

## דו"ח מיני פרויקט – אנליזה של מידע

### שמות המגישים:

רן עזרא, ת.ז. 302779657

נדב חנמ, ת.ז. 313260648

### שאלה 1

בגרסת המימוש השנייה, בחרנו לבדוק עבור כל פיקסל נתון לשכלל גם את הפיקסלים שבסביבתו.

לדוג', עבור פיקסל 324, נרצה לחשב בנוסף לערכו גם את ערכם של הפיקסלים הבאים:

295	296	297
323	<u>324</u>	325
351	352	353

לצורך כך, נדרשנו לחשב את העמודה והשורה שמס' פיקסל נתון מייצג, ולחשב את מספרי הפיקסלים של העמודות והשורות הסמוכים.

כמובן שנדרשנו לחישוב ספציפי עבור מקרים שונים בהם פיקסל נתון הוא בשולי הדף.

לדוג' עבור פיקסל 5, קיימים 5 תאים סמוכים בלבד, שכן לא קיימת שורה מעליו:

4	<u>5</u>	6
32	33	34

אל מול רעיון מרכזי זה, ניסינו מגוון וריאציות שונות, ע"מ לקבל את הפתרון האפקטיבי ביותר.

במהלך גיבוש הרעיון, ניסינו לבחון את השפעת התנאים הבאים על אחוז השגיאה :

### (1) שיטת החישוב:

- א. Average – בשיטה זו, מתבצע חישוב ממוצע הערכים של סביבת פיקסל נתון ובדיקה האם הערך שמתקבל גדול ממס' קבוע כלשהו – "MinAvg"
- ב. Counter - בשיטה זו, מגדירים מס' מינימום "MinCount", כך שהתנאי מתקיים עבור כל פיקסל שבסביבתו יש לכל הפחות "MinCount" פיקסלים הגדולים מ – 128.
- ההבדל בין השיטות מתבטא במשקל שמקבלים ערכים "גבוליים" של פיקסלים נתונים. כלומר, פיקסלים שערכם קרוב יחסית לערך 128 משפיעים באופן שונה בכל אחת מהשיטות.
- לדוג', נוכל להתבונן בסביבת הפיקסל הבאה :

120	120	120
120	255	120
120	120	120

לפי שיטה א', עבור  $\text{MinAvg} = 128$ ,

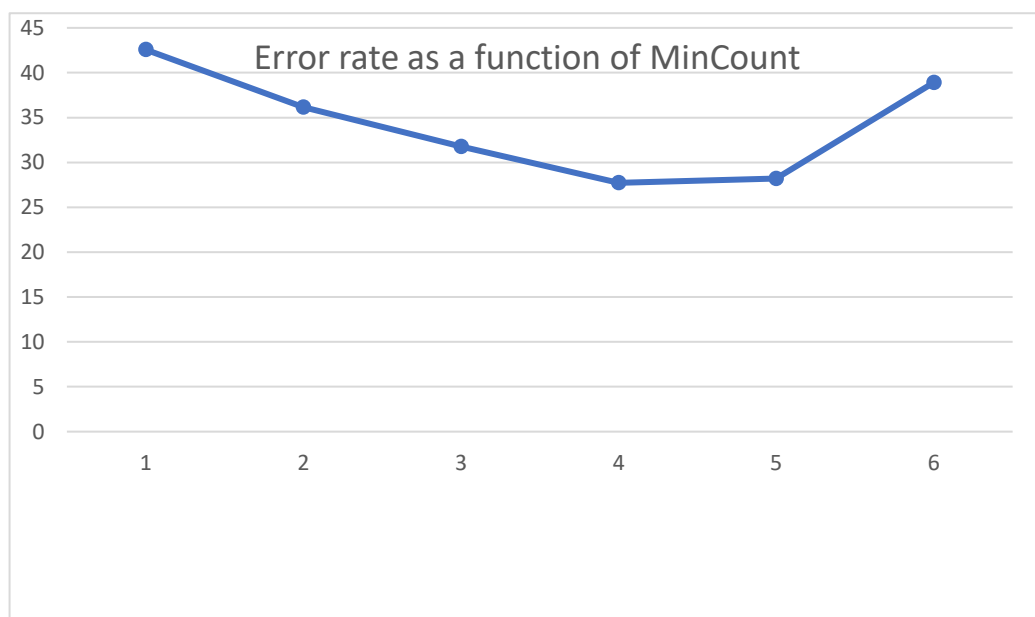
הערך שיתקבל עבור ממוצע סביבת הפיקסל הנתון יהיה  $\frac{(8 \cdot 120) + 255}{9} = 135$ .

כלומר, **סביבת הפיקסל תעמוד בשיטה א'.**

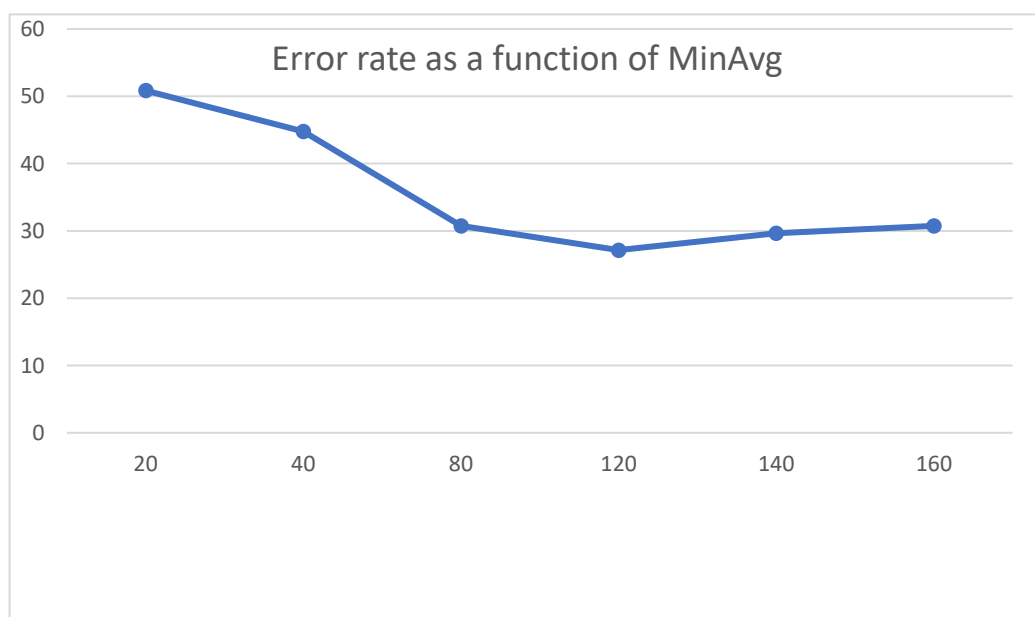
לעומת זאת, לפי שיטה ב', עבור ערך  $\text{MinCount} = 4$ ,

