

הוראות לנבחן בצידו השני של הדף

אין לכתוב מעבר לשוליים משני צידי הדף

חלק ב' /

$$\frac{28}{32} - 1 = \frac{26}{30}$$

שאלה 1 -

$$\frac{0.5}{10}$$

שאלה 2 -

$$\frac{8}{10}$$

מס' כיתה 306 בנין 1532

מס' נבחן 020313

20/02/2012 1345302 1 12/12

הנדסת מערכות מידע

אחזור מידע וספריות דיגיטליות

03721440601101




לשימוש המרצה הבודק

יחידות ועשרות ומאות

<input type="checkbox"/>	<input type="checkbox"/>	0
<input type="checkbox"/>	<input checked="" type="checkbox"/>	1
<input type="checkbox"/>	<input checked="" type="checkbox"/>	2
<input type="checkbox"/>	<input checked="" type="checkbox"/>	3
<input type="checkbox"/>	<input checked="" type="checkbox"/>	4
<input type="checkbox"/>	<input checked="" type="checkbox"/>	5
<input type="checkbox"/>	<input checked="" type="checkbox"/>	6
<input type="checkbox"/>	<input checked="" type="checkbox"/>	7
<input type="checkbox"/>	<input type="checkbox"/>	8
<input type="checkbox"/>	<input type="checkbox"/>	9

ציון הבחינה 76

שם המרצה קסם שמש

חתימה

תאריך

המחלקה למע' מ'ק' שנה 2

תאריך בחינה 20/2/12

מקצוע בחינה אמצ' מ'ק'ע





## אוניברסיטת בן-גוריון בנגב

### הוראות לנבחן

1. בהגיעך למקומך יש להניח את כרטיס הנבחן ותעודה מזהה על שולחןך.
2. אסור להביא למקום הבחינה תיקים, ספרים, מחברות, טלפון נייד או רשימות פרט למותר על פי שאלון הבחינה.
3. עזב תלמיד את האולם אחרי חלוקת השאלונים, דינו כדין "נבחן" בבחינה.
4. אסור לנבחן לשוחח בזמן הבחינה, או לעזוב את מקומו ללא נטילת רשות.



איחזור מידע תשע"א – 372.1.4406

סמסטר חורף מועד א' 31.01.11 - גירסא ב'

פרופ' ברכה שפירא, אורלי מורנו

משך המבחן : שעתיים וחצי

חומר עזר: מותר (לא מחשב נייד), מותר מחשבון

יש להחזיר את גיליון הבחינה. המבחן כולל 4 דפים

חלק א' - 30% - יש לענות על גבי הטופס

סמן תשובה אחת נכונה

1. 5% נתון מאגר עם 4 המסמכים הבאים:

$D_1$ : alpha bravo charlie delta echo foxtrot golf

$D_2$ : golf golf golf delta alpha

$D_3$ : bravo charlie bravo echo foxtrot bravo

$D_4$ : foxtrot alpha alpha golf golf delta

לשאלתא alpha bravo

- א. מנוע שמבוסס על vector space יחזיר את כל המסמכים ואת מסמך  $D_1$  במקום הראשון
- ב. מנוע שמתבסס על vector space יחזיר את כל המסמכים ואת מסמך  $D_3$  במקום הראשון
- ג. מנוע שמתבסס על מודל בוליאני טהור יחזיר את מסמך  $d_1$
- ד. מנוע שמתבסס על מנוע בוליאני טהור לא יחזיר את מסמך  $d_2$
- ה. ג+ד נכונים
- ו. ב+ג נכונים

2. 4% הפרמטר d בנוסחת Pagerank

- א. מאפשר למדל באופן נכון את התנהגות ה Random surfer (הגולש האקראי) ולפתור את בעיית ה spider trap
- ב. מאפשר לחשב pagerank באופן מהיר
- ג. מאפשר לנרמל את תוצאת ה Pagerank
- ד. הוא הווקטור העצמי של מטריצת המעברים
- ה. מאפשר למשתמש לעבור לדף באופן אקראי
- ו. א+ה נכונים

3. 4% מדדי ההערכה הבאים מתאימים להערכת מנועי חיפוש אינטרנטיים:

- א. DCG , F-measure, precision@k , R-precision
- ב. MAP , interpolated average precision , precision@k , DCG
- ג. Reciprocal Rank , precision@k , MAP , DCG
- ד. MAP , precision@10 , Fallout , NDCG
- ה. כל המדדים שהוזכרו בהרצאות יכולים להתאים גם למנועי חיפוש אינטרנטיים
- ו. אף אחד מהמדדים שהוזכרו בהרצאות לא יכול להתאים בלי שינויים בהגדרות





