

ממן 21

סעיף ד'

בהמשך לסעיפים א ו-ב ערכו סקירה השוואתית לכלל החלופות האפשריות (לפחות 4 חלופות) לביצוע כריית מידע. בתשובתכם יש להתייחס ליתרונות/חסרונות כל אחת מהחלופות בהקשר לבעיה הנתונה.

יתרונות הסוג	חסרונות הסוג	מדד הפיצול	הסבר כללי	יתרונות	חסרונות	ספציפית בנוגע לבעיה הנתונה
<ul style="list-style-type: none"> יכול לקבל ערכים קטנוניים וזוגיים רציפים. לא דורש נרמול של d_{atan} קל לחסרה לזכרון ולא מהתחום ולכן לדוגמה במקרה שלנו יהיה קל להעביר מידע זה לצבור ולמנות השמנת יתר באוכלוסיה. לא מתחייב שהקשר בין התכונה לסיווג יהיה ליניארי. בניגוד לרגרסיה ליניארית תוצאות חריגות לא משפיעות את התוצאות לחלוטין. 	<ul style="list-style-type: none"> תמיד בביצוע החלטות שלו זמן האימון הנדרש ארוך יחסית לסוגים אחרים בדרך כלל המסקנות דורשות חישוב ארוך יותר בשביל להגיע אליהם יש נטייה ל $overfitting$ בניגוד לאלגוריתמים אחרים. לא מספיק טוב בשביל חיזוי ערכים רציפים. מצריך סט אימון גדול יחסית בשביל להגיע לתוצאות טובות 	<div>C4.5</div>	<ol style="list-style-type: none"> בכל צומת של העץ, בוחר את התכונה שמפצלת בצורה היעילה ביותר את קבוצת הסיווג לקבוצות משנה המפורטות ביותר לסיווג הנוכחי או לאחר. קריטריון הפיצול הוא רווח המידע המנומל (הפרש באנטרופיה). התכונה עם רווח המידע הנורמלי הגבוה ביותר נבחרת כדי לקבל את ההחלטה. חזרה עבור העלים שנוצרו לעץ לעיל. 	<ul style="list-style-type: none"> מבצע גיזום ולכן פחות סבירות ל $overfitting$ מסוגל להתמודד עם ערכים חסרים ופחות נוקשה בהכנת הנתונים שלו. יכול לבצע חלוקה של העץ עם תכונות שלהם יש עלויות שונות 	<ul style="list-style-type: none"> יוצר עצי החלטה לא מאוזנים בגלל שתכונותיהם מביאות ל"רווח" הטוב ביותר. 	<ul style="list-style-type: none"> הגיזום יביא לחיזוי טוב יותר לעומת $id3$ יקל עלינו בהכנת המידע שכן לא נדרש להעביר את הנתונים הרציפים לבדידים ולכן יתן חיזוי טוב יותר.
<ul style="list-style-type: none"> חשוב האנטרופיה של כל תכונה של ערכת הנתונים חלוקה לקבוצות משנה באמצעות התכונה שעבורה האנטרופיה המתקבלת היא אידיאלית יצירת עץ החלטה המכיל תכונה זו. חזרה על כל קבוצות המשנה עם התכונות הנותרות 	<div>ID3</div>	<div>Information gain</div>	<ol style="list-style-type: none"> חשוב האנטרופיה של כל תכונה של ערכת הנתונים חלוקה לקבוצות משנה באמצעות התכונה שעבורה האנטרופיה המתקבלת היא אידיאלית יצירת עץ החלטה המכיל תכונה זו. חזרה על כל קבוצות המשנה עם התכונות הנותרות 		<ul style="list-style-type: none"> פתרון לא אופטימלי אין גיזום בניגוד ל- $c4.5$ ול- $cart$ לא תומך בערכים חסרים/רציפים 	<ul style="list-style-type: none"> יקשה עלינו מאוד את הביצוע מכיוון שלא תומך בערכים רציפים, ולנו בנתונים יש כמה וכמה ערכים (כאלו(גיל, משקל וכו'))