כמות הפיקסלים בסביבת הפיקסל הנתון שערכם 128 ומעלה היא 1, ולכן סביבת הפיקסל **לא תעמוד בשיטה ב'.**

הנחנו כי ערך MinCount האידיאלי יהיה כמחצית מכמות הפיקסלים בסביבתו של כל פיקסל. למעט בשולי התמונה, בסביבתו של כל פיקסל נמצאים 8 פיקסלים, ולכן הערכנו כי ערך MinCount האופטימלי יהיה 4, אך בחנו גם ערכים אופציונליים אחרים, המוצגים בגרף הבא, כאשר  $P=10$ ,  $L=10$  וגודל המדגם הוא 10,000.



בחנו גם את שיטת AVG עבור ערכי MinAvg שונים, המוצגים בגרף הבא, כאשר  $P=10$ ,  $L=10$  וגודל המדגם הוא 10,000.

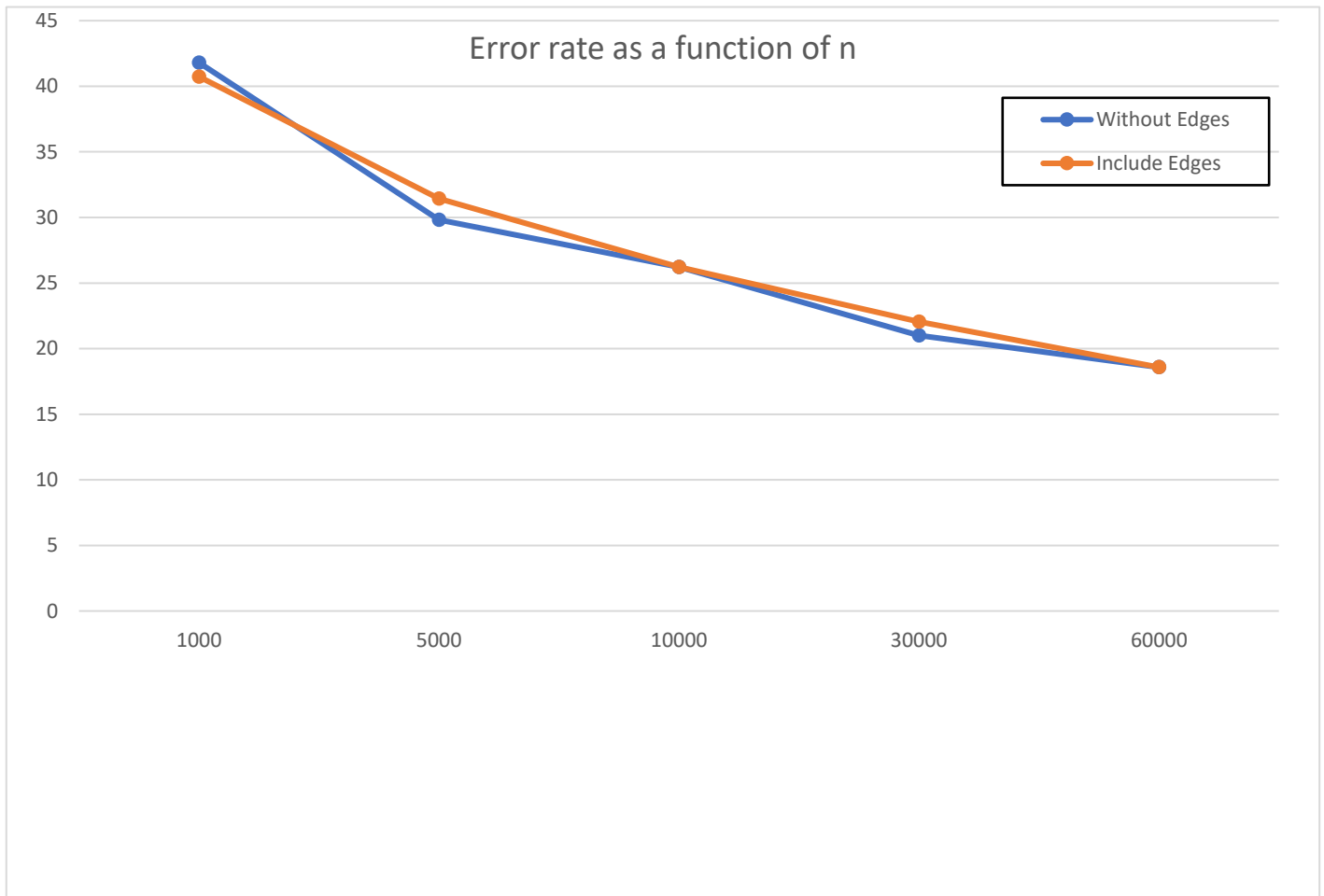


## (2) סיווג הפיקסלים השונים בהקשר ההשפעה שלהם על התוצאה:

פיקסלים שמייצגים מידע ב"קצוות התמונה" למשל, משמעותיים הרבה פחות אל מול פיקסלים שנמצאים במרכז התמונה.

ניסינו לבחון 2 גישות – דילוג על חישוב "סביבת הפיקסלים" בקצוות התמונה וחיסוב ערך הפיקסלים בלבד, אל מול חישוב סביבת הפיקסלים החלקית בקצוות התמונה.

בחרנו לעבוד לפי שיטת Counter,  $L=10$ ,  $P=10$  וגודל מדגם משתנה.



ראינו כי קיים הבדל מינורי בין השיטות, כאשר העדיפה מביניהן היא השיטה שלא כוללת את בדיקת סביבתם של הפיקסלים בקצוות התמונה.

