצת תכונות נתונה, נסכום את ערכי ה-score על פני כל התכונות בקבוצה- לכל תכונה, נחשב את ערך ה-score של קבוצת התכונות ללא התכונה, יחד עם התכונה.

אלגוריתם חיפוש בגרף

הקדמה

נרצה להציע אלגוריתם מתוחכם שפותר את הבעיה שהצגנו. לשם כך, ניעזר באלגוריתמי חיפוש בגרפים שלמדנו בחלקו הראשון של הקורס.

נרצה להסתכל על הבעיה המוצגת כבעיית חיפוש בגרף. לשם כך, נגדיר גרף מכון שצומת בו הוא קבוצת תכונות אפשריות. הקשתות היוצאות מכל צומת הם הצמתים המייצגים קבוצות תכונות השונות מקבוצת התכונות שצומת האב מייצג בתכונה אחת בלבד, כך שדרגת היציאה של כל צומת כנ"ל שווה למספר התכונות שלא חלק מקבוצת התכונות שהצומת מייצג. משקל הקשתות הוא מחיר התכונה המבדילה בין הצמתים וההיוריסטיקה על כל צומת היא לפי פונקציית ה-score שהגדרנו בפרקים הקודמים. צומת המטרה יהיה הצומת הכולל את מספרי כל התכונות האפשריות.

הגדרה פורמלית של מרחב החיפוש

ראשית, לצורך פשטות, נקבע סימונים:

- K - קבוצת התכונות הנתונות לאלגוריתם (ניתן כקלט לאלגוריתם בשלב ה-predict).
- D - מספר התכונות האפשריות לכל דוגמה (מימד המרחב הוקטורי שכל הדוגמאות נלקחות ממנו).
- C - וקטור בגודל D אשר בכל כניסה $0 \leq i < D$ שלו מחיר התכונה ה- i . נסמן את מחיר התכונה ה- i כ- $C[i]$.
- \mathcal{P} - סימון לקבוצת חזקה, למשל $\mathcal{P}(A) = \{B \mid B \subseteq A\}$, כלומר $\mathcal{P}(A)$ היא קבוצת כל תתי הקבוצות של A .
- $[D]$ - סימון לקבוצת כל המספרים הטבעיים מ-0 ועד $D-1$: $[D] = \{n \in \mathbb{N} \mid n < D\}$ (D מספר טבעי).
- \setminus - סימון "להפרש קבוצות", למשל $A \setminus B = \{x \mid x \in A, x \notin B\}$, כלומר $A \setminus B$ היא הקבוצה שמכילה את כל איברי A שלא שייכים ל- B .
- $| \cdot |$ - סימון לגודל קבוצה, כלומר $|A|$ הוא גודל הקבוצה A .

נגדיר כעת את מרחב החיפוש $\mathcal{S} = (S, O, I, G)$, כאשר-

- **קבוצת המצבים:**

קבוצת המצבים היא קבוצת כל קבוצות התכונות האפשריות שניתן להוסיף לקבוצת התכונות הנתונות.

$$S = \{P \cup K \mid P \in \mathcal{P}([D] \setminus K)\}$$

- **קבוצת האופרטורים:**

ניתן לעבור ממצב אחד לעוקבו בתנאי שקבוצת התכונות המיוצגת ע"י המצב הראשון וקבוצת התכונות המיוצגת ע"י המצב השני נבדלים זה מזה בתכונה אחת בלבד.

$$O = \{o^{s_1, s_2} \mid s_1, s_2 \in S \wedge |s_2 \setminus s_1| = 1\}$$

- **המצב ההתחלתי:**

המצב הכולל את קבוצת התכונות הנתונות בלבד.

$$I = K$$

- **קבוצת המצבים הסופיים:**

ישנו מצב סופי אחד- המצב הכולל את קבוצת כל התכונות האפשריות (מימד המרחב הוקטורי שכל הדוגמאות נלקחות ממנו).

$$G = \{[D]\}$$

לאחר שהגדרנו את מרחב החיפוש, נגדיר בנוסף:

- **עלות אופרטור:**

נגדיר את פונקציית העלות $c(s_1, s_2)$ עבור מעבר ממצב אחד s_1 למצב עוקב שלו s_2 כמחיר התכונה הנוספת, כלומר מחיר התכונה שנמצאת בהפרש של הקבוצות המיוצגות ע"י המצבים s_1, s_2 .

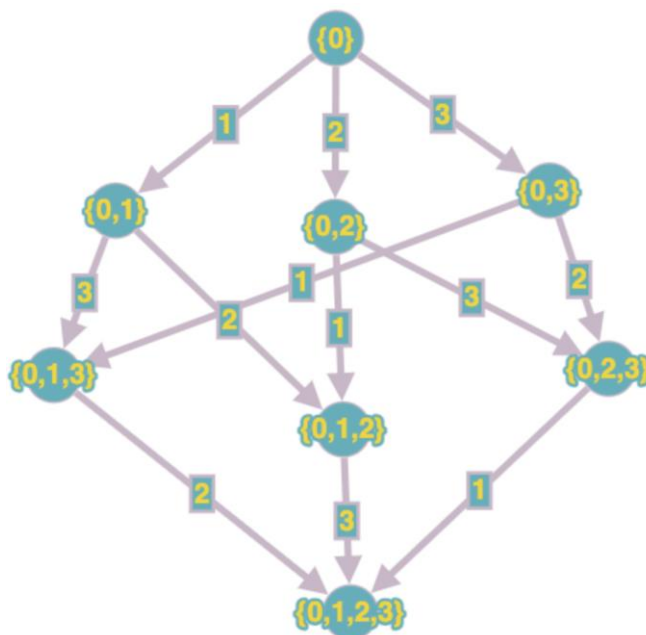
$$c(s_1, s_2) = C[s_2 \setminus s_1] \quad , c: S \times S \rightarrow \mathbb{R}$$

- **היוריסטיקה:**

נגדיר היוריסטיקה Sc עבור המצב s_1 כערך אחת מפונקציות ה-score עבור הצומת $Sc(s_1)$ כאשר את פונקציות ה-score הגדרנו בפרקים הקודמים.

דוגמה

נניח שקבוצת התכונות האפשריות היא $D = 4$, $[D] = \{0, 1, 2, 3\}$.
כמו כן, נניח שקבוצת התכונות הנתונות בתחילת האלגוריתם היא $K = \{0\}$ וקבוצת מחירי התכונות:
 $C = \{1, 2, 3\}$ (מחיר תכונה 1 הוא 1, מחיר תכונה 2 הוא 2 ומחיר תכונה 3 הוא 3).
הגרף המתקבל:



קיבלנו מצומת $\{0\}$ שלוש קשתות לצמתים $\{0,1\}$, $\{0,2\}$, $\{0,3\}$ וזאת מכיוון שכל אחד מהצמתים הללו נבדל מהצומת $\{0\}$ בתכונה אחת בלבד. נשים לב שעל הקשתות המחברות בין הצמתים משקל השווה למחיר הוספת התכונה המבדילה בין הצמתים.

הגדרת האלגוריתם

נציג כעת את האלגוריתם בצורה פורמלית. לצורך כך נבחר אלגוריתם למידה כלשהו (כדוגמת KNN או $ID3$) ואלגוריתם לחיפוש בגרף (כדוגמת BFS , DFS , A^* וכו').

אימון:

• קלט:

(1) קבוצה $\{X, Y\}$ כאשר $X = \{\bar{x}_1, \dots, \bar{x}_m \mid \bar{x}_i \in \mathbb{R}^D\}$ קבוצת דוגמאות (D מספר התכונות)

ו- $Y = \{\bar{y}_1, \dots, \bar{y}_m \mid \bar{y}_i \in \{0,1\}\}$

(2) קבוצת מחירים לתכונות $C = \{\bar{c}_1, \dots, \bar{c}_D \mid \bar{c}_i \in \mathbb{Z}\}$ בהתאמה.

1. שמור את קבוצת הדוגמאות והתיוגים $\{X, Y\}$.

2. שמור את קבוצת המחירים לתכונות C .

סיווג:

• קלט:

(1) דוגמה $\hat{x} \in \mathbb{R}^k$ כאשר $k < D$ ו- D מספר התכונות האפשרי לכל דוגמה.

(2) תקציב מקסימלי \tilde{C} .

(3) קבוצת אינדקסים K של התכונות הנתונות.

• פלט:

(1) סיווג הדוגמה.

1. בנה גרף מכון $G = (V, E)$ המתאים למרחב החיפוש שהגדרנו ($V = S$, $E = O$). על קשתות הגרף E

הוסף משקל לפי פונקציית העלות c כפי שהגדרנו.

2. הרץ את האלגוריתם לחיפוש בגרף שבחרת על הגרף G . במידת הצורך, השתמש בהיוריסטיקה Sc כפי

שהגדרנו. נסמן את המסלול הקצר ביותר שקיבלנו ב- P (וקטור צמתיים המסודרים לפי סדר המסלול הקצר ביותר שמצא האלגוריתם).

3. עבור על הצמתיים ב- P לפי הסדר ועבור כל זוג צמתיים עוקבים $p_i, p_f \in P$:

3.1. חשב את התכונה שנוספה בין הצמתיים: $p_f \setminus p_i \rightarrow f$.

3.2. החסר את מחיר התכונה שנוספה לצומת מהמחיר המקסימלי: $\tilde{C} - C[f] \rightarrow \tilde{C}$.

3.3. אם המחיר המקסימלי $\tilde{C} \leq 0$:

3.3.1. הוסף את התכונה שנוספה לקבוצת התכונות הנתונות ("רכוש את התכונה"):

$$K \cup \{f\} \rightarrow K$$

3.4. אחרת:

3.4.1. עצור את המעבר על הצמתים ב- P .

4. הרץ את שלב האימון של אלגוריתם הלמידה שבחרת (כאשר התכונות הנתונות לכל דוגמה הן התכונות בקבוצת התכונות הנתונות שחישבנו בשלב 3) עבור קבוצת הדוגמאות והתיוגים ששמרת בשלב האימון.

5. הרץ את שלב הסיווג של אלגוריתם הלמידה שבחרת.

6. החזר את הפלט שהתקבל.

מימשנו את האלגוריתם בקובץ graph_search_algorithm.py.