<ul style="list-style-type: none"> יכול לקבל ערכים קטנוניים וזוגיים רציפים. קל לחסרה לזכרון ולא מהתחום ולכן לדוגמה במקרה שלנו יהיה קל להעביר מידע זה לצבור ולמנות השמנת יתר באוכלוסיה. לא מתחייב שהקשר בין התכונה לסיווג יהיה ליניארי. בניגוד לרגרסיה ליניארית תוצאות חריגות לא משפיעות את התוצאות לחלוטין. 	<ul style="list-style-type: none"> תמיד בביצוע החלטות שלו זמן האימון הנדרש ארוך יחסית לסוגים אחרים בדרך כלל המסקנות דורשות חישוב ארוך יותר בשביל להגיע אליהם יש נטייה ל $overfitting$ בניגוד לאלגוריתמים אחרים. לא מספיק טוב בשביל חיזוי ערכים רציפים. מצריך סט אימון גדול יחסית בשביל להגיע לתוצאות טובות 	<div>Cart</div>	<ol style="list-style-type: none"> בכל צומת של העץ, בוחר את התכונה שמפצלת בצורה היעילה ביותר את קבוצת הסיווג לקבוצות משנה המפורטות ביותר לסיווג הנוכחי או לאחר. קריטריון הפיצול הוא לפי מדד גיני התכונה עם רווח המידע הנורמלי הגבוה ביותר נבחרת כדי לקבל את ההחלטה. חזרה עבור העלים שנוצרו לעץ לעיל. 	<ul style="list-style-type: none"> בניגוד ל $c4.5$ תומך בחיזוי רציף. לא מחשב ערכות כללים (פחות סבירות ל $overfitting$) חסין לערכים חסרים ופחות נוקשה בהכנת הנתונים שלו. יכול לבצע חלוקה לפי עלויות תכונה שונות 	<ul style="list-style-type: none"> עץ ההחלטה המתקבל הוא בינארי 	<ul style="list-style-type: none"> הגיזום יביא לחיזוי טוב לעומת $id3$ יקל עלינו בהכנת המידע שכן לא נדרש להעביר את הנתונים הרציפים לבדידים ולכן יתן חיזוי טוב יותר. אפשרי לבצע חיזוי יותר מדויק מאשר זה שאנו מנסים ליישם בכך שנחזה השמנת יתר שתלוי ב bin וכן בנתונים נוספים שניתנו לנו, ולהגיע למסקנות רלוונטיות יותר(רציפות)
<ul style="list-style-type: none"> אמון מהיר יותר ביחס לעצים נותן תוצאות מדויקות לתחום 	<ul style="list-style-type: none"> מניח שיש התפלגות אחידה של הנתונים. רשים למתחילים חריגים ביחס ל $id3$ 	<div>לינארית</div>		<ul style="list-style-type: none"> עובד רק אם מדובר בקשר ליניארי. 	<ul style="list-style-type: none"> רגיש מאוד לחריגים 	<ul style="list-style-type: none"> מכיוון שהבעיה הנתונה היא בדידה הישום של כלי זה יכול להיות אפשרי אבל יקשה עלינו מאוד את התהליך. מכיוון שיש לנו ערכים רציפים ובדידים התוצאות של עץ החלטה ינטו להיות טובות יותר.
<ul style="list-style-type: none"> יכול להתמודד עם נתונים לא מאוזנים. וכן מונע $overfitting$ 	<div>חסרונות:</div> <ul style="list-style-type: none"> זמן חישוב ארוך יחסית האלגוריתמים 	<div>יתרונות:</div> <ul style="list-style-type: none"> זמן חישוב ארוך יחסית האלגוריתמים 	<p>מטא מסווג (כמה עצי סיווג אקראיים) כאשר כל עץ סיווג תוך שימוש בתת רשימה של מאפיינים(אקראיים) מתוך כלל המאפיינים</p>		<p>יכול להיות יעיל מאוד לבעיה הנתונה בעקבות רמת הפירוט שיש לכל רשומה וכך שבעץ אקראי יבחנו בכל שלב רק מספר תכונות לסיווג ולא כלל התכונות ובעזרת כך להביא לתוצאות מיטביות</p>	