### (3) מציאת מכנה משותף בין רשומות עם תוויות זהות:

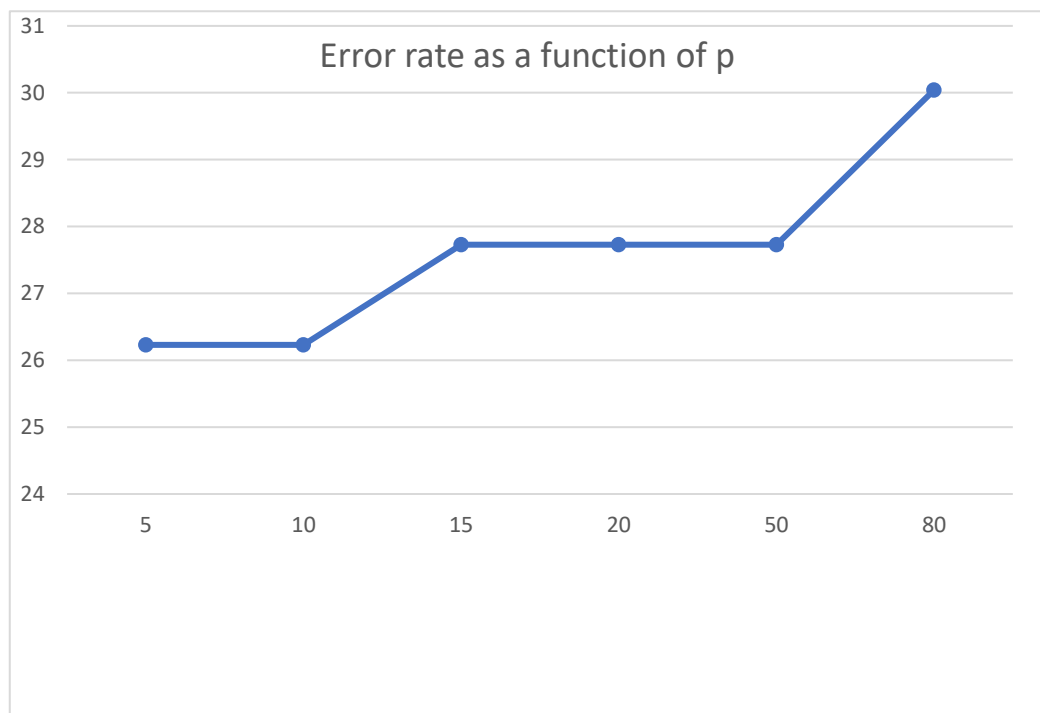
ניסינו לאפיין האם קיים מכנה משותף כלשהו המייחד כל אחת מהספרות, ולהטמיע את התובנות שלנו כתנאים. לאחר ניסיונות שונים ומגוונים, הבנו כי ללא הנחות נוספות על הקלט (מיקום הספרה ביחס לדף שעליו היא מופיעה למשל, או גודל הספרה ביחס לדף שעליו נכתבה), יהיה לנו קשה עד בלתי אפשרי ליצור עץ החלטה שיהיה רלוונטי לכל קלט.

### (4) שינוי גודלו של P:

בחנו את תפקודם של ערכי P שונים.

בחרנו לבחון את השפעתו של P לפי שיטת Counter,  $L=10$ , וגודל מדגם של 10,000.

קיבלנו את הגרף הבא:

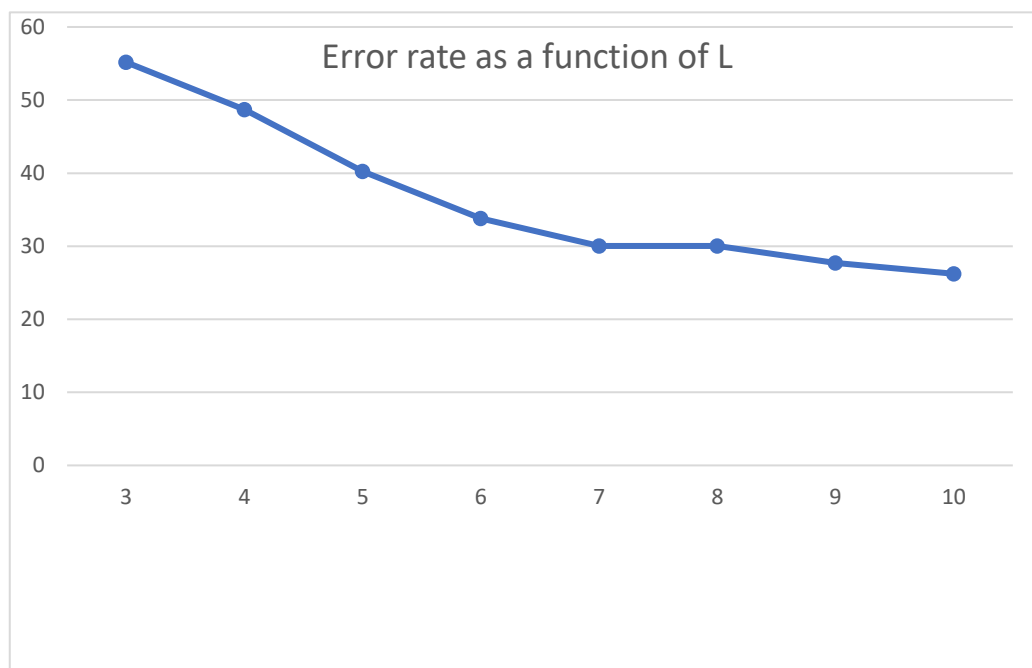


(5) שינוי גודלו של L:

בחנו את תפקודם של ערכי L שונים.

בחרנו לבחון את השפעתו של L לפי שיטת Counter,  $P=10$  וגודל מדגם של 10,000.