יתרונות וחסרונות האלגוריתם

ניתן לראות שבעת שימוש באלגוריתם ה"ל", אנו משתמשים בכל התקציב שברשותנו על מנת "לרכוש" תכונות ולהוסיף מידע. למעשה, על מנת "לרכוש" את התכונות "הטובות ביותר", אנו נעזרים בהיוריסטיקה המשכללת את הקורלציה שבין התכונות והמחיר שלהם. כך אנו מתחשבים ב-Trade-off בין אופטימליות מספר התכונות לתרומתן- אנו בוחרים את התכונות המוסיפות את המידע הרב ביותר, כל עוד יש לנו מספיק תקציב לכך. יתרה מכך, אנו מוצאים מסלול מוגדר המספק לנו דרך להוסיף תכונות עד כמה שניתן, ללא כל תלות בתקציב המקסימלי שאנו מקבלים כקלט- כך נוכל למצוא את המסלול הטוב ביותר לכל תקציב מקסימלי שנקבל. כך אנו למעשה מתעלמים מהצורך "לנבא" את התקציב המקסימלי שנקבל ומוצאים פתרון המתאים לכל תקציב מקסימלי.

חסרונות של האלגוריתם הם הזמן והמקום, שכן כעת, מלבד הרצת אלגוריתם למידה כלשהו (שעלול לדרוש הקצאה גדולה של זמן ומקום), אנו בונים גרף ומריצים עליו חיפוש (מה שמגדיל את הזמן והמקום הנדרשים). על מנת להתמודד עם החסרונות הללו, נבצע ניסויים במטרה להחליט על אלגוריתם החיפוש וסוג ההיוריסטיקה (פונקציית ה-score) שיביא אותנו לתוצאות הטובות ביותר בזמן סביר.

בנוסף, נשים לב שייתכן מצב שבו על אף שנוותר לנו די תקציב "לרכוש" תכונה נוספת, אנו לא עושים זאת מכיוון שהתכונה הבאה שהאלגוריתם ה"ל" יציע להוסיף (הצומת הבא במסלול הקצר ביותר) יקר יותר מהתקציב שנוותר לנו. למעשה, אנו עלולים לוותר על הוספת מידע כלשהו (גם מידע קטן) מכיוון שאנו לא יכולים להוסיף את התכונה המוסיפה את המידע הרב ביותר באותה נקודת זמן. למעשה, נוכל בעתיד אף לבצע שיפורים באלגוריתם- כמו למשל מציאת דרך לגיזום מוקדם של מסלולים במרחב

החיפוש (למשל, גיזום מסלול שמחירו עובר את התקציב המקסימלי הנתון וכ') או מציאת דרך לנצל את כל התקציב שברשותנו (יצירת מספר עצים כתלות בתקציב הנוטר לנו לאחר כל שלב או מציאת היוריסטיקה מתאימה).

שיפור האלגוריתם: אלגוריתם לחיפוש מקומי

נוכל להציע שיפור לאלגוריתם הנ"ל שיתמודד עם החסרונות שהזכרנו. כעת, במקום לבנות גרף מלא כנזכר לעיל, נשתמש במרחב החיפוש שהגדרנו לצורך פתרון בעיית אופטימיזציה.

נשתמש בסימונים שהגדרנו על מנת להגדיר בצורה פורמלית את הבעיה. ראשית, נרחיב את הגדרת ההיוריסטיקה. לכל מצב $s \in S$ נגדיר ערך היוריסטי:

$$H(s) = \begin{cases} Sc(s), & \text{if } \sum_{f \in s} c[f] \leq \tilde{c} \\ -\infty, & \text{else} \end{cases}$$

כלומר, הערך ההיוריסטי של כל מצב $s \in S$ יוגדר להיות ערך פונקציית ה-score שבחרנו אם המחיר הכולל של כל התכונות במצב s הוא לכל היותר התקציב המקסימלי המתקבל כקלט בשלב ה-predict, אחרת $-\infty$.

כעת, בעיית האופטימיזציה שלפנינו הינה-

$$\max_{s \in S} H(s)$$

במילים אחרות, נחפש את המצב $s \in S$ עבורו הערך ההיוריסטי מקסימלי.

נפתור את בעיית האופטימיזציה ע"י הרצת אלגוריתם לחיפוש היוריסטי מקומי, כפי שלמדנו בקורס (כדוגמת Steepest Ascent hill-climbing, Stochastic hill-climbing, beam search וכ').

נציג את האלגוריתם המשופר בצורה פורמלית. לצורך כך נבחר אלגוריתם למידה כלשהו (כדוגמת KNN או $ID3$) ואלגוריתם לחיפוש מקומי.

אימון:

• קלט:

(1) קבוצה $\{X, Y\}$ כאשר $X = \{\bar{x}_1, \dots, \bar{x}_m \mid \bar{x}_i \in \mathbb{R}^D\}$ קבוצת דוגמאות (D מספר התכונות)

ו- $Y = \{\bar{y}_1, \dots, \bar{y}_m \mid \bar{y}_i \in \{0,1\}\}$

(2) קבוצת מחירים לתכונות $C = \{\bar{c}_1, \dots, \bar{c}_D \mid \bar{c}_i \in \mathbb{Z}\}$ בהתאמה.

1. שמור את קבוצת הדוגמאות והתיוגים $\{X, Y\}$.

2. שמור את קבוצת המחירים לתכונות C .

סיווג:

• קלט:

(1) דוגמה $\hat{x} \in \mathbb{R}^k$ כאשר $k < D$ ו- D מספר התכונות האפשרי לכל דוגמה.

(2) תקציב מקסימלי \tilde{C} .

(3) קבוצת אינדקסים K של התכונות הנתונות.

• פלט:

(1) סיווג הדוגמה.

1. הרץ את האלגוריתם לחיפוש מקומי שבחרת החל מהמצב ההתחלתי K ומצא את המצב s^* עבורו הערך

היוריסטי $H(s^*)$ מקסימלי.

2. עבור על הצמתים ב- s^* ועבור כל צומת ב- s^* שאינו ב- K , הוסף את התכונה שנוספה לקבוצת התכונות

הנתונות ("רכוש את התכונות"): $K \cup s^* \rightarrow K$.

3. הרץ את שלב האימון של אלגוריתם הלמידה שבחרת (כאשר התכונות הנתונות לכל דוגמה הן התכונות

בקבוצת התכונות הנתונות שחישבנו בשלב 2) עבור קבוצת הדוגמאות והתיוגים ששמרת בשלב האימון.

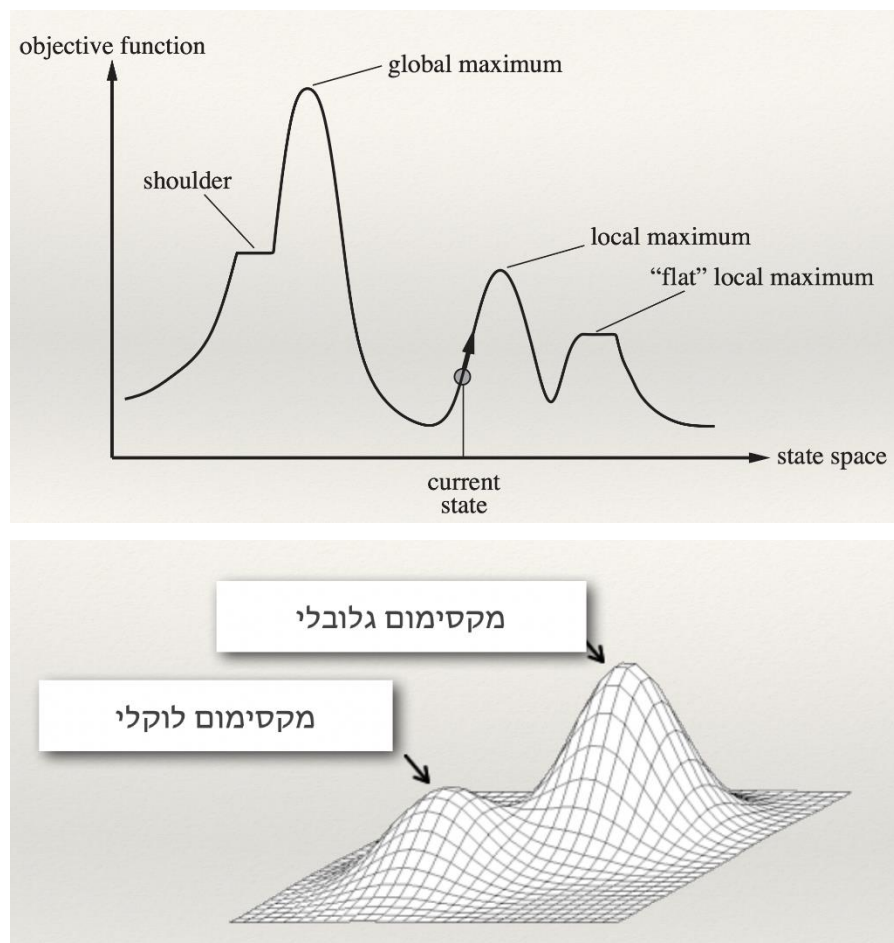
4. הרץ את שלב הסיווג של אלגוריתם הלמידה שבחרת.

5. החזר את הפלט שהתקבל.

מימשנו את האלגוריתם בקובץ `graph_search_algorithm.py`.

ניתן לראות שהאלגוריתם הנ"ל מתמודד עם החסרונות של אלגוריתם החיפוש בגרף- כעת אנו לא מפתחים את כל הגרף, אלא לכל היותר מצב אחד וכל שכניו ובכך חוסכים בזיכרון. בנוסף, אנו מנצלים את כל התקציב שברשותנו ומוסיפים כמה שיותר מידע (תכונות), שכן מאופן הגדרת ההיוריסטיקה, האלגוריתם הנ"ל תמיד ימצא מצב עם מספר אופטימלי של תכונות.

עם זאת, לאלגוריתם הנ"ל חסרונות אחרים. פונקציית ה-score הקובעת את ההיוריסטיקה איננה בהכרח קעורה- כלומר, האלגוריתם לחיפוש מקומי לא בהכרח ימצא מקסימום (או לחילופין, עלול להיתקע במקסימום מקומי ולא גלובלי)- לכן ביצועי האלגוריתם עלולים להיפגע (כמובן שנוכל להשתמש באלגוריתמי חיפוש מקומי מתוחכמים יותר, אלא שבעיית הקעירות תישאר קיימת).



אלגוריתם חיפוש בהשראת אלגוריתמים גנטיים

הקדמה

נרצה להציע אלגוריתם נוסף שבבסיסו עומד אלגוריתם גנטי. לשם כך נרצה להרחיב קצת בנושא של אלגוריתמים גנטיים. אלגוריתמים גנטיים מבצעים תהליכי ברירה מלאכותית המבוססים על עקרונות ביולוגים בהשראת האבולוציה.