4. 4% עדכון של האינדקס מיד עם הגעה של מסמך חדש מה crawler
- מאט את פעולת האינדוקס
  - משפר בהכרח את ה precision של שאילתות רלוונטיות
  - משפר בהכרח את ה recall של שאילתות רלוונטיות
  - מעלה את העדכניות של האינדקס אך מסבך את הניהול שלו
  - משפר את המהירות של ביצוע השאילתא.
5. 5% ההשפעה של חוק zipf באה לידי ביטוי באורך רשימות postings של terms על פי הפירוט הבא (רשימות posting הן הרשימות הכוללות אינפורמציה על מופעים של Terms במסמכים):
- אין קשר בין הדברים
  - רשימות של מעט מילים תהיינה מאוד ארוכות ושל הרבה מילים תהיינה מאוד קצרות
  - רשימות של הרבה מילים תהיינה מאוד ארוכות ושל מעט מילים תהיינה מאוד קצרות
  - החוק יבוא לידי ביטוי רק אם יופעל תהליך stemming על המסמכים
  - החוק יבוא לידי ביטוי רק אם יופעל תהליך של הסרת stopwords
  - ב+ד נכונים
  - ב+ה נכונים
  - ג+ה נכונים
- סמן נכון או לא נכון:
6. 4% כדי לאפשר מענה לשאילתות של ביטויים האינדקס חייב להיות במבנה של biword index . נכון/לא נכון
7. 4% Relevance feedback יכול להתבצע באמצעות האלגוריתם של Rocchio אשר מעדכן את ווקטור השאילתא שהמשתמש הקליד על פי תגובת המשתמש לתוצאות שהמנוע החזיר לשאילתא בהתאם לערכי הפרמטרים שקובעים את רמת ההתחשבות בשאילתת המשתמש לעומת תגובתו לתוצאות השאילתא . נכון/לא נכון





חלק ב' 70%

ענה על השאלות הבאות:

1. 30%

הנח שאילתת q1 שלה 5 מסמכים רלוונטים במאגר על פי הפירוט הבא,

מסמך d13 רלוונטי באופן מושלם (ציון 3 על סקאלה של 0-3 כאשר 0 הוא לא רלוונטי ו-3 רלוונטי באופן מושלם), מסמכים d10 ו d9 רלוונטים מאוד (ציון 2), מסמכים d1 ו d3 רלוונטים באופן סביר (ציון 1) וכל שאר המסמכים במאגר לא רלוונטים. (הערה: כל המסמכים שקיבלו ציון גבוה או שווה ל-1 נחשבים רלוונטים). מנוע E1 הפעיל את שאילתת q1 על המאגר הנ"ל והחזיר 20 מסמכים. ידועים רק 9 המסמכים הראשונים שהמנוע החזיר ואת הסדר שבו החזיר אותם (משמאל לימין):

✓ × × × ✓ ✓ ✓ ✓ ×  
d13, d2, d8, d15, d3, d10, d1, d9, d4.....

ענה על השאלות הבאות:

- 2% מהו ה precision ב 10 מסמכים של המנוע, e1 על פי השאילתת הנ"ל
- 3% אם המנוע היה אידיאלי לשאילתת הנ"ל (כלומר היה מחזיר את כל המסמכים הרלוונטים ראשונים, מדורגים לפי רמת הרלוונטיות שלהם), מה היה אז ה precision ב 10 מסמכים?
- 4% מהו ה r-precision, ומהו ה Reciprocal Rank של מנוע e1 על פי השאילתת הנ"ל
- 4% מהו ה precision ב 30% recall?
- 5% חשב DCG ב 5, השתמש ב  $1/\log_2(\text{rank})$  discount של
- 5% מנוע e1 הריץ שאילתת q2 שלה 15 מסמכים רלוונטים במאגר, המנוע החזיר 15 תוצאות על פי הסדר הבא (משמאל לימין, על כל מסמך שחזר מסומן אם הוא רלוונטי או לא, כאשר "ר" מסמן מסמך רלוונטי ו "ל" מסמן מסמך לא רלוונטי):  
ל,ר,ל,ל,ר, ל, ל, ל,ר, ר,ר,ל,ר,ר,ר,

חשב את interpolated average precision על פי שתי השאילתות.

- 3% חשב MAP של המנוע על פי שאילתת q1.
- 6% איך ישתנה ה MAP אם מתווסף מסמך רלוונטי נוסף למאגר והמנוע מחזיר אותו בין ה 20 המסמכים הראשונים. (הראה את ה MAP המקסימלי והמינימאלי האפשריים בעקבות השינוי).

- 5% במערכת סינון שמקבלת באופן קבוע זרימה של מסמכים חדשים, לכל מסמך שמגיע למערכת מחושב הדמיון שלו עם פרופיל המשתמש. נניח שכל מסמך מיוצג כווקטור על פי גישת  $tf \cdot idf$  וכך גם פרופיל המשתמש. מהו idf במקרה זה, כלומר לאיזה מאגר אפשר לייחס את חישוב ה idf?