קיבלנו את הגרף הבא:

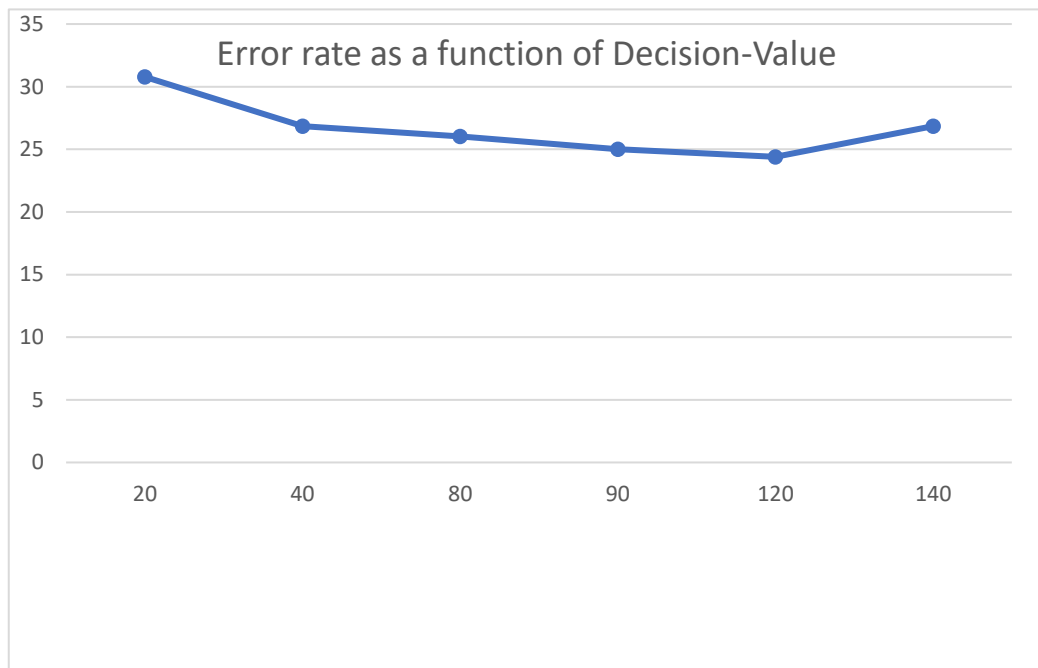


**(6) שינוי הערך "128" :**

ניסינו לבחון את התנהגות האלגוריתם עבור ערכים שונים, בעיקר נמוכים יותר מ-128.

בחרנו לבחון את השפעתו של ערך זה לפי שיטת Counter,  $L=10$ ,  $P=10$  וגודל מדגם של 10,000.

קיבלנו את הגרף הבא :



אל מול התוצאות שהניבו שיטת Average ושיטת Counter, בחרנו להשתמש בסופו של דבר בשיטת Counter, עם הפרמטרים הבאים:

$\text{MinCount} = 4$  ושימור הערך 128.

למרות העובדה שהפרמטרים  $L$  ו  $P$  השפיעו על אחוז השגיאה, הם מתקבלים כקלט לאלגוריתם ע"י המשתמש ולכן לא מוגדרים כחלק מהאלגוריתם.

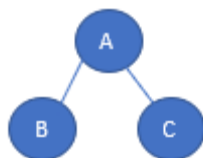
## שאלה 2

במהלך כתיבת הפרויקט, נתקלנו באתגרים שונים. אלו האתגרים העיקריים:

בחירת מבני נתונים אידיאליים:

### (1) פלט אלגוריתם בניית העץ –

התלבטנו רבות כיצד לייצא את העץ שנוצר ע"י אלגוריתם הלמידה לאלגוריתם החיזוי. ניסינו לעשות זאת באמצעות ייצוא אובייקט של עץ, אך התקשינו לממש זאת באמצעות ספריות קיימות, והבנו שהאובייקט כולו שומר מידע רב שלא רלוונטי לאלגוריתם הפענוח. לכן, בחרנו ליצור ערך String בעלות זיכרון מינימאלית, המייצג את ערכי הקודקודים בעץ בלבד ללא שאר הערכים שמבנה נתונים זה שומר, וליצור פונקציה באלגוריתם הפענוח שמפענחת מחרוזת זו לעץ בינארי. לדוג, את העץ הבא



נייצג ע"י המחרוזת  $(A,(B),(C))$

### (2) שמירת מבנה נתונים יחיד עבור כלל התמונות ופניה אליו לפי אינדקס –

האלגוריתם מקבל תחילה מספר רב של תמונות (60,000) ולכל תמונה, מערך בן 784 תאים המייצגים את הפיקסלים שלו. במהלך בניית האלגוריתם שמנו לב שתמונות רבות מועתקות פעמים רבות, והדבר עולה בזמן ריצה וזיכרון שניתן לחסוך. על מנת לפתור בעיה זו העברנו את כל מבני הנתונים להחזיק אינדקס של תמונה בלבד ולא את התמונה עצמה, על מנת לפנות לערך התמונה נפנה למאגר הנתונים הראשי במיקום האינדקס. יש לציין שמאגר התמונות הוא יחיד וכל שאר המאגרים מחזיקים .int



### (3) ייצוג Node :

בחרנו לייצג Node ע"י השדות הבאים :

- (א) מאגר אינדקסי התמונות הרלוונטיות – לכל Node יש מאגר אינדקסים המייצג את האינדקסים של כלל התמונות המגיעות אליו במצב העץ הנוכחי. באופן זה, בעת פעולת הרחבה של העץ, נבדקים התנאים רק על תמונות אלו ולא על כלל התמונות במאגר.
- (ב) תנאים רלוונטיים – ע"מ לחסוך כפילויות של תנאים, כל Node מחזיק אך ורק מאגר של תנאים הרלוונטיים אליו, ללא תנאים שכבר השתמשנו בהם בדרך אליו (כך נחסכת בדיקה כפולה של תנאי).
- (ג) IG – כל Node מתחזק את ערך ה-IG האופטימאלי עבורו, שחושב פעם אחת לכל Node ומחזיק בנוסף את החלוקה לבנים לפי התנאי שהניב ערך IG זה.

### זמן ריצה :

לאור העובדה שהאלגוריתם נדרש לעבד כמות גדולה יחסית של מידע (כ- 60,000 רשומות המכילות כל אחת 784 אינדקסים), כל פעולה שבתכנות רגיל לא מקבלת "טיפול מיוחד" בהקשר זמני ריצה כמו מעבר על כלל הפיקסלים של תמונה נתונה, הופכת למשמעותית.

הבנו שהמפתח לקיצור זמן הריצה הינו ע"י ביצוע Trade-Off אל מול משאבי זיכרון.

לכן, בחרנו להשתמש בכמות גדולה יחסית של זיכרון ע"מ להרוויח זמני ריצה יעילים יותר בפעולות שונות.