רגע של ביולוגיה – צ'ארלס דרווין התחיל את פיתוח התורה האבולוציונית בבסיסה עומדים עקרונות מנגנוני ההורשה והגנטיקה. העקרון המנחה הוא תהליך הברירה הטבעית. תהליך הברירה הטבעית אומר כי תכונות תורשתיות של יצורים נעשות נפוצות יותר ויותר מדור לדור ככל שהן תורמות יותר לרביית הפרט ויכולת הישרדותו. באוכלוסיית יצורים, למי שמעמיד הרבה צאצאים יש סיכוי גדול יותר להוריש להם את התכונות התורשתיות שלו.

העקרון באלגוריתמים גנטיים אומר שמבין מרחב הפתרונות, בדיוק כמו בטבע, רק "המתאימים שורדים". נוכל לשלב אלמנטים של פתרונות אפשריים לבעיה תוך הפעלת הליכים של ברירה מלאכותית כדי "לבחור" מועמדים לדור הבא. ההצלחה של רעיון זה נובעת מההצלחה של האבולוציה בפתרון בעיות אמיתיות. בפשטות, ניקח פתרונות אפשריים לבעיה, נבחר פתרונות שנגדיר כ"מתאימים", נערבב אותם ונוסיף רעש נקבל דור חדש של פתרונות הקרוב צעד אחד יותר לפתרון הבעיה.

הגדרות

- גנום (*genome*) – פתרון אפשרי לבעיה. נהוג לייצג כווקטור בינארי, כאשר האינדקס ה- i הוא 1 אם ורק אם בפתרון יש שימוש בתכונה ה- i .
- אוכלוסייה (*population*) – אוסף של גנומים. מייצג את אוסף הפתרונות שאנו בוחנים בדור הנוכחי.
- פונקציית כשירות (*Fitness*) – פונקציה שבהינתן גנום מחזירה את "טיב" הגנום.
- תהליך זיווג (*crossover*) – לקיחת שני גנומים מהאוכלוסייה ויצירת גנום חדש. הגנום החדש מכיל מידע גנטי של שני הוריו (בדרך כל מכיל את הגנים המובחרים של הוריו).
- יצירת מוטציות (*mutation*) – שינוי מספר תכונות של גנום יחיד כדי למנוע חזרה של חלקים שלמים במחרוזת הביטים. מתבצעת באחוזים נמוכים (בהסתברות נמוכה) משום שאם הייתה מתבצעת בסבירות גבוהה, האלגוריתם היה הופך להיות אקראי.

- שיערוך (evaluation) – בחירת הגנומים שישתתפו בזיווג לפי פונקציית הכשירות.

באופן פורמלי, אם $f: \mathbb{R}^n \rightarrow \mathbb{R}$ פונקציות הכשירות אותה אנו רוצים למקסם :

- $\bar{x} = (x_1, \dots, x_n)$ הוא גנום במודל ושייך לתחום של פונקציית הכשירות.
- x_1, \dots, x_n הם כרומוזומים.
- $P_i = \{\bar{x}_1, \dots, \bar{x}_k\}$ היא האוכלוסיה בדור i .
- $f(\bar{x}) \in \mathbb{R}$ הוא השיערוך של הגנום \bar{x} .
- בהינתן \bar{x}_1, \bar{x}_2 , הזיווג שלהם הוא גנום \bar{x}_3 המורכב מהכרומוזומים של הוריו.
- מוטציה של \bar{x} הוא שינוי בהסתברות קטנה כלשהי p של חלק מהכרומוזומים של \bar{x} .

נקביל את המונחים ההגדרות שהצגנו לבעיה שלנו:

- גנום (genome) – וקטור בינארי, כאשר האינדקס ה- i הוא 1 אם ורק אם בפתרון יש שימוש ב-feature ה- i . נציין כי אנו נשתמש אך ורק בפתרונות חוקיים לבעיה (עלות הפתרון \geq העלות המקסימלית).
- אוכלוסייה (population) – אוסף של גנומים. מייצג את אוסף הפתרונות שאנו בוחנים בדור הנוכחי. כל הפתרונות חוקיים.
- שיערוך (evaluation) – נשתמש במסווג בסיס מסוג ריגרסור לינארי. נייצר מהדוגמאות שניתנו לנו בשלב ה-train סט אימון וסט ולידציה, וכל פעם נאמן מסווג עם סט האימון עבור התכונות שהפתרון הנוכחי מייצג, ונעריך את איכותו לפי ה-accuracy שיקבל על סט הולידציה.
- נציין שבספרות שימוש באלגוריתמים גנטיים הוא דרך פופולרית לביצוע feature selection (באופן כללי זו דרך טובה לפתרון בעיית בחירת תת קבוצה מיטבית).
- בבסיסה של הבעיה אנו מנסים לפתור עומד אלמנט של בחירת תכונות, ולכן יש הגיון בבחירת שילוב אלגוריתם זה בפתרון שאנו מציעים.
- בפתרון שלנו אנו משתמש בספריט deap שמממשת אלגוריתם גנטי (esSimple) שמיועד לפתור את בעיית בחירת תת קבוצה מיטבית.

אימון:

• קלט:

(1) קבוצה $\{X, Y\}$ כאשר $X = \{\bar{x}_1, \dots, \bar{x}_m \mid \bar{x}_i \in \mathbb{R}^D\}$ קבוצת דוגמאות (D מספר התכונות)

ו- $Y = \{\bar{y}_1, \dots, \bar{y}_m \mid \bar{y}_i \in \{0,1\}\}$

(2) קבוצת מחירים לתכונות $C = \{\bar{c}_1, \dots, \bar{c}_D \mid \bar{c}_i \in \mathbb{Z}\}$ בהתאמה.

1. שמור את קבוצת הדוגמאות והתייגים $\{X, Y\}$.

2. שמור את קבוצת המחירים לתכונות C .

סיווג:

• קלט:

(1) דוגמה $\hat{x} \in \mathbb{R}^k$ כאשר $k < D$ ו- D מספר התכונות האפשרי לכל דוגמה.

(2) תקציב מקסימלי \tilde{C} .

(3) קבוצת אינדקסים K של התכונות הנתונות.

• פלט:

(1) סיווג הדוגמה.

1. הרץ אלגוריתם גנטי למציאת תת קבוצה של קבוצת התכונות ($HOF - hall of fame$).

1.1. האוכלוסייה של האלגוריתם תאותחל עם פתרונות חוקיים בלבד.

1.2. בפונקציית הכשירות יפסלו פתרונות חדשים שאינם חוקיים (יוכלים להגיע מהתליך הזיווג).

1.3. פונקציית הכשירות תעריך את הדיוק של סיווג באמצעות התכונות הנוכחיות בלבד.

1.4. בכל פתרון אפשרי יהיה שימוש בכל התכונות הניתנות (בהסתברות גבוהה מאוד) מאופן איתחול

האוכלוסייה וביצוע זיווג בין פתרונות חוקיים (בהסתברות כי יש מוטציות).

2. מהפתרונות החוקיים ב- HOF צור תת קבוצה של כל הפתרונות שמוערכים בכשירות הגבוהה ביותר.

3. בחר אחד מהפתרונות באופן רנדומלי.

4. הרץ את שלב הסיווג של אלגוריתם הלמידה שבחרת עם התכונות שנמצאות בפתרון שבחרת.

5. החזר את הפלט שהתקבל.

מימשנו את האלגוריתם בקובץ `genetic_algorithm.py`.

יתרונות וחסרונות האלגוריתם

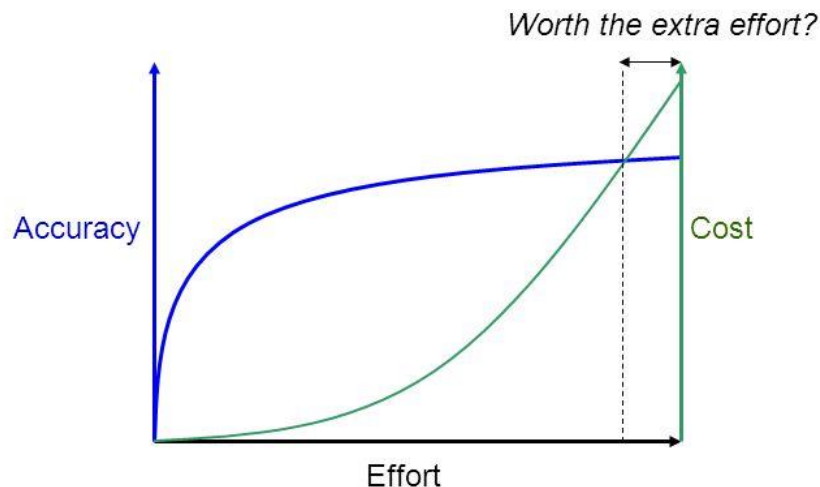
האלגוריתם הנ"ל משתמש בכל התקציב שברשותו על מנת לרכוש תכונות ולהגדיל את המידע שברשותו. האלגוריתם משתמש בפתרון בעיה כללית יותר – מציאת תת קבוצה מיטבית כקופסה שחורה. תת הקבוצה המיטבית מוגדרת לפי פונקציית המחיר שהגדרנו.

חסרון בולט מאוד של שיטה זו הוא זמן הריצה (ינותח באופן מלא בשלב הניסויים). בכל הרצה, אנו מפתחים 10 דורות לאוכלוסייה, ובוחנים פתרונות מתוך קבוצה בגודל $|\#features|^2$, כלומר הפתרון הוא מעריכי בגודל הקלט, ורחוק מלהיות אופטימלי במונחים של זמן ריצה.

אנו מצפים שלא אלגוריתם זה יהיו ביצועים מצוינים, ולעומת האלגוריתם האחרים, שזמן ריצתו יהיה משמעותית גדול יותר.

בכך אלגוריתם זה מהווה דוגמה ל- accuracy effort tradeoff . מונח זה מוכר מאוד בתחום ובעצם מדגים את ה-tradeoff ההגיוני בין יעילות האלגוריתם לבין "ההשקעה" באלגוריתם, שבמקרה הזה באה לידי ביטוי בזמן ריצה.

The effort-accuracy trade-off



מתודולוגיית הניסויים

בשלב זה נרצה לאמוד את הצלחת הפתרונות שהצענו. אם נחזור לבעיה שממנה התחלנו, בעולם האמיתי, הפתרון יוגדר להיות מוצלח אם הבאים מתקיימים:

- הפתרון שמיש – ניתן לקבל את הכרעת הפתרון בזמן שיוגדר להיות סביר (לדוגמא, אם נחזור לבעיית הרופא – חולה, אנו לא יכולים לתת לאלגוריתם שלנו לרוץ כל זמן שירצה, התשובה צריכה להיות תוך פרק זמן שיהיה רלוונטי לביצוע הבדיקות בזמן).
- הפתרון בעל איכות תוצאה המצדיקה את השימוש בו – אל מול ה-base lines שהגדרנו, האלגוריתמים שלנו משתמשים במשאבים רבים (זיכרון וזמן), ולכן שימוש בהם צריך לתת יעילות שאינה מגיעה בפתרונות הנאיביים. נרצה להצדיק את הטריידאוף שבין "יוקר" הפתרון אל בין איכותו.