- 15% לצורך תיקון שגיאות של שאילתות, מחשבים מרחק בין מילת שאילתת שמשתמש הקליד ואיננה נמצאת באינדקס ומילים פוטנציאליות שאליהן התכוון המשתמש. הצע **נוסחא** לחישוב מרחק בין מילה שגויה לבין מילה פוטנציאלית מתוקנת שתתחשב במרחק בין אותיות במקלדת (בין אותיות שאולי הוחלפו בין המילים לתיקון) וכן בהפרש בין אורך המילה המקורית לאורך המילה המוצעת כתיקון. הראה באופן מפורט את כל חלקי הנוסחא הנדרשים, למשל הראה כיצד מחושב מרחק בין האותיות במקלדת. (כתשובה יש להראות נוסחה ולא הסברים מילוליים)





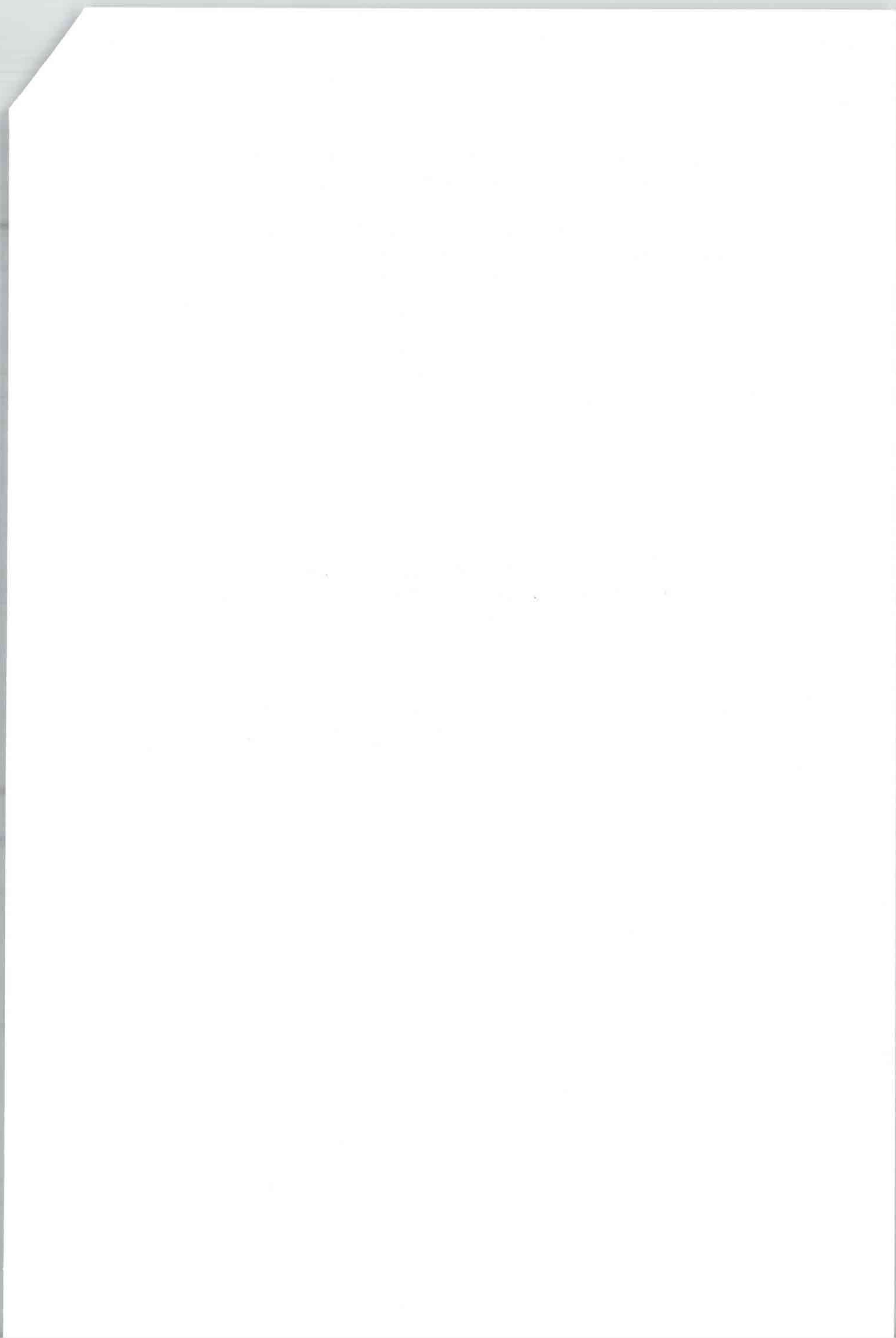


4. 10% קיימות שיטות שונות לקודד מילים על פי איך שהן נשמעות. השיטה הידועה ביותר היא שיטת ה Soundex – שהיא שיטה המבוססת על פונטיקה שהומצאה בשנת 1918. כל מילה מקודדת לאות שאחריה 3 ספרות. האות היא האות הראשונה של המילה, והספרות נקבעות על ידי כללים שונים שמקבצים כמה אותיות דומות לאותה ספרה. כך למשל, N i M מקודדות לסיפרה 5. V i , P , F , B מקודדות לסיפרה 1. אותיות ניקוד אינן מיוצגות, אלא אם כן האות הראשונה של המילה היא אות ניקוד. להלן כמה דוגמאות למילים מקודדות בשיטת soundex:

הקוד S-530 הוא הקוד של המילים: smythe-i smith  
הקוד a-450 הוא הקוד של המילים: Allan, Allen, Alan, Allyn ועוד.....  
Retrieve i retrieval מקודדות ל r-361 (משום שמתעלמים מאותיות עודפות מעבר לשלוש ספרות קוד).  
ציין שני יתרונות ושני חסרונות (שונים) לשימוש ב Soundex כשיטה לייצוג המילים באינדקס- במקום להשתמש ב Terms עצמם (או ב stem שלהם).

5. 10% מנועי חיפוש משתמשים בתוצאות של שאילתות זהות קודמות כדי לשפר תוצאות של שאילתא נוכחית. כלומר, המנועים שומרים במאגר מיוחד את השאילתות של המשתמשים, את התוצאות שהם החזירו וגם את התנהגות המשתמש עם התוצאות, כלומר אילו מהתשובות הוא אהב. כאשר משתמש מקליד שאילתא, המנוע מחפש במאגר המיוחד שאילתא זהה (אחת או יותר) ומשתמש בתוצאות של השאילתות הזוהות כדי לשפר את תוצאות השאילתא הנוכחית. הבעיה היא שכ 50% מהשאילתות שנשלחות למנועי חיפוש הן ייחודיות (כלומר, אין להן שאילתות זהות). הסבר כיצד אפשר להשתמש בתוצאות של שאילתות קודמות, גם כאשר השאילתא שהמשתמש הקליד אינה זהה לשאילתא שקיימת במאגר המיוחד. (אין צורך להציג אלגוריתם מדויק בפסדו-קוד אלא לספק הסבר ברור על השיטה שאתם מציעים)

בהצלחה - ברכה ואורלי





12/12

$\frac{28}{38}$

א' ג' ח' ט' י' י"א

7.7.7

1

Precision	recall	rel?	doc#
1	0.2	✓	1
0.5	"	✗	2
0.33	"	✗	3
0.25	"	✗	4
0.4	0.4	✓	5
0.5	0.6	✓	6
0.57	0.8	✓	7
0.625	1	✓	8
0.55	"	✗	9
0.5	"	✗	10
		✗	

✓

0.5

10

✓

0.5

7

✓

$r_{precision} = 0.4$

2

$$ARR = \frac{1}{1} = 1 \quad \checkmark$$

ה' ח' ט' י' י"א recall - ה' ח' ט' י' י"א  
 30% דיוק  
 (ח) ~~אנדרבולטיות~~

3

7

$$DCG_5 = 3 + 0 + 0 + 0 + \frac{1}{\log_2 5} = 3.43 \quad \checkmark$$

1

לפי דיוק

Handwritten text on the left margin, partially visible and cut off by the edge of the page.



recall 15 - q2 ①

Position	recall	val?	best
1	0.06	✓	1
1	0.13	✓	2
1	0.2	✓	3
	"	X	4
0.8	0.26	✓	5
0.83	0.33	✓	6
0.85	0.4	✓	7
	"	X	8
	"	X	9
	"	X	10
0.63	0.46	✓	11
	"	X	12
	"	X	13
0.57	0.53	✓	14
	"	X	15

<del>cut-off</del> Avg precision	q2 Precision	q7 precision	interpolated
1	1	1	0
1	1	1	0.1
1	1	1	0.2
0.625	0.85	0.4	0.3
0.625	0.85	0.4	0.4
0.53	0.57	0.5	0.5
0.25	0	0.5	0.6
0.28	0	0.57	0.7
0.28	0	0.57	0.8
0.31	0	0.625	0.9
0.31	0	0.625	1

interpolated average precision = 
$$\frac{1 + 1 + 1 + 0.625 \cdot 2 + 0.53 + 0.25 + 0.28 \cdot 2 + 0.31 \cdot 2}{11} =$$