היישום של עיקרון זה בא לידי ביטוי, למשל, בהגדרת המחלקה "Node" – בשונה מעץ בינארי רגיל, הוספנו למחלקה שדות נוספים שישמרו מידע שכבר חושב בעבר, למשל רשימת האינדקסים של רשומות שכבר הגיעו לאותו Node ספציפי.

בנוסף, מימשנו את הפתרונות הבאים :

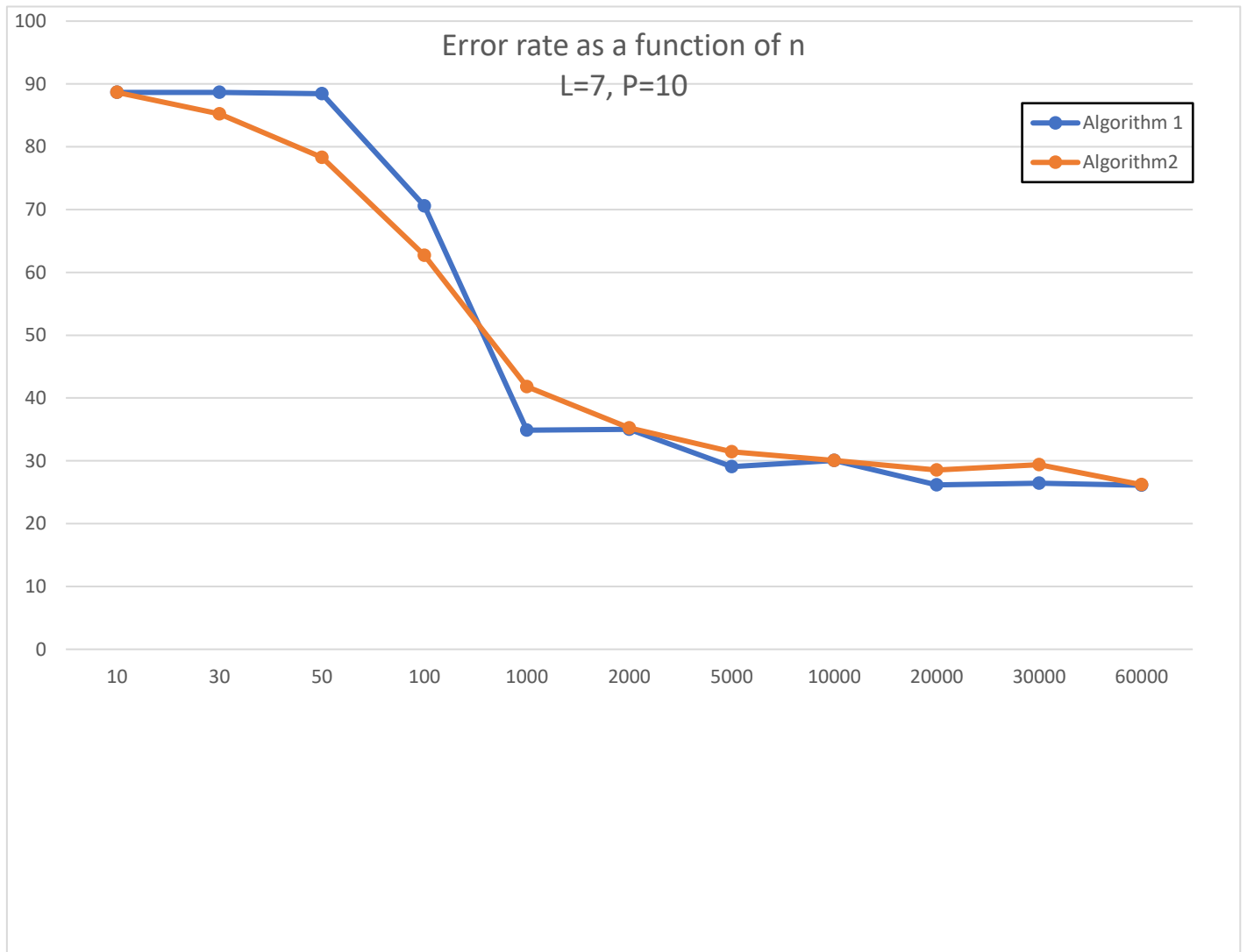
(1) **מערך דו – מימדי השומר את תוצאות התנאים** – בהינתן קבוצת תנאים, נחשב מבעוד מועד בפעם אחת, עבור כל תמונה את הפעלת קבוצת התנאים על כל פיקסל שבה ואת התוצאה נשמור במערך דו – מימדי גלובאלי. בכל פעם שנרצה לחשב IG נעבור על כל התנאים עבור כל התמונות לקבלת התנאי המקסימאלי, נוכל לגשת ב  $O(1)$  למאגר גלובאלי זה ולגלות את תוצאת תנאי ספציפי על פיקסל ספציפי. במהלך ביצוע בדיקות, הוספת עדכון זה צמצמה בכמעט חצי את זמן הריצה.

(2) **שימוש בערמת מקסימום עבור ערכי IG** – בכל איטראציה של בניית העץ, אנו נדרשים לבחור להרחיב את העלה בעל ה-IG המקסימאלי. ע"מ לעשות פעולה זו ביעילות, בכל איטראציה אנו מחזיקים ערמת מקסימום עבור כל העלים הממוינת לפי ערכי ה-IG שלהם. בצורה זו, שליוף ה-Node – בעל ערך ה-IG המקסימאלי מתבצעת ב  $O(1)$ .

### שאלה 3

#### סעיף א'

עבור כל ערך  $n$ , נדגמו 3 הרצות של כל אלג' ונבחר ערך השגיאה הממוצע.



## סעיף ב'

באופן כללי, ניתן לראות בצורה מובהקת ב – 2 האלגוריתמים, כי קיים יחס הפוך בין גודל המדגם לבין אחוז השגיאה.

כלומר, ככל ש-  $n$  גדול יותר, גודל השגיאה קטן.

בדיוק כמו שהיינו מניחים באופן אינטואיטיבי, ככל שעץ החיזוי נבנה ע"ב כמות נרחבת יותר של "תמונות אימון", איכות התוצר משתפרת.

ההבדלים הניכרים לעין הם שעבור  $n = 500$ , 2 האלגוריתמים נפגשים ב – Error rate דומה. עבור ערכים קטנים מ – 500, האלגוריתם הראשון מדויק פחות, ואילו עבור ערכים גדולים מ – 500, האלגוריתם הראשון מדויק יותר.