נרצה להגדיר את טיב האלגוריתמים שלנו. לשם כך נגדיר פרמטרים לאיכות פתרון אותם נבחן בניסויים:

1. זמן ריצה – נרצה לבחון את זמן הריצה של האלגוריתמים שלנו ביחס ל-base line שיצרנו. נרצה להשוות את שתי הגישות שהצענו במובן של זמני ריצה, וכן לבחון אותן אל מול האלגוריתמים הנאיביים לפתרון הבעיה שהצגנו בהתחלה. נציין כי אנו צופים מראש כי זמני הריצה של הפתרונות שלנו יהיו משמעותית ארוכים יותר מזמני הריצה של הפתרונות הנאיביים.
* באלגוריתם הגנטי אנו בוחנים מצבים חוקיים מתוך קבוצה בגודל $|features|^2$, כלומר זמן הריצה הוא מעריכי בכמות התכונות.
* באלגוריתם החיפוש אנו מפתחים גרף עם דרגת יציאה התחלתית (בשכבה הראשונה) של $|features - given features|$, בשכבה העוקבת דרגת יציאה הקטנה ב-1 וכן הלאה. לכן ככל שכמות התכונות תהיה גדולה יותר, ככה פיתוח הגרף יהיה ארוך יותר.
באלגוריתם זה, צפינו את החסרון הזה ולכן הצענו להשתמש באלגוריתם חיפוש מקומי במקום. נרצה לבחון את שיפורו אל מול אלגוריתם החיפוש המפתח את הגרף כולו.
2. דיוק הפתרון – בסופו של דבר, נרצה שהאלגוריתם שהצענו יגיע לתוצאת דיוק גבוהה. נשים לב שכל אחד מהאלגוריתמים שלנו בנויים על גבי מספר פרמטרים. ראשית נבחן את תרומת כל אחד מהם לדיוק הפתרון.
לאחר מכן נשקלל מכל התוצאות את הקומבינציות הטובות ביותר ונבדוק מה דיוק כל אלגוריתם.

דברים שנרצה לבדוק בחלק זה :

- היפרפרמטרים לפונקציית score. מזכיר שבפונקציית scoreA השתמשנו בסקלר הכופל את אחד הגורמים במשוואה כדי למשקל את החשיבות שלהם. נרצה לבדוק מי מהגורמים הינו בעל השפעה מכרעת יותר בקביעת איכות התכונה, ובשביל זה נבצע כיוון להיפרפרמטר זה. שלב זה יגודר כולידציה.
 - נרצה לבדוק איזו מבין פונקציות score שהצענו מביאה לדיוק גבוה יותר בפתרון. כזכור, פונקציית scoreA באה להגדיר תכונה "טובה" כאחת שקורולטיבית לתגית ולא קורולטיבית לתכונות שכבר קיימות בסל התכונות. מנגד, פונקציית scoreB משתמשת בהגדרה שלנו לוודאות הפתרון. נרצה לראות אילו מבין השניים מצליחה לממש את מטרתה בצורה טובה יותר.
 - האלגוריתמים שלנו משתמשים במסווגי בסיס (classifiers). נרצה לראות האם קיים קשר בין מסווג הבסיס לבין הצלחת אחד האלגוריתמים. אנו צופים כי לא יהיה קיים קשר ישיר, אך במובני זמן ריצה יהיו הבדלים (צופים הצלחה ל-logistic regression).
 - עבור האלגוריתם שמבוסס חיפוש בגרף, נרצה לראות מה אלגוריתם החיפוש שמיטיב עם הפתרון. אנו צופים כי תהיה השפעה בעיקר במונחים של זמני ריצה, ולא רואים סיבה לקשר ישיר לטיב האלגוריתם.
- בסיום חלק זה נשקלל את התוצאות כדי לבצע את הניסוי הסופי – בחינת הפתרונות שהצענו ומתן ניקוד לכל אחד מהם.

אספנו שלושה data sets מעולמות שונים ומקורות שונים. חילקנו כל dataset לשלושה –

- קטן – 80 דוגמאות
- בינוני – 200 דוגמאות
- גדול – 600 דוגמאות

את הניסויים נריץ על שלושת הסטים, ובסופו של דבר נשקלל את התוצאות משלושת הסטים, על מנת לצמצם את התלות בסוג מאגר נתונים ספציפי ומאפיינים שלו (מספר התכונות, מספר דוגמאות, קורולציה קיימת בין עמודות או אי התאמה של עמודות לתגית).

מימשנו את הניסויים בקובץ `General/experiments.py`. על מנת להריץ, יש להריץ את הקובץ `experiments.py` באמצעות ה-termial ע"י הרצת הפקודה `python ./General/experiments.py` או דרך סביבת עבודה IDE כלשהי לבחירתכם. הרצת הניסויים בצורה זו תפלוט את הגרפים המתאימים. כעת נפרט את תוצאות הניסויים, נצרף את הפלטים הרלוונטים וכן ננתח את הממצאים.

ולידציה: היפרפרמטר לפונקציית המחיר A

רקע

ניזכר כי בפרק על פונקציית המחיר הגדרנו את הפונקציה :

def ScoreFunctionA(features, new_feature, γ):

return $\frac{a(\text{new_feature}) + \gamma \cdot b(\text{new_feature}, \text{features})}{c(\text{new_feature})}$

כאשר הפרמטרים הם-

(a) קורולטיביות לתגית.

(b) קורולטיביות לתכונות קיימות.

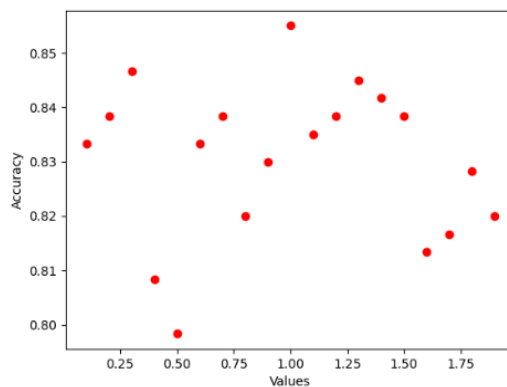
(c) מחיר.

כעת נרצה לכוון את הפרמטר γ שמציין את היחס בין a ל- b.

אופן ביצוע

הרצנו ולידציה באמצעות kfold כפי שלמדנו בקורס על dataset גדול כדי לכוון את ההיפרפרמטר הנ"ל. נבדוק את הקשר של ערך זה לדיוק האלגוריתם.

תוצאות



קיבלנו כי הערך $\gamma = 1$ הינו האופטימלי, והחל מעתה נשתמש בו בכל הניסויים שלנו.

ניסוי: פונקציית המחיר הטובה ביותר

רקע לניסוי

במהלך הפרויקט פיתחנו תיאוריה המגדירה ניקוד לכל תכונה בהינתן תת קבוצה קיימת של תכונות.

המטרה הייתה לבדוק את היעילות לאלגוריתם של הוספת תכונה חדשה לסט התכונות הקיימות.

כדי לממש מטרה זו, היה עלינו להגדיר מה היא תכונה "טובה".

נקטנו בשתי גישות שהובילו לפיתוח שתי פונקציות מחיר – score function שונות.

1. Score function A – פונקציה זו מסתכלת על שלושה מדדים.

a. תכונה היא "טובה" אם היא קורולטיבית לתגית.

b. תכונה היא "טובה" אם היא לא קורולטיבית לתכונות הקיימות.

c. תכונה היא "טובה" אם היא זולה.

2. Score function B – פונקציה זו מסתכלת על שני מדדים.

a. תכונה היא "טובה" אם הוספתה מגדילה את רמת הוודאות (שהוגדרה בפרק זה בדו"ח) של

האלגוריתם.

b. תכונה היא "טובה" אם היא זולה.

כעת, נרצה לבדוק איזו מהפונקציות הנ"ל היא הטובה ביותר. נגדיר את הפונקציה הטובה ביותר במובן של אחוזי דיוק.

אופן ביצוע הניסוי

נבדוק את הפונקציות השונות על ה-dataset (קטן, בינוני וגדול) וכך נשווה את אחוזי הדיוק של כל

האלגוריתמים. שתי הפונקציות ניתנו לאלגוריתם זהה ואותחלו עם פרמטרים זהים.

השערות

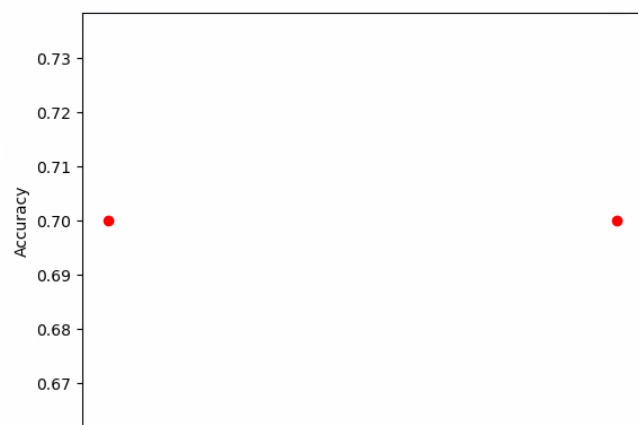
אנו מניחים שבמובנים של הצלחה בניסוי כפי שהוגדרו (אחוזי דיוק), פונקציה B אמורה להיות טובה יותר כי היא נותנת ניקוד לפי ודאות סיווג.

בנוסף, אנו צופים ביצועים שונים לפונקציה A על גדלים שונים של דאטה סטים בשל הרחבת האפשרות

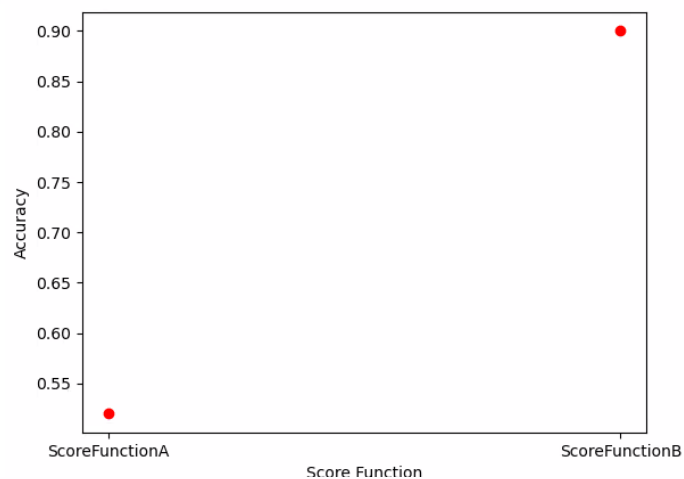
לקורולטיביות להיות משמעותית יותר (ככל שהעמודה גדולה יותר מדד זה מדויק יותר).