= 0.56

recall //

132

5

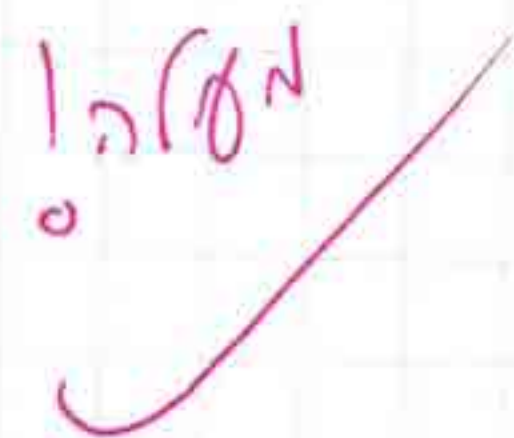
$$MAP(q_r) = \frac{1 + 0.4 + 0.5 + 0.57 + 0.625}{5} = 0.619$$



ה- MAP הממוצע של המיקום (המקום) ה-20, נקבע על ידי

7

$$\frac{MAP_{max}}{MAP} = \frac{1 + 1 + \frac{3}{6} + \frac{4}{2} + \frac{5}{8} + \frac{6}{9}}{6} = 0.72$$



הממוצע של המיקום ה-20, נקבע על ידי

$$\frac{MAP_{min}}{MAP} = \frac{1 + 0.4 + 0.5 + 0.57 + 0.625 + \frac{6}{20}}{6} = 0.56$$









הסדר לנסח:

קניין - נ"מ א זקנה (פ"ס) לעמוד המידע  
המקנה מ"מ המוגדר ב"ה זהו נ"מ א זקנה המידע  
מקנה מ"מ א זקנה (פ"ס) המוגדר ב"ה זהו נ"מ א זקנה  
מ"מ א זקנה (פ"ס) המוגדר ב"ה זהו נ"מ א זקנה

הנסח עמוד פ"ס א זקנה המוגדר ב"ה זהו נ"מ א זקנה  
מ"מ א זקנה (פ"ס) המוגדר ב"ה זהו נ"מ א זקנה  
מ"מ א זקנה (פ"ס) המוגדר ב"ה זהו נ"מ א זקנה

מקנה מ"מ א זקנה (פ"ס) המוגדר ב"ה זהו נ"מ א זקנה  
מ"מ א זקנה (פ"ס) המוגדר ב"ה זהו נ"מ א זקנה







Handwritten text on the left margin, partially visible and oriented vertically.

Handwritten text in the top section of the page.

Handwritten text in the middle section of the page.

Handwritten text in the lower middle section of the page.

Handwritten text in the bottom section of the page.

Handwritten text in the bottom section of the page.

Handwritten text in the bottom section of the page.

Handwritten text in the bottom section of the page.

Handwritten text in the bottom section of the page.