## סעיף ג'

ניתן לראות שעבור קלטים גדולים, האלגוריתם בגרסתו השנייה מדויק פחות מהאלגוריתם בגרסתו הראשונה. לאור העובדה שההפרשים בין אחוזי השגיאה הם מינוריים בין 2 הגרסאות, להנחתנו, סביבת הפיקסל אינה מעידה בצורה משמעותית על ערכו של פיקסל מסוים.

## סעיף ד'

גירסה ראשונה: גודל המדגם = 10,000,  $P = 10\%$ ,  $L = 7$

9	8	7	6	5	4	3	2	1	0	
12	9	21	59	47	5	11	25	0	791	0
0	36	11	9	9	1	2	10	1057	0	1
34	145	38	78	17	24	19	618	25	0	2
62	32	29	20	185	5	598	43	16	20	3
106	14	30	61	12	701	35	13	5	5	4
27	14	19	51	452	36	225	24	5	39	5
13	23	0	748	28	50	8	51	10	27	6
73	19	813	8	18	29	3	49	11	5	7
54	538	25	167	29	44	29	38	23	27	8
679	21	124	18	51	47	45	8	8	8	9

- ניתן לראות כי כמות ניכרת מהתמונות שערך 3 פוענחו להיות הערך 5 (185 תמונות) ומנגד, המון תמונות שערך 5 פוענחו להיות הערך 3 (225 תמונות). מה שמצביע על כך שספרות אלו דומות מבחינת האלגוריתם הראשון.  
(תוצאה דומה יצאה גם בבדיקה על האלגוריתם השני).
- הערך 4 סווג להיות הערך 9 פעמים רבות אך מנגד, הערך 9 לא סווג להיות הערך 4 כמות פעמים ניכרת  
(תוצאה דומה יצאה גם בבדיקה של האלגוריתם השני).
- הערך 8 סווג להיות הערך 6 פעמים רבות אך מנגד, הערך 6 לא סווג להיות הערך 8 כמות פעמים ניכרת.
- הערך 2 סווג להיות הערך 8 פעמים רבות אך מנגד, הערך 8 לא סווג להיות הערך 2 כמות פעמים ניכרת.

גירסה שנייה: גודל המדגם = 10,000,  $P = 10\%$ ,  $L = 7$

9	8	7	6	5	4	3	2	1	0	
35	7	11	31	27	16	23	10	0	820	0
0	24	10	5	8	0	18	33	1037	0	1
22	52	58	62	25	18	34	711	20	30	2
21	65	20	9	130	4	714	20	5	22	3
131	20	13	45	37	673	34	19	3	7	4
48	42	12	50	523	36	113	21	5	42	5
9	53	16	708	31	37	7	40	9	48	6
68	12	843	4	12	14	26	42	6	1	7
74	592	21	54	65	38	47	57	17	9	8
758	25	32	14	47	61	48	9	5	10	9

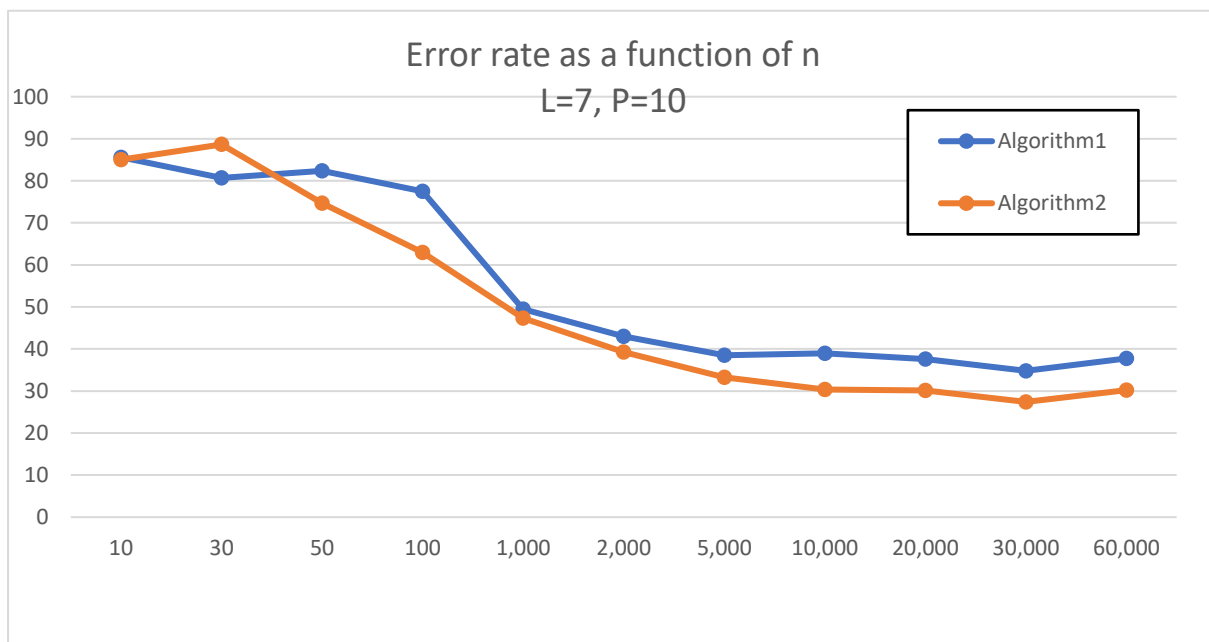
- ניתן לראות כי כמות ניכרת מהתמונות שערך 3 פוענחו להיות הערך 5 (185 תמונות) ומנגד, המון תמונות שערך 5 פוענחו להיות הערך 3 (225 תמונות). מה שמצביע על כך שספרות אלו דומות מבחינת האלגוריתם השני.
  - הערך 4 סווג להיות הערך 9 פעמים רבות אך מנגד, הערך 9 לא סווג להיות הערך 4 כמות פעמים ניכרת.
- מעיון בתוצאותיהם של 2 האלגוריתמים, ניתן להבחין כי התקבלו ע"י האלגוריתם השני ערכים נמוכים יותר עבור מרבית המספרים.
- הבלבול המשמעותי והדו – כיווני בין הספרות 3 ו – 5 עדיין נשמר, אך ב"עוצמה נמוכה יותר", כמו גם הבלבול הניכר בסיווג הערך 4 לערך 9.
- אך מנגד, כמות הפעמים שהערך 8 סווג להיות הערך 6 והערך 2 סווג להיות הערך 8 ירדה משמעותית.
- מכאן, ניתן להסיק שבדיקה בסביבת הפיקסל משפיעה באופן ניכר על פיזור השגיאה בין ערכים שונים ולכן האלגוריתם השני עדיף מבחינה זו.