תוצאות

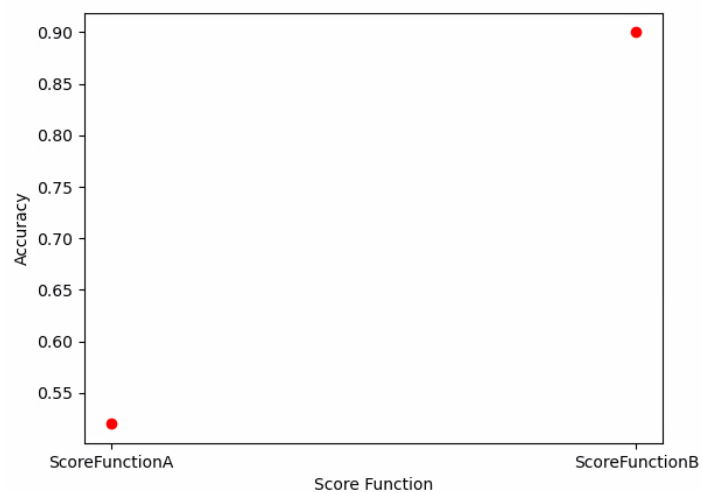
dataset קטן:



dataset בינוני:



dataset גדול:



ניתוח תוצאות

מהתוצאות עולה כי עבור datasets קטנים הביצועים של הפונקציות במובנים של דיוק הינו זהה.

אך כאשר מספר הדוגמאות עולה, פונקציה B מנצחת באופן משמעותי.

מכך אנו מסיקים מספר מסקנות –

- ההגדרה שביצענו עבור ודאות של שלב באלגוריתם הינה הגיונית ועובדת. אכן קיימת התאמה בין כמות הודאות הגדלה באלגוריתם איטרטיבי לבין אחוז הדיוק.
- ב-datasets קטנים פונקציה A משיגה אחוז דיוק דומה לזה של B. אנו יכולים לשער שזה קורה כי עבור עמודות תכונות קטנות יותר, האקספרסיביות של פונקציה B קטנה יותר כי האימון מתבצע על dataset קטן יותר.

ניסוי: מסווג הבסיס הטוב ביותר

רקע לניסוי

בכל האלגוריתמים שפיתחנו אנו עושים שימוש במסווג בסיס. מסווג בינארי זה יכול להיות פרמטר המשפיע על ביצועי האלגוריתם. נרצה לחקור את השפעת מסווג הבסיס על ביצועי האלגוריתם מבחינת דיוק.

אופן ביצוע הניסוי

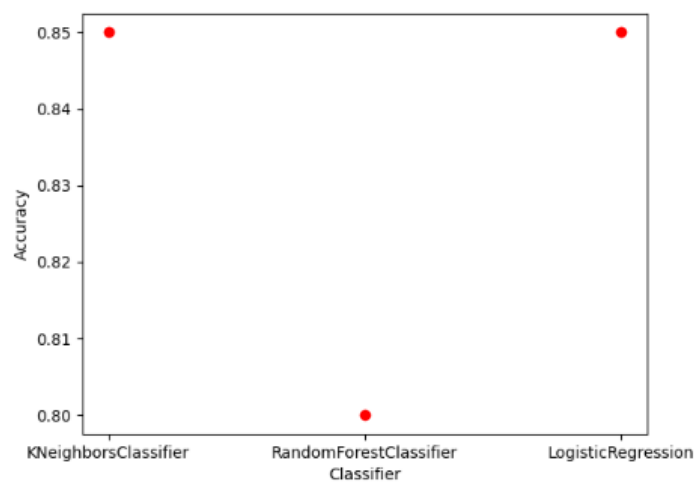
נבדוק את הפונקציות השונות על ה-dataset (קטן, בינוני וגדול) וכך נשווה את אחוזי הדיוק של כל האלגוריתמים. שתי הפונקציות ניתנו לאלגוריתם זהה ואותחלו עם פרמטרים זהים. נציין שניסוי זה התרחש על שלושה datasets שונים ובלתי תלויים והתוצאות הן ממוצע של התוצאות (מכל dataset הפקנו שלושה datasets בגדלים שונים). הסיבה היא שלמעשה שלכל מסווג בסיס יש ביצועים שונים בהתאם להתפזרות ה-data (כפי שראינו בקורס).

השערות

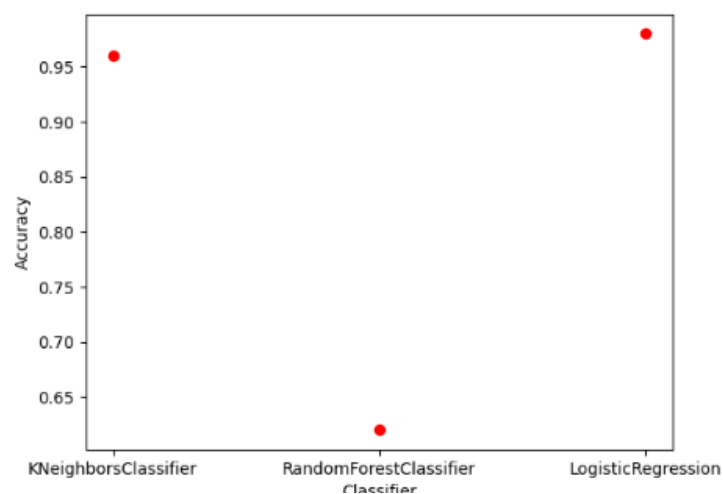
אין לנו השערות ספציפיות לגבי איזה מסווג אמור לתת ביצועים מרשימים באופן מובהק, אך מניסיון העבר שלנו אנו משערים שה-Linear Regressor אמור לבצע עבודה טובה כי הוא פחות מוגבל לעומת מסווגים אחרים שבדקנו (KNN מצפה ל-data שפריד מסוג מסוים, וכן יער שבנוי מעצים המפרידים data לינארית בכל מימד).

תוצאות

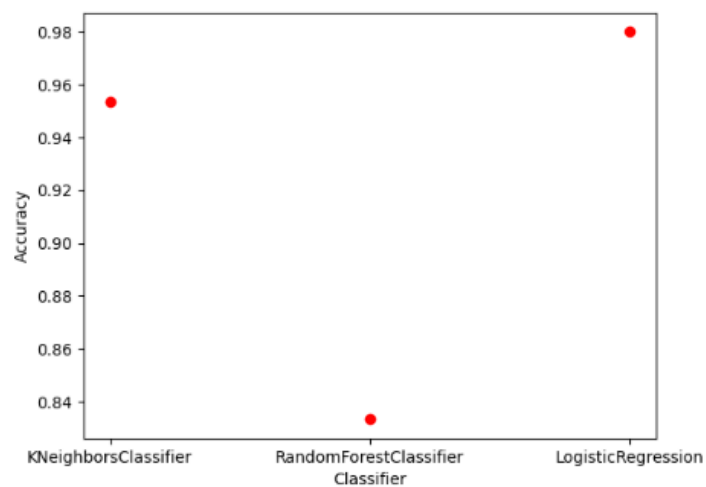
dataset קטן:



dataset בינוני:



dataset גדול:



ניתוח תוצאות

מהתוצאות ניתן להסיק כי כמעט תמיד ה-Random Forest משיג תוצאות פחות טובות משאר המסווגים וכן כי ה-Linear Regressor ביצע על כל סוגי ה-data וכן על כל הגדלים את התוצאות הטובות ביותר. נחזור ונאמר כי אנו מבינים כי ניסוי זה אינו חד משמעי וכי יכול להיות שעבור סוגי data שונים יתאימו מסווגים אחרים בשל תכונות ה-data ואופן התפלגותו. עדיין ראינו לנכון לבצע ניסוי זה על ה-data שהשתמשנו בו בניסויים כדי לדעת באיזה מסווג בסיס בדאי לנו להשתמש.

ניסוי: אלגוריתם הלמידה הטוב ביותר

רקע לניסוי

במהלך הפרויקט הצענו מספר אלגוריתמי למידה תחת אילוצים:

אלגוריתמים נאיביים:

1. "האלגוריתם הריק" - סיווג לפי התכונות הנתונות בלבד.
2. "האלגוריתם הרנדומלי" - בוחר תכונות רנדומליות עד אשר מחיר התכונות הנבחרות עובר את התקציב הנתון.
3. "האלגוריתם האופטימלי" - בוחר תכונות לפי מחירן, מהזולה ביותר ועד היקרה ביותר, עד אשר מחיר התכונות הנבחרות עובר את התקציב הנתון.

אלגוריתם ביניים:

4. "אלגוריתם שונות מקסימלית" - בוחר תכונות לפי גודל השונות שלהן, מהגדול ביותר ועד הקטן ביותר, עד אשר מחיר התכונות הנבחרות עובר את התקציב הנתון.

אלגוריתמים מתוחכמים:

5. "אלגוריתם חיפוש בגרף" - בוחר תכונות לפי בניית גרף והרצת אלגוריתם חיפוש על הגרף.
6. "אלגוריתם חיפוש מקומי" - בוחר תכונות לפי הרצת אלגוריתם חיפוש מקומי על פונקציית ה-score.
7. "האלגוריתם הגנטי" - בוחר תכונות לפי הרצת אלגוריתם גנטי (התכונות המתאימות ביותר "שורדות").

כעת, נרצה לבדוק איזה מהאלגוריתמים הנ"ל הוא הטוב ביותר. נגדיר את האלגוריתם הטוב ביותר במובן של אחוזי דיוק וזמן הריצה.