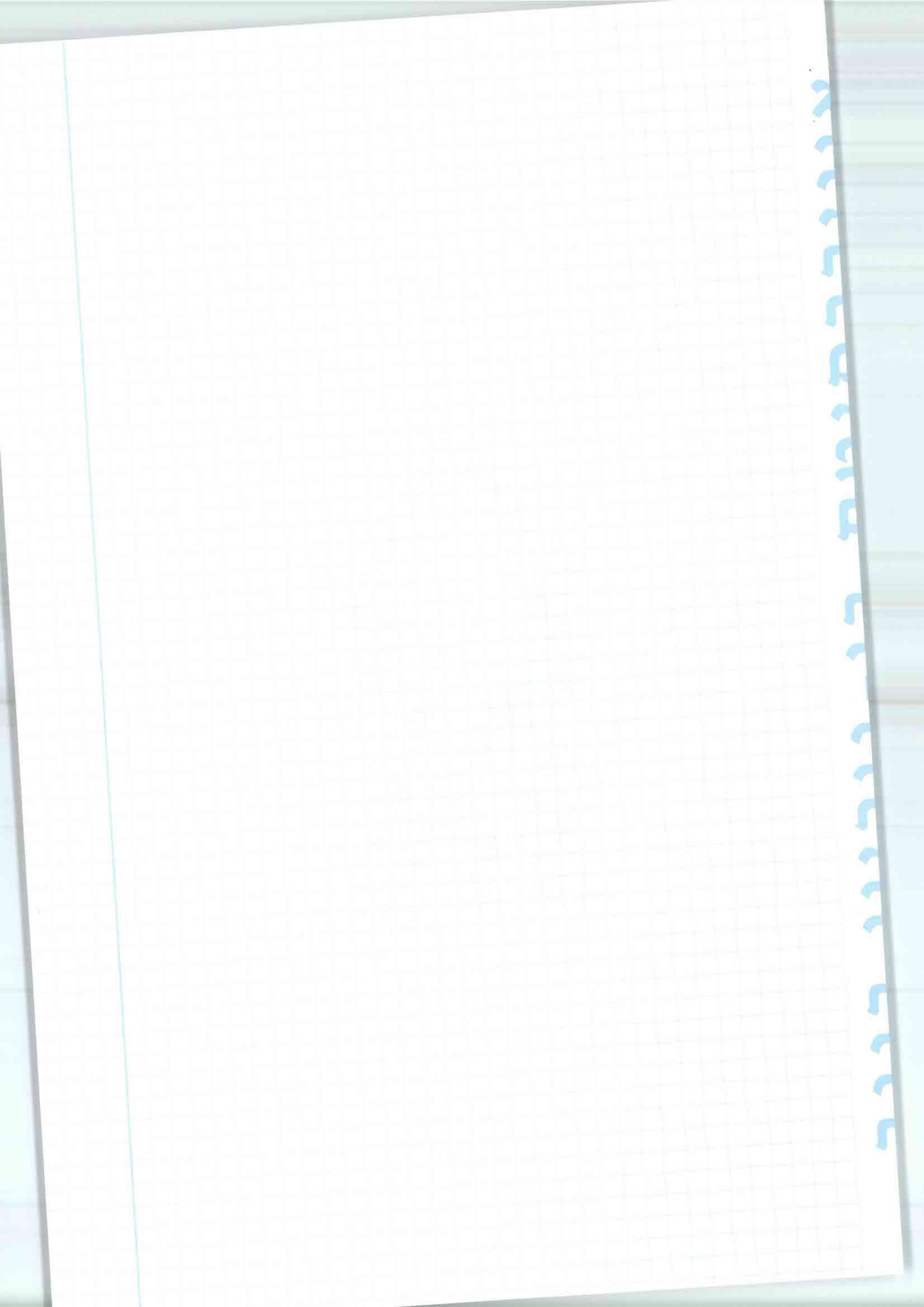


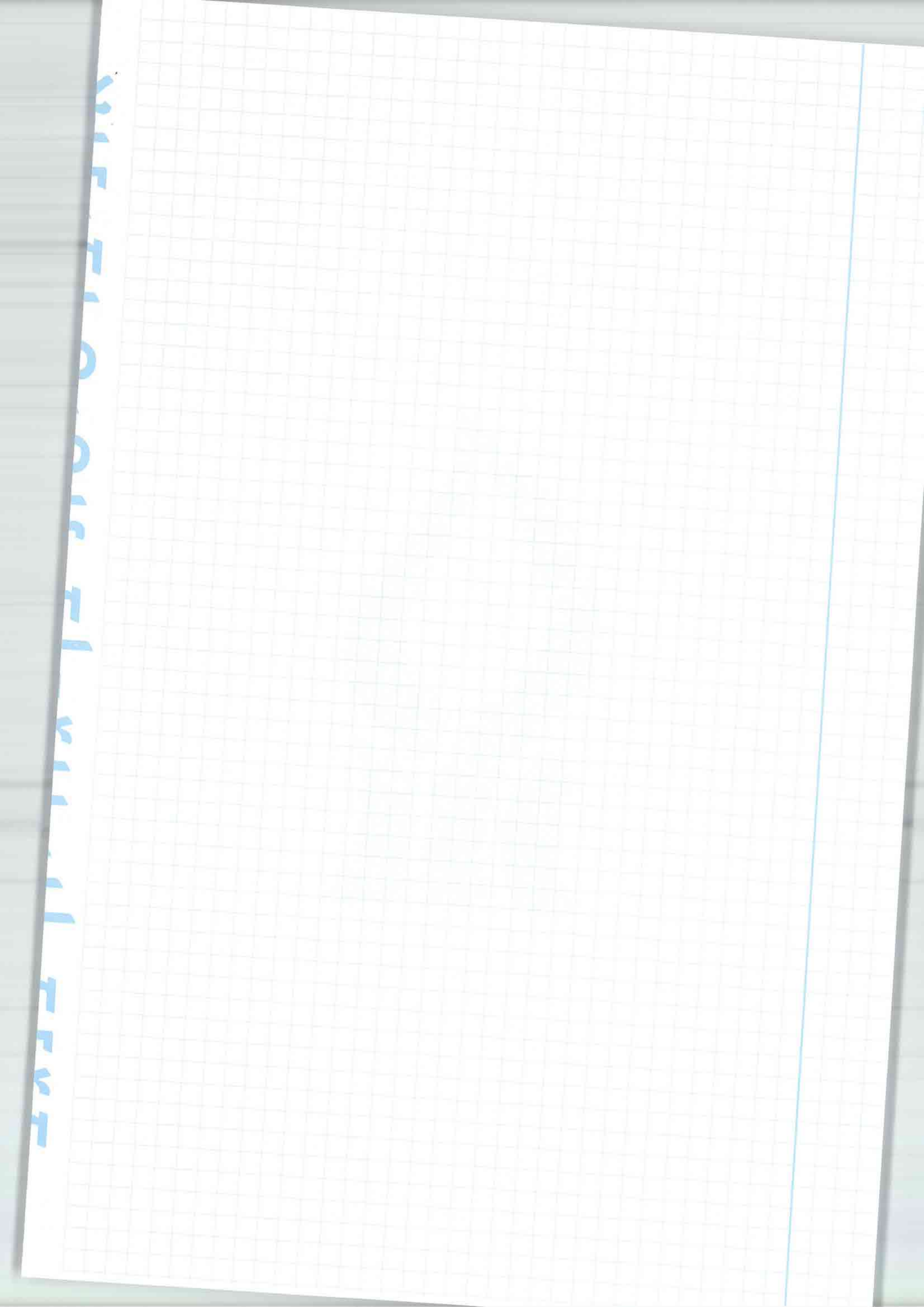


Mathematics











2.5.1.1

$$\text{dist}(\text{word1}, \text{word2}) = \sum_{i=1}^{\min(|w_1|, |w_2|)} \text{dist}(w_1[i], w_2[i])$$

word

d o g  
d o b

1  
[apple  
[apple]

$$\text{dist}(w_1, w_2) = \frac{\sum_{i=1}^{\min(|w_1|, |w_2|)} m - \text{dist}}{m \cdot \min(|w_1|, |w_2|)} + (1 - d) \frac{\min(\quad)}{\max(\quad)}$$

$$d = 0.9$$

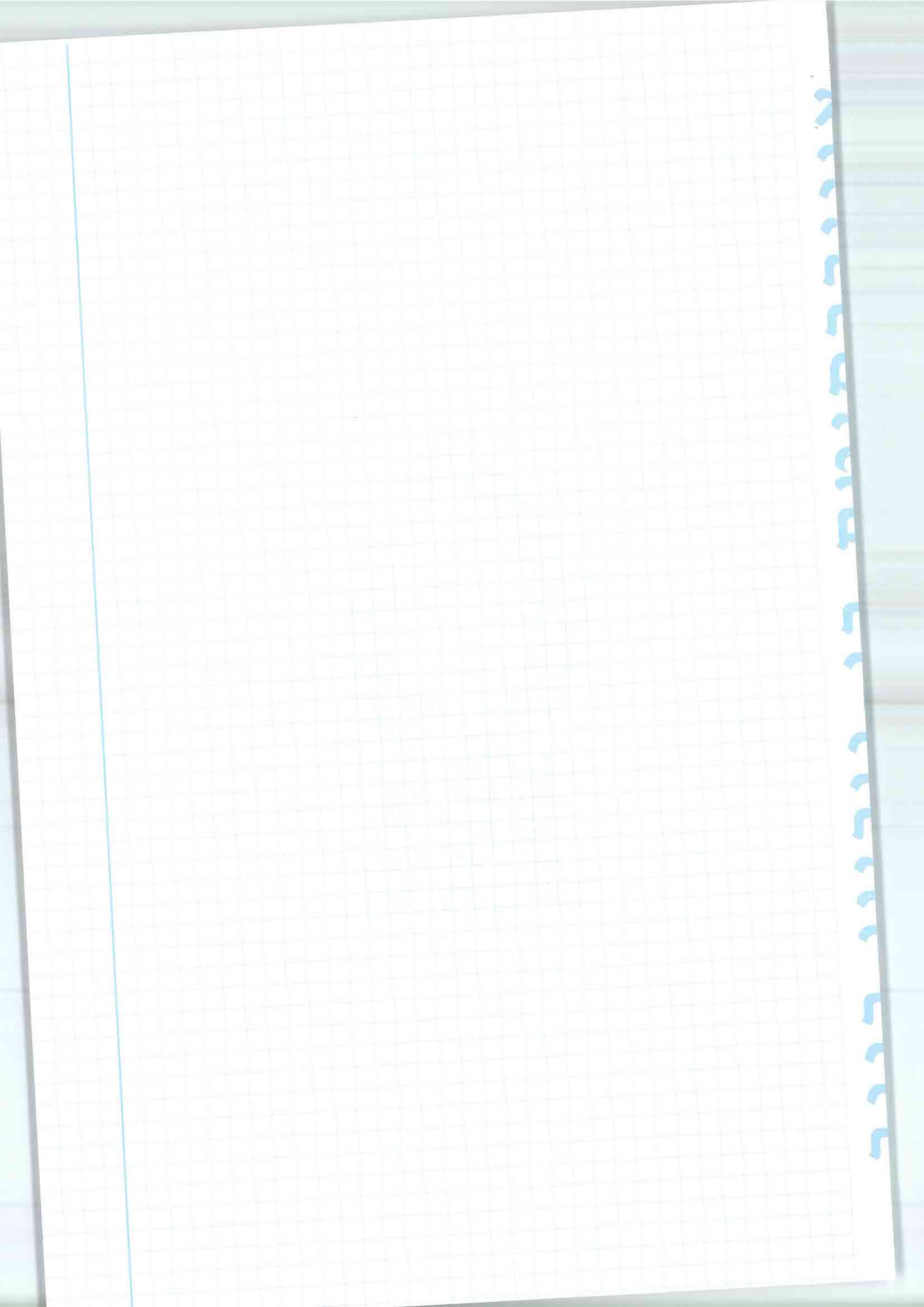
Handwritten text in blue ink, likely a title or header, partially visible on the left edge of the page. The text is oriented vertically and appears to be in a stylized font.