#### שאלה 4

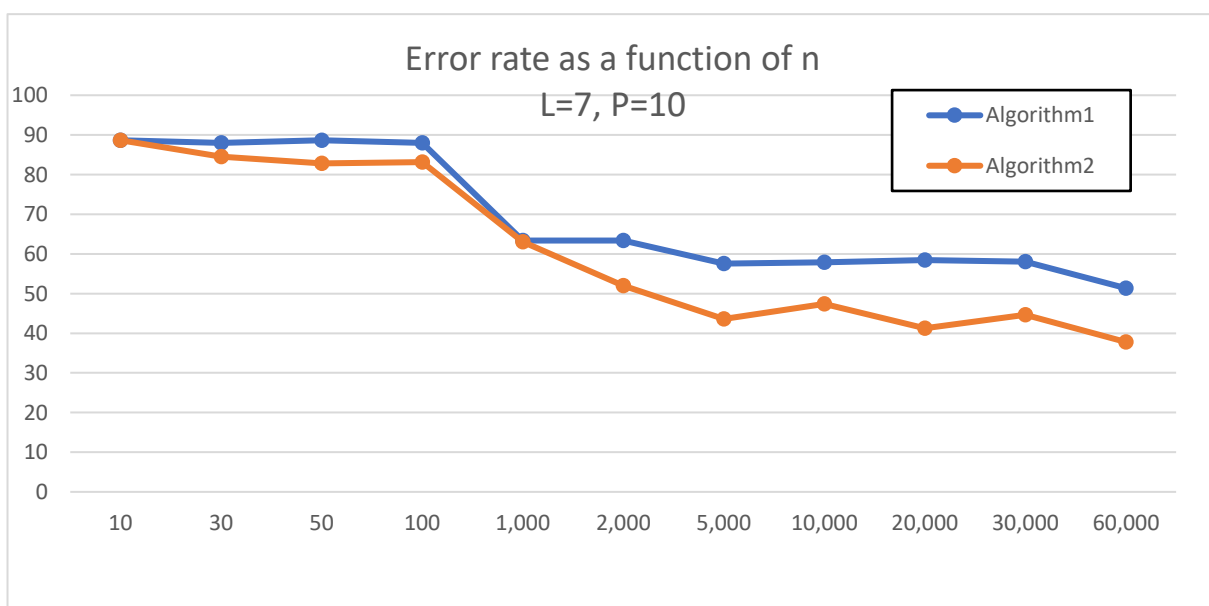
#### סעיף א'

עבור כל ערך  $n$ , נדגמו 3 הרצות של כל אלג' ונבחר ערך השגיאה הממוצע.

$\alpha = 0.1$



$\alpha = 0.3$



$$\alpha = 0.6$$



### סעיף ב'

ניתן לראות כי עבור ערכי  $\alpha$  נמוכים יחסית, אלגוריתם 2 בעל אחוז שגיאה נמוך יותר מאלגוריתם 1.

בנוסף לכך, ב – 2 האלגוריתמים עבור ערך  $\alpha$  נמוך יחסית, ניתן לראות כי אחוז השגיאה יורד ככל שגודל המדגם גדל.

מנגד, עבור ערך  $\alpha$  גבוה (0.6), ניתן לראות כי אחוז השגיאה עולה באופן דרמטי ב – 2 האלגוריתמים, ונשאר גבוה גם כאשר גודל המדגם גבוה.

### סעיף ג'

להנחתנו, ההבדל באחוז השגיאה בין אלגוריתם 1 ואלגוריתם 2 עבור רמות רעש נמוכות יחסית נובע מכך שהבדיקה בסביבתו של כל פיקסל נתון "מנטרלת" את הרעש בכך שהיא דוגמת עבור כל פיקסל את הנקודות הצמודות אליו, שהסתברותית (כאשר מקדם הרעש נמוך יחסית) לא השתנו ובכך האפקט של הרעש קטן.

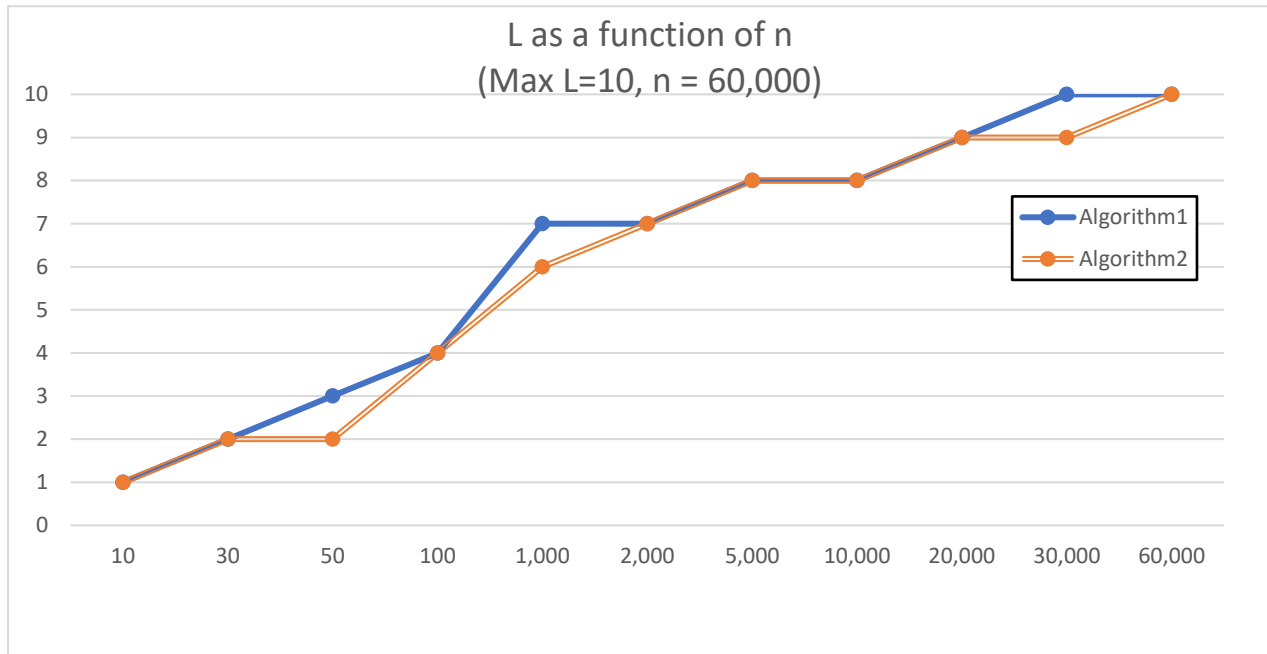
מאידך, כאשר מקדם הרעש הוא גבוה מדי, גם סביבתו של הפיקסל מורעשת מספיק ולכן, איכות הדיוק של 2 האלגוריתמים נפגעת באופן דומה.