אופן ביצוע הניסוי

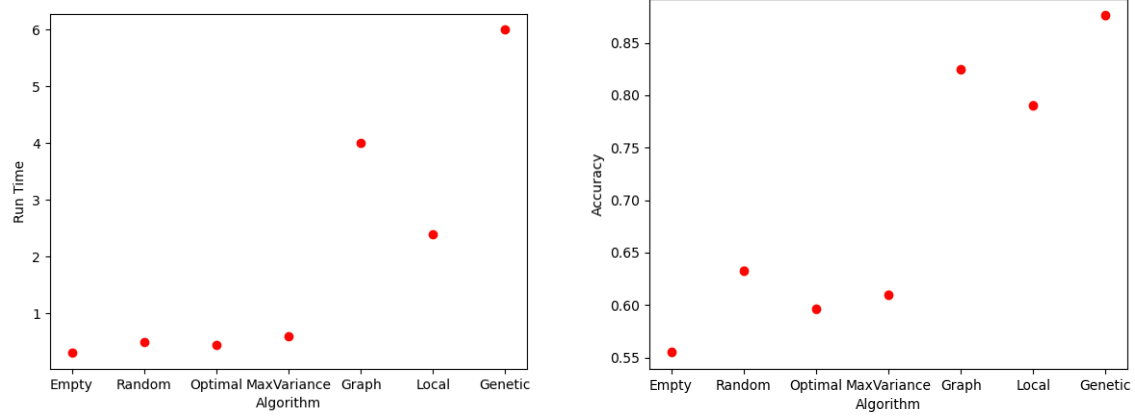
נבדוק את האלגוריתמים השונים על אותו dataset וכך נשווה את אחוזי הדיוק וזמן הריצה של כל האלגוריתמים כתלות בגודל ה-dataset בו אנו משתמשים (קטן, בינוני וגדול).
על מנת לבדוד את המשתנים בניסוי ולהימנע מרעשים והשפעות הנובעים מהפרמטרים שכל אלגוריתם מקבל, נבחר פרמטרים קבועים לכל האלגוריתמים:

- נקבע שהמסווג שכל האלגוריתמים ישתמשו בו הוא Linear Regressor.
- נקבע את התכונות הנתונות, מחירי התכונות והמחיר המקסימלי להיות זהים בין כל האלגוריתמים.
- עבור כל האלגוריתמים המשתמשים בפונקציית score, נבחר כברירת מחדל את ScoreFunctionB.

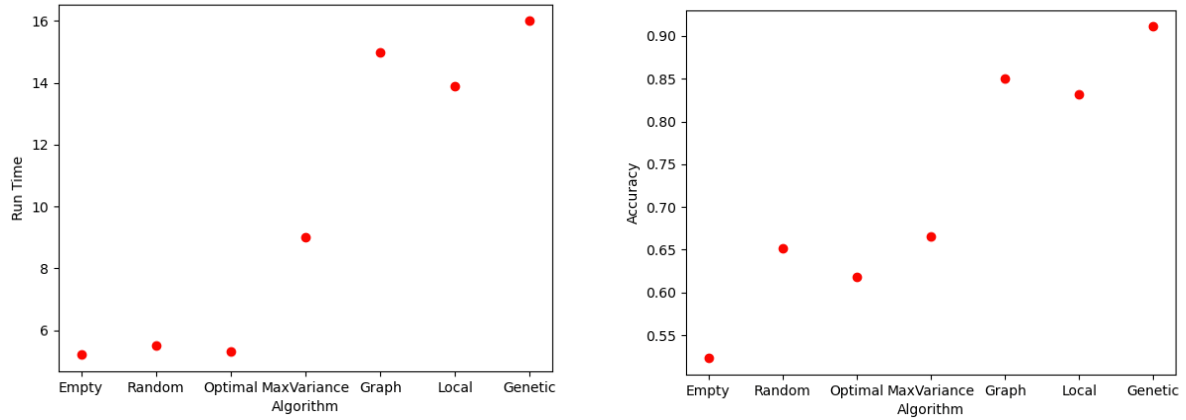
השערות

אנו מצפים שאחוזי הדיוק הטובים ביותר יתקבלו עבור האלגוריתמים המתוחכמים יותר, העושים שימוש בהיוריסטיקות, בכך בוחנים בצורה טובה יותר את הקשר שבין התכונות ומנצלים את כל המידע שביכולתם להשיג. נצפה שהדיוק הנמוך ביותר יתקבל עבור "האלגוריתם הריק", שכן הוא עושה שימוש במספר התכונות הקטן ביותר ובכך מסתפק במידע מאד בסיסי.
לעומת זאת, אנו מצפים שזמן הריצה של האלגוריתמים הנאיביים יהיה קצר משמעותית מזמן הריצה של האלגוריתמים המתוחכמים, שכן האלגוריתמים המתוחכמים כוללים בניית גרף, יצירת מצבים ושימוש בפונקציית ה-score, שכן כל אלו עלולים לבוא לידי ביטוי בזמן ריצה ארוך.
כמו כן, נצפה שזמן הריצה של כל האלגוריתמים יעלה כתלות בגודל ה-dataset, בפרט נצפה להבדלים גדולים בזמן הריצה של האלגוריתמים המתוחכמים במיוחד.

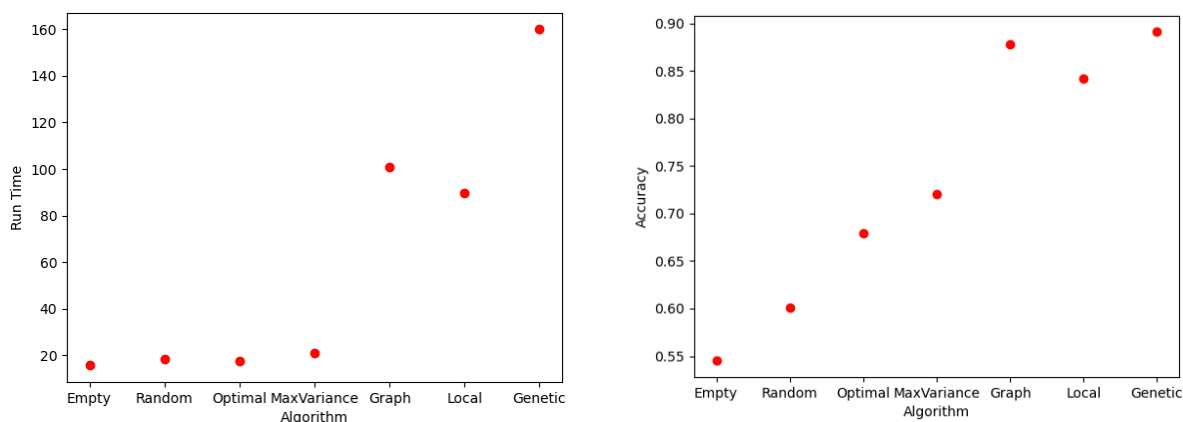
Dataset קטן



Dataset בינוני



Dataset גדול



ניתוח תוצאות

ראשית, נשים לב שמבחינת הדיוק, לא התקבלו הבדלים משמעותיים בהרצה על dataset בגדלים שונים. כצפוי, "האלגוריתם הריק" השיג את אחוזי הדיוק הנמוכים ביותר בכל אחת מההרצות- זאת מכיוון שהוא מסווג לפי מספר התכונות הקטן ביותר מבין האלגוריתמים. כמו כן, ניתן לראות שכל האלגוריתמים הנאיביים השיגו אחוזי דיוק קרובים יחסית זה לזה- "האלגוריתם האופטימלי" אמנם מיצה את התקציב הניתן לו לרכישת תכונות, אלא שלא הצליח למנף זאת- כנראה בשל העובדה שבחר תכונות קורלטיביות זו לזו (אם כי אחוזי הדיוק שלו אכן משתנים בין ההרצות על dataset בגדלים שונים). אחוזי הדיוק של "האלגוריתם הרנדומלי" אכן משתנים במקצת בין dataset שונים, אך נוכל לקשור את העובדה הזאת בעובדה שבחירת התכונות נעשית באופן רנדומלי ולכן ההבדל באחוזי הדיוק טמון בכך ולא בגודלו של ה-dataset.

ניתן לראות שהאלגוריתמים המתוחכמים יותר השיגו אחוזי דיוק גבוהים יותר בהשוואה לאלגוריתמים הנאיביים, כאשר "האלגוריתם הגנטי" בעל אחוזי הדיוק הגבוהים ביותר, בכל אחת מההרצות.

מעניין לראות "שאלגוריתם החיפוש בגרף" השיג אחוזי דיוק גבוהים יותר מזה של "האלגוריתם לחיפוש מקומי" על אף "שהאלגוריתם לחיפוש מקומי" הוא שיפור של "האלגוריתם לחיפוש בגרף". ניתן להסביר תופעה זו בכך "שהאלגוריתם לחיפוש מקומי" נתקע במקסימום מקומי או שפונקציית ה-score איננה קעורה- נרצה להתעמק בתופעה הזאת בהמשך כאשר נריץ ניסוי להשוואת סוגי אלגוריתמים לחיפוש מקומי.

ביחס לזמן הריצה, ניתן לראות שכצפוי- האלגוריתמים הנאיביים היו המהירים ביותר ואף תוצאותיהם קרובות. לעומתם, האלגוריתמים המתוחכמים יותר דרשו זמן ריצה רב יותר, בייחוד האלגוריתם הגנטי שלו זמן הריצה

הגדול ביותר. מעניין להבחין שככל שה-dataset גדול יותר, ההבדלים בין זמני הריצה של האלגוריתמים הנאיביים והמתוחכמים גדולים יותר.

נשים לב שבכל שלושת הגדלים של ה-dataset, האלגוריתם לחיפוש בגרף דרש זמן ריצה גדול יותר מאשר האלגוריתם לחיפוש מקומי, שכן האלגוריתם לחיפוש מקומי לא דורש פיתוח של גרף שלם והרצת אלגוריתם חיפוש עליו- הוא מפתח בכל איטרציה את השכנים ומבצע עליהם את הבדיקות. בכך, נוכל להסיק שהשיפור שלנו עבור האלגוריתם לחיפוש בגרף אכן עובד.

לסיכום, נראה שמבחינת ה-Trade Off בין זמן הריצה לאחוזי הדיוק- אלגוריתם שדרש זמן ריצה גדול יותר הצליח להגיע לאחוזי דיוק טובים יותר.

ניסוי: זמן ריצה כתלות בכמות התכונות

רקע לניסוי

בהמשך לניסוי הקודם, נרצה כעת לבחון את זמן הריצה כתלות בכמות התכונות שב-dataset. במיוחד, נרצה לבחון את ההשפעה על פני האלגוריתמים המתוחכמים שהצענו: Graph Search, Local Search, Genetic.