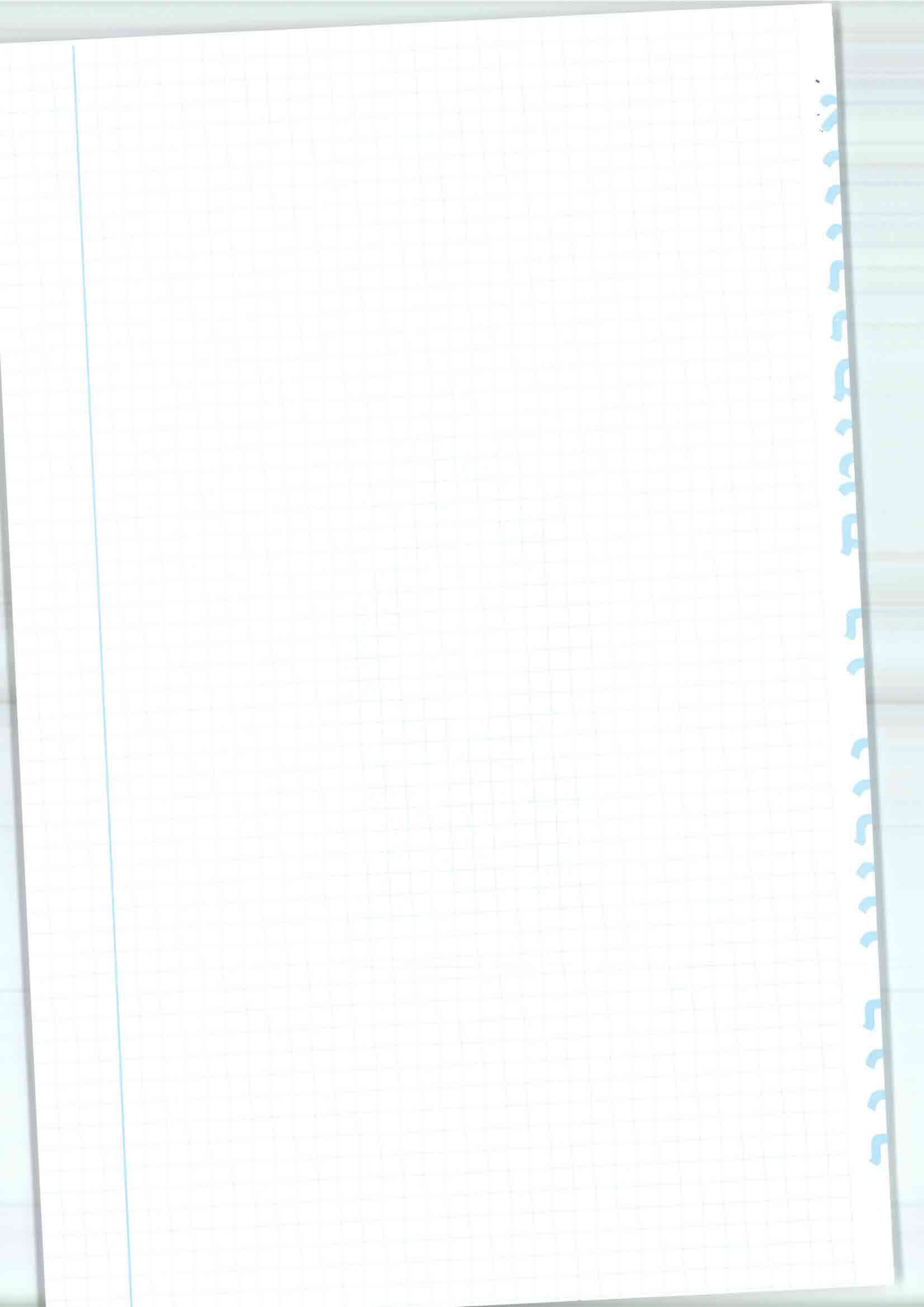






Handwritten text on the left margin, partially visible and rotated 90 degrees counter-clockwise. The text appears to be a list or index of items, with some words like "Handwritten", "List", "Index", "Handwritten", "List", "Index" visible.







Handwritten text on the left margin, partially visible and rotated 90 degrees counter-clockwise. The text appears to be a list or index of items, possibly related to a project or study, but is mostly illegible due to the angle and rotation.