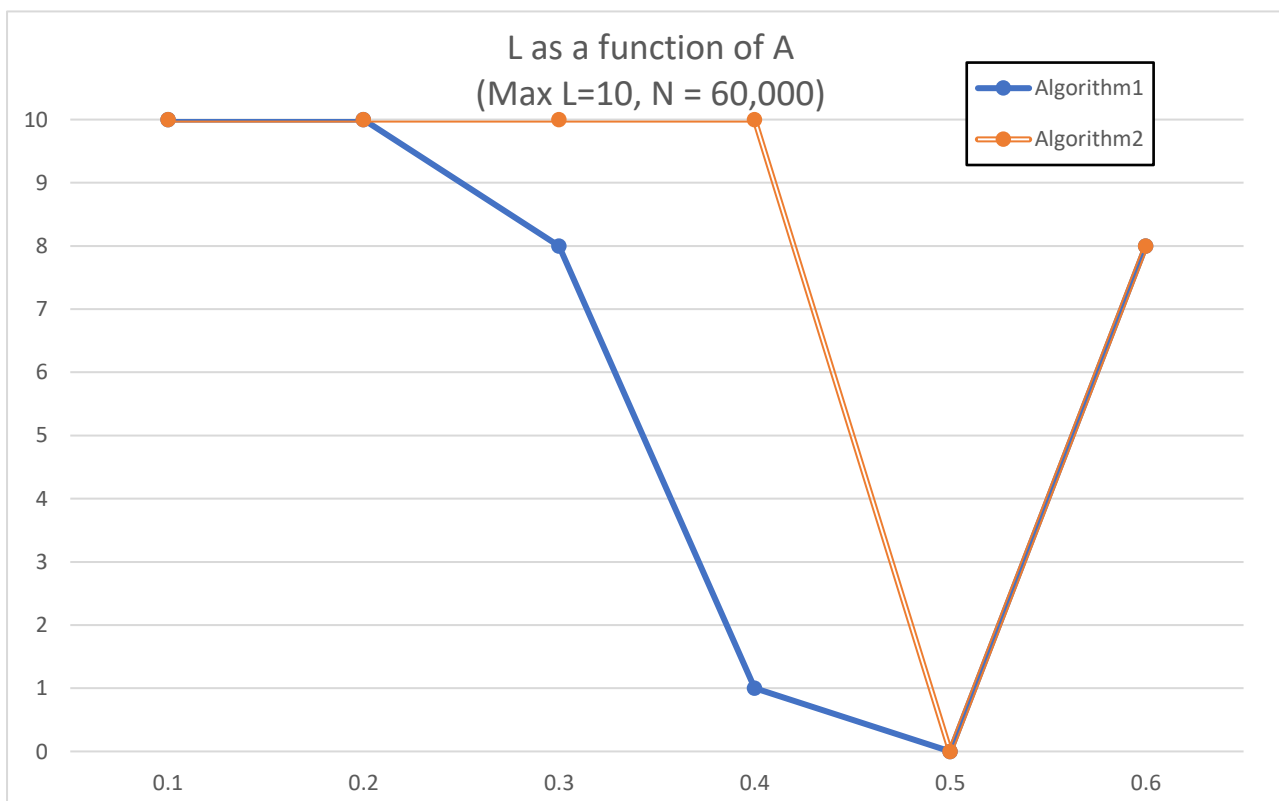
## שאלה 5

### סעיף א'

גודל העץ כתלות בגודל המדגם:



גודל העץ כתלות במקדם הרעש:



## סעיף ב'

מהגרף המתאר את גודל העץ כתלות בגודל המדגם, ניתן לראות עבור 2 האלגוריתמים כי ככל שגודל המדגם גדול יותר, מתקבל עץ גדול יותר.

מהגרף המתאר את גודל העץ כתלות במקדם הרעש, ניתן לראות כי ככל שערכי  $\alpha$  גדלים, גודל העץ קטן.

בנוסף, ניתן לראות מהגרף כי אלגוריתם 2 מייצר עץ גדול יותר גם עבור מקדמי רעש גבוהים יחסית, שבאלגוריתם 1 גרמו ליצירת עץ בגודל קטן.

## סעיף ג'

להבנתנו, תוצאותיו של הגרף המתאר את גודל העץ כתלות בגודל המדגם הגיוניות, שכן עבור מס' דוגמאות גבוה יותר, מדגם האימון מגוון יותר, ויש לאלגוריתם יותר מידע לעבד ולנתב, לכן ככל שמס' צמתי ההחלטה שהאלגוריתם מייצר גדול יותר, כך המידע מעובד בצורה מדויקת יותר.

עבור גרף מקדם הרעש, התוצאות מתכתבות עם התוצאות שהתקבלו בסעיפים הקודמים, המצביעות על כך שהאלגוריתם בגרסתו השנייה פחות רגיש לרעשים (ניתן לראות זאת במסקנותיו של סעיף 3 בשאלה הקודמת). נזכור כי מדדנו את השפעת מקדם הרעש על גודל מדגם מקסימאלי ולכן, הגיוני שיתקבלו תוצאות השקולות לתוצאות המתקבלות עבור גודל מדגם גבוה יחסית ללא מקדם רעש כלל.