אופן ביצוע הניסוי

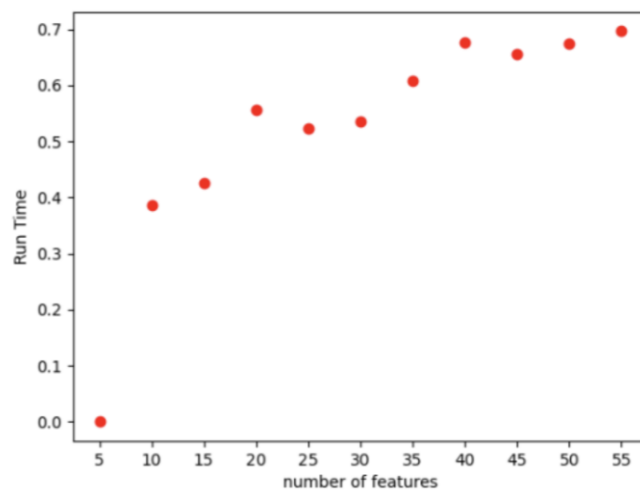
נבדוק את השימוש באלגוריתמים השונים על ה-dataset הקטן כך שבכל פעם נגדיל את מספר התכונות שב-dataset וכך נשווה את זמני הריצה המתקבלים כתלות במספר התכונות (כאמור, את מספר הדוגמאות שב-dataset נשאיר קבוע).
על מנת לבודד את המשתנים בניסוי ולהימנע מרעשים והשפעות הנובעים מהפרמטרים שכל אלגוריתם מקבל, נבחר פרמטרים קבועים לכל האלגוריתמים:

- נקבע שהמסווג שכל האלגוריתמים ישתמשו בו הוא Linear Regressor.
- נקבע את התכונות הנתונות, מחירי התכונות והמחיר המקסימלי להיות זהים בין כל האלגוריתמים.
- עבור כל האלגוריתמים המשתמשים בפונקציית score, נבחר כברירת מחדל את ScoreFunctionB.

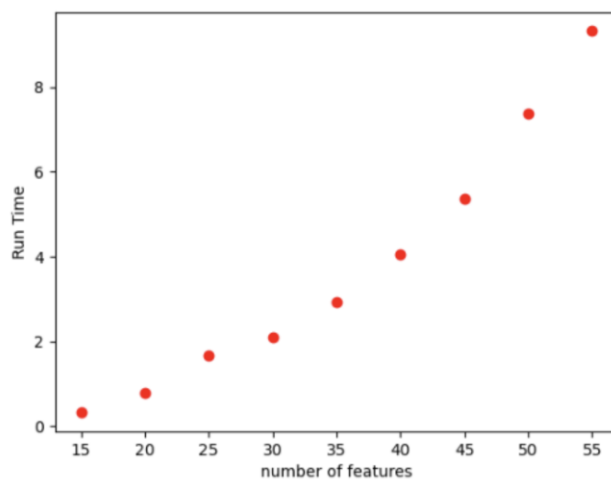
השערות

אנו מצפים שעבור כל האלגוריתמים הנ"ל, נקבל גרף עולה ממש כתלות במספר התכונות.

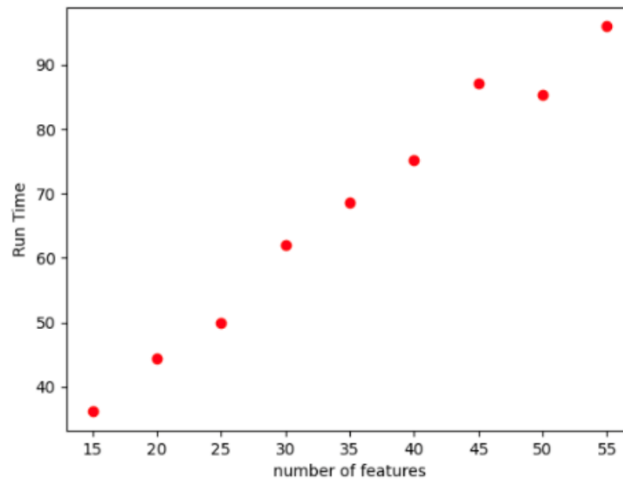
Graph Search Algorithm



Local Search Algorithm



Genetic Algorithm



ניתוח תוצאות

כצפוי, ניתן לראות שזמן הריצה של כל האלגוריתמים עולה כתלות בכמות התכונות. ניתן לראות שאמנם מגמת כל הגרפים דומים, אך עם זאת קיים הבדל אשר מתבטא בגרף של Local Search Algorithm שהינו קעור בעוד הגרפים של שני האלגוריתמים האחרים קמור מכאן שהעלייה בכמות התכונות אמורה להשפיע בצורה דרסטית יותר על האלגוריתם החיפוש המקומי (אם כי האלגוריתם הגנטי הוא בעל זמן הריצה הגדול ביותר וכן בעל השיפוע הגדול ביותר).

ניסוי: אלגוריתם החיפוש המקומי הטוב ביותר

רקע לניסוי

במהלך הפרויקט הצענו אלגוריתם המפתח גרף מכוון ממושקל (כאשר כל צומת בו הוא קבוצת תכונות אפשרית לסיווג על פיהם) ומריץ עליו אלגוריתם חיפוש לגרף על מנת למצוא את המסלול הקצר ביותר בגרף, מצומת ההתחלה לצומת המטרה (צומת הכולל את קבוצת כל התכונות האפשריות). לפי מסלול זה, הבנוי מקשתות בין צמתים שקבוצות התכונות המתאימות להם נבדלות בתכונה אחת בלבד, נוכל לבחור את קבוצת התכונות לפיהם נסווג בכך "שנרכוש" את התכונה הנוספת בין כל זוג צמתים במסלול הקצר ביותר עד אשר התקציב שברשותנו נגמר. לגרף המוזכר לעיל הגדרנו פונקציית score שתשמש אותנו להערכת טיב הצומת. פונקציה זו משמשת אותנו כהיוריסטיקה עבור הגרף.

בהמשך, הצענו שיפור לאלגוריתם הנ"ל שמבצע חיפוש מקומי על מרחב החיפוש במטרה למקסם את פונקציית ה-score שהגדרנו. לאחר שנמצא מצב עבורו ערך ה-score מקסימלי, נוכל להשתמש בקבוצת התכונות שהמצב מתאר על מנת לסווג על פיהם.

כעת, נרצה לבחון איזה מאלגוריתמי החיפוש המקומי שלמדנו בקורס הוא הטוב ביותר, במונחי אחוזי דיוק, לשמש אותנו במלאכת הסיווג של האלגוריתם שהצענו.

בקורס למדנו מספר אלגוריתמים שימושיים לחיפוש בגרפים שאותם נרצה כעת לבחון:

1. hill-climbing - החיפוש מתחיל במצב התחלתי נתון. בכל שלב מפתחים את כל השכנים ובוחרים את השכן הטוב ביותר - זה בעל ערך ה-score הגדול ביותר. אם ישנם כמה שכנים טובים ביותר, האלגוריתם יבחר אקראית. האלגוריתם ייעצר אם הגיע למצב מטרה או אם השכן הטוב ביותר איננו טוב יותר מהמצב הנוכחי.
2. Stochastic hill-climbing - אלגוריתם חמדן שבוחר תמיד את הצעד התלול ביותר, אך מוכן לצעוד בכיווני שיפור שאינם התלולים ביותר. האלגוריתם בוחר אקראית מבין הצעדים המשפרים עם הסתברות פרופורציונית לעוצמת השיפור.
3. Simulated-annealing - אלגוריתם שמתיר צעדים בכיוון מטה (כאלה שמרעים את הערך ההיוריסטי), בכך מונעים היתקעות במקסימום מקומי. ההסתברות לבחירת צעדים כאלו יורדת ככל שמתקדם החיפוש.
4. Random hill-climbing - שיפור של אלגוריתם ה-hill climbing. האלגוריתם מריץ את hill-climbing בצורה איטרטיבית ובכל פעם שהחיפוש נתקע, הוא מתחיל מחדש תוך הגרלת מצב התחלה רנדומלי. האלגוריתם שומר את המקסימום הטוב ביותר שמצא ובכך מונע היתקעות במקסימום מקומי.
5. Local beam search - אלגוריתם המבצע חיפוש אלומה ללא התעניינות במסלול.

אופן ביצוע הניסוי

נבדוק את השימוש באלגוריתמים השונים על אותו dataset וכך נשווה את אחוזי הדיוק המתקבלים בעת השימוש באלגוריתמים השונים.

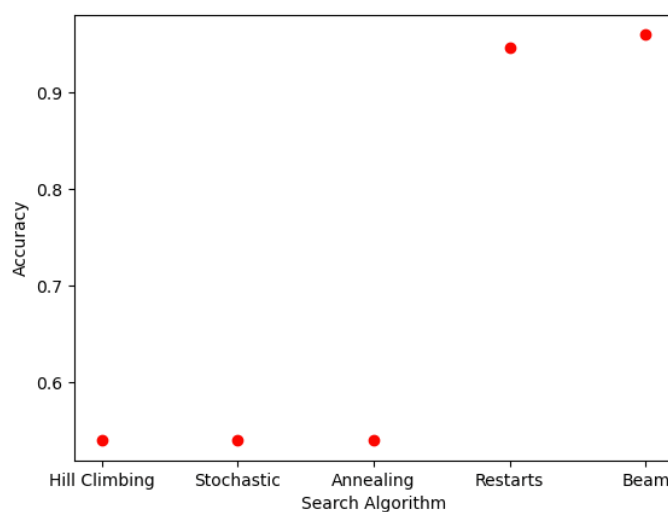
על מנת לבודד את המשתנים בניסוי ולהימנע מרעשים והשפעות הנובעים מהפרמטרים שכל אלגוריתם מקבל, נבחר פרמטרים קבועים לכל האלגוריתמים:

- נקבע שהמסווג שכל האלגוריתמים ישתמשו בו הוא Linear Regressor.
- נקבע שהתכונות הנתונות, מחירי התכונות והמחיר המקסימלי הניתנים לאלגוריתם שהצענו יהיו זהים בכל אחת מההרצות.
- נבחר כברירת מחדל את ScoreFunctionB כפונקציית ה-score בה נשתמש.

השערות

אנו מצפים שאחוזי הדיוק הטובים ביותר יתקבלו עבור האלגוריתמים שמונעים היתקעות במקסימום מקומי, זאת משום שאלגוריתמים כאלו מאפשרים מציאה של מקסימום גלובלי ובכך את קבוצת התכונות הטובה ביותר שבדאי לסווג לפיה.

תוצאות



ניתוח תוצאות

כפי שצפינו, נראה כי האלגוריתמים המונעים היתקעות במקסימום מקומי הניבו תוצאות טובות יותר (מלבד Simulated-annealing שהינו אלגוריתם בעל פן הסתברותי ובשל כך רנדומלי- ניתן לחשוב שהתוצאה הפחות טובה התקבלה בשל הרנדומליות). האלגוריתם המוצלח ביותר לשימוש באלגוריתם שהצענו הינו Local Beam Search שנראה שהגיעו לקרוב ל-98% דיוק, בעוד ש-Random hill-climbing הגיע לתוצאה לא רחוקה ממנו.

סיכום הפרויקט

סיכום

את הרעיון לפרויקט שאבנו מבעיה רפואית מציאותית מחיי היום-יום: חולה מגיע לרופא ומבקש ממנו לאבחן מחלה שקיימת או לא קיימת אצלו תחת מסגרת תקציב לביצוע בדיקות רפואיות. לאחר חשיבה וקריאה נוספת של מידע ומחקרים, החלטנו להרחיב את הפרויקט ולהפכו לפרויקט אלגוריתמי כללי. מטרתנו בפרויקט זה הייתה למצוא אלגוריתם למידה העובד תחת אילוצים של תכונות חסרות (מימד חסר של מרחב התכונות). במהלך הפרויקט הצענו מספר אלגוריתמים שונים, ניתחנו אותם וסיכמנו את יתרונותיהם וחסרונותיהם - כך התקדמנו והצענו אלגוריתמים נוספים במטרה להתמודד עם החסרונות, לשמר את היתרונות ולשפר את האלגוריתמים אחד אחר השני. מאלגוריתמים נאיביים, הגענו לאלגוריתמים מתוחכמים שמבצעים שימוש בכלים שונים שלמדנו במהלך הקורס "מבוא לבינה מלאכותית". לצד הכלים שרכשנו בקורס, ביצענו ניסויים במטרה לאמוד את טיב המסווגים שהצענו ובחנו את השימוש באלגוריתמים השונים שלמדנו על טיב האלגוריתמים שהצענו. כמו כן, במהלך הפרויקט נחשפנו לכלים טכנולוגיים מתקדמים וספריות מיוחדות שעוסקות בלמידה (כפי שניתן לראות בקובץ `environment.yml`), בכך התמקצענו בתחום והתנסינו בכתיבת קוד טכנולוגי ויעיל המשלב מסווגים ולמידה באופן כללי.

דיון בתוצאות

לאחר ביצוע הניסויים וניתוח התוצאות, מצאנו שהאלגוריתם הטוב ביותר לפתרון הבעיה בה התעסקנו, מבחינת אחוזי דיוק, הינו האלגוריתם הגנטי (המתבסס על Linear Regressor). למדנו שהאלגוריתם הגנטי משמש לבחירת תת הקבוצה "החזקה" ביותר (כאשר "חוזק" של קבוצה מתבטא בפונקציית `fitness` של האלגוריתם), מה שקרוב באופן יחסי לבעיה אותה ניסינו לפתור. נראה שלצד אחוזי הדיוק הרבים שהאלגוריתם הניב, זמן הריצה הארוך שלו (בעיקר בשלב הסיווג) יכול להוות עקב אכילס לשימוש בו. אנו סבורים ובטוחים שהתבססות על האלגוריתם הנ"ל לצד שדרוג שלו, יכול להוות פריצת דרך של ממש בתחום הסיווג של אילוצי הגודל של מרחב התכונות. יחד עם זאת, גם האלגוריתמים המבצעים חיפוש במרחב חיפוש הובילו לתוצאות טובות וגם הם יכולים להוות אבן דרך של ממש בתחום. נשים לב שמבין פונקציות ה-`score` שהצענו, פונקציית `scoreB` העושה שימוש בהגדרת הוודאות שפיתחנו הינה העדיפה (אם כי זמן החישוב שלה ארוך יותר).

כיוונים להמשך המחקר

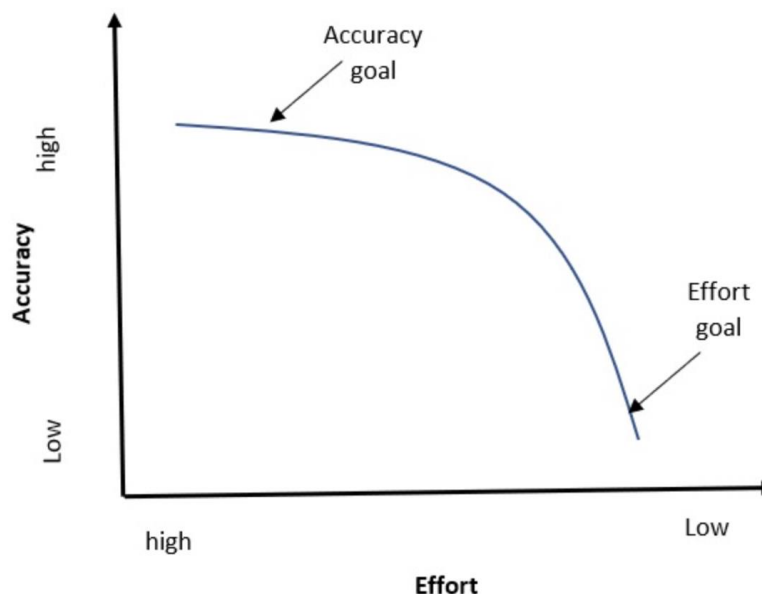
הוא מאפשר לבצע סיווג תחת אילוצי תכונות חסרות, מה שיכול להתברר כחשוב וכבסיס בתחום. בפרט, ישנם אפיקים נוספים שאנו חושבים שיכולים להוביל לתוצאות מעניינות ולתת מענה לבעיות מסוג זה:

1. פיתוח אלגוריתמים חדשים: ניתן להמשיך ולהרחיב את המחקר שביצענו על ידי הצעת אלגוריתמים מתוחכמים נוספים שמתמודדים עם החסרונות באלגוריתמים המתוחכמים השונים שהצענו. למשל, ניתן לפתח אלגוריתמים נוספים שמשתמשים באלגוריתמים אחרים שלא דווקא נלמדים במסגרת הקורס "מבוא לבינה מלאכותית".
לדוגמה, להציע אלגוריתם המתבסס על אלגוריתמי Selector Feature Sequential שהינם משפחה של אלגוריתמים חמדניים המשמשים להורדת מימד מרחב התכונות (כלומר, בהינתן dataset בעל מרחב תכונות מממד d , האלגוריתמים הללו פולטים את מרחב התכונות הטוב ביותר מממד $k < d$). על מנת לפתח אלגוריתם כנ"ל, נדרש ללמוד את אופן פעולת האלגוריתם בהרחבה ולהציע דרך להשתמש בפלט לצד שמירה על מקסימום קורלטיביות בתת קבוצת התכונות המתקבלת כפלט.
בנוסף, ניתן לשלב ולחשוב על אלגוריתם שמתבסס על כלים בלמידה עמוקה (למידה מונחית, בלתי מונחית ולמידת חיזוק) ובפרט ברשת ניורונים. זהו תחום מתפתח שמציע שיפורים ושדרוגים טכנולוגיים רבים וניתן אולי לחשוב על דרך לרתום את התחום הנ"ל לצורך סיווג תחת אילוצי הגודל של מרחב התכונות.
2. שיפור אלגוריתמים קיימים: ניתן לחשוב על דרכים לשפר את הדיוק וזמן הריצה של האלגוריתמים שכבר הצענו. למשל, מציאת דרך לגיזום בעת בניית הגרף באלגוריתם החיפוש בגרף, דבר שיכול להוביל לזמן ריצה קצר יותר. ניתן למשל לגזום מסלולים העוברים בקבוצות תכונות שמחירן הכולל הוא מעל התקציב המותר, או לחילופין לגזום קבוצות תכונות שקורלטיביות זו לזו מעל סף מסוים. ניתן בפרט לחשוב על דרך להפחית את זמן ריצת האלגוריתם הגנטי שסיפק את התוצאה הטובה ביותר.
3. שיפור פונקציית ה-score: אפיק נוסף להמשך מחקר טמון בפונקציית ה-score. על מנת לשפר את אלגוריתמי הגרף והחיפוש המקומי שהצענו, ניתן לשפר בפרט את פונקציית ה-score, כך למשל לחשוב על פונקציית score קעורה בהכרח (לה מובטח מקסימום), פונקציית score שלה בהכרח מספר מועט של מקסימום מקומי או חשיבה על פונקציית score חדשה שלוקחת בחשבון פרמטרים אחרים. ניתן בפרט לשפר את זמן הריצה הדרוש לחישוב ה-score (בפונקציית ה-score השנייה שהצענו, אנו נדרשים לעשות שימוש במסווג קיים מה שעלול להוביל לזמן ריצה ארוך במיוחד).

4. סיווג שאיננו בהכרח בינארי: על מנת לפשט את הפרויקט ולנתח את העומד בבסיסו, ניתחנו בעיות סיווג עבור תיוגים בינאריים (Classification). ניתן כעת לנתח ולבחון את השפעת האלגוריתמים שהצענו על בעיות סיווג שאינן בינאריות, כלומר בעיות סיווג למספר רב של תיוגים.

5. עיבוד מקדים לתכונות: ניתן לחשוב על ביצוע שלב מקדים טרם הרצת האלגוריתם (Preprocessing) ובמהלכו לבצע Data exploration, דבר שיכול להוביל "לזריקת" תכונות לא רלוונטיות (עקב מציאת תכונה לא רלוונטית לתגית המטרה, או תכונה שהינה קומבינציה לינארית של תכונות אחרות) טרם הרצת האלגוריתם- בכך לתרום לזמן הריצה של האלגוריתם (הוכחנו במהלך הניסויים שיש קשר בין כמות התכונות לזמן הריצה).

6. ניסיון לסיווג האלגוריתמים על קשת Accuracy-Effort- Trade Off: אנו במהלך הניסויים התייחסנו לטיב האלגוריתמים במונחי אחוזי דיוק, אלא שאולי נרצה להגדיר את "האלגוריתם הטוב ביותר" כאחוזי דיוק תחת מגבלת זמן ריצה נתון. לכן נרצה לדעת היכן ממוקמים האלגוריתמים שהצענו על העקומה המצורפת ובעתיד לחשוב על שיפורים לאלגוריתמים במונחי Accuracy-Effort ולא רק במונחי אחוזי דיוק.



בנימה אישית

הפרויקט שלנו, שהתחיל מחשיבה על פתרון בעיה רפואית מציאותית מחיי היום-יום, התרחב לפרויקט אלגוריתמי בו ניתן להשתמש בתחומים רבים בחיים במטרה לשפר את החיים ולתרום לאנושות. הפקנו מהפרויקט המון, הן מהפן המחקרי והן מהפן הטכנולוגי. ראינו כיצד הנלמד בקורס "מבוא לבינה מלאכותית" קורם עור וגידים ומשמש אותנו למחקר ופיתוח בחיי היום-יום. אין צל של ספק שהפרויקט בו התנסינו פתח לנו שער למחקר עתידי ועניין רב בתחום, יהווה בסיס להצעות מחקר שנגיש בעתיד ואף ילווה אותנו רחוק בהמשך הקריירה. אנו מרגישים שזכינו לקחת חלק בתחום טכנולוגי מתפתח וכי יש עוד אפיקים רבים למחקר עתידי.