

פרק תשיעי

כימות המפגשים (תכונה א)

כאן נעסוק במפגשים של צפני ELS על פני טבלאות דו-ממדיות. אנו עוקבים אחרי תכונה א (הנזכרת לעיל):

נטייה למפגש בין צפני ELS המייצגים "ביטוי א", לבין צפני ELS המייצגים "ביטוי ב", כאשר קיים קשר מושגי מובהק בין "ביטוי א" ל"ביטוי ב".

אנו זקוקים לכלי מדידה כדי לקבוע אם קיימת תכונה כזאת. נתקדם כדלקמן. נניח שעומד לרשותנו מחשב ובו טקסט נתון ותוכנה מתאימה. עתה יש להקליד למחשב צמד ביטויים – "ביטוי א" ו"ביטוי ב" - שיש קשר מושגי ביניהם.

- ראשית, המחשב סורק את הטקסט, מוצא את המד"שים של "ביטוי א" ושל "ביטוי ב" השייכים ל"אוסף" (כלומר, הם צפני ELS).
- אחר כך, המחשב מכין טבלאות דו-ממדיות עבור כל מד"ש שמצא. אם המד"ש מופיע בדילוג d , אזי בונה המחשב מן הטקסט סדרה של טבלאות הנקבעות על פי גודל דילוגו של מד"ש זה: טבלה ובה d טורים, טבלה ובה $[d/2]$ טורים, טבלה ובה $[d/3]$ טורים, וכן הלאה.

בטבלאות אלו יופיע המד"ש בצורה מכונסת. כדי שמפגש בינו לבין המד"ש של הביטוי השני יהיה 'מוצלח', צריך גם המד"ש השני להופיע בצורה מכונסת על פני הטבלה (זו הסיבה שאנו מתעניינים דווקא בטבלאות הנקבעות על ידי המד"שים: שהרי אם קיים מפגש מכונס ביניהם – הוא חייב להתרחש על אחת מן הטבלאות הללו).

בצעד הבא צריך לכמת (מלשון כמות) את המפגשים בין המד"שים על פני הטבלאות דו-ממדיות שנתקבלו בצעדים הקודמים. דבר זה נעשה על ידי הגדרה של "מידת הקרבה" בין שני מד"שים. "מידת הקרבה" נותנת ערך מספרי (כמותי) לכל מפגש בין צמד מד"שים המייצגים צמד ביטויים. אנו מבטאים בכימות זה את המאפיינים הבאים:

- טיב הנפגשים:** המפגש מתקיים בין צפני ELS, דהיינו, מד"שים שהם מינימליים בקטעים גדולים בספר.
- טיב המפגש:** המפגש נראה כמקבץ מכונס; כלומר, המד"שים מופיעים על פני הטבלה כאשר הם
 - קרובים זה לזה.
 - אינם מפוזרים.

כימות נכון של התכונה שבמעקב צריך לתאר את המרכיבים האלה ולהכיל שני חלקים: חלק המתאר את טיב הנפגשים, וחלק המתאר את טיב המפגש.

חלק הכימות, המתאר את **טיב הנפגשים**, יעניק ציון "טוב יותר" למפגשים בין מד"שים שהם מינימליים בקטעים גדולים יותר בספר. כך, למשל, הוא יעניק ציון "טוב יותר" למפגש בין צמד מד"שים, שכל אחד מהם מינימלי בכל הספר, לעומת מפגש בין צמד מד"שים שאחד מהם (או שניהם) מינימלי רק במחצית הספר.

חלק הכימות, המתאר את **טיב המפגש**, יעניק ציון "טוב יותר" למפגשים, שבהם מופיעים המד"שים כשהם קרובים זה לזה ולא מפוזרים. זה החלק המודד את "דחיסות" המפגש. המודד "לדחיסות" לוקח בחשבון את המרחק בין הביטויים, ואת מידת הפיזור שלהם – כל זאת על פני הטבלה. כדי להבהיר זאת, נשוב ונתבונן בטבלה 1 מן הפרק הקודם:

טבלה ט-1

ב	ר	א	ש	י	ת	ב	ר	א	א	ל	ה	י	מ	י	א	ת	ה	ש	מ	י	מ	ו	א	ת	א	ה	א			
ר	צ	ו	ה	א	ר	צ	ה	י	ת	ה	ה	ו	ו	ב	ו	ה	ו	ו	ח	ש	כ	ע	ל	פ	נ					
י	ת	ה	ו	מ	ו	ר	ו	ר	ו	ה	ל	ה	י	מ	י	מ	ר	ח	פ	ת	ע	ל	פ	נ	י	ה	מ			
י	מ	ו	י	א	מ	ר	א	ל	ה	י	מ	י	ה	י	א	י	ה	י	א	ו	ר	ו	ר	ו						
י	ר	א	א	ל	ה	י	מ	י	ה	א	ת	ה	א	ו	ר	כ	י	ט	ו	ב	ו	י	ב	ד	ל	א	ל			
ה	י	מ	ב	י	נ	ה	א	ו	ר	ו	ב	י	נ	ה	ח	ש	כ	ו	ק	ר	א	א	ל	ה						
י	מ	ל	א	ו	ר	י	ו	מ	ו	ל	ח	ש	כ	ק	ה	א	ל	י	ל	ה	ו	י	ה	י	ע					
ר	ב	ו	י	ה	י	ב	ק	ר	י	ו	מ	א	ח	ד	ו	י	א	מ	ר	א	ל	ה	י	מ	י					
ה	י	ר	ק	י	ע	ב	ת	ו	כ	ה	מ	י	מ	ו	י	ה	י	מ	ב	ד	י	ל	ב	י	נ					
מ	י	מ	ל	מ	י	מ	י	ו	י	ע	ש	א	ל	ה	י	מ	א	ת	ה	ר	ק	י	ע	ו	י	ב				
ד	ל	ב	י	נ	ה	י	מ	מ	י	מ	א	ש	ר	מ	ת	ח	ת	ל	ר	ק	י	ע	ו	ב	י	נ				
מ	י	מ	א	ש	ר	מ	ע	ל	ל	ר	ק	י	ע	ו	י	ה	י	כ	נ	ו	י	ק	ר	א	א	ה				
ל	ה	י	מ	ל	ר	ק	י	ע	ש	מ	י	מ	ו	י	ה	י	ע	ר	ב	ו	י	ה	י	ב	ק					
ר	י	ו	מ	ש	נ	י	ו	י	א	מ	ר	א	ל	ה	י	מ	י	ק	ו	ו	ה	מ	י	מ						
ת	ח	ת	ה	ש	מ	י	מ	א	ל	מ	ק	ו	מ	א	ח	ד	ו	ת	ר	א	ה	ה	י	ב	ש					
ה	ו	י	ה	י	כ	י	כ	נ	ו	י	ק	ר	א	א	ל	ה	י	מ	ל	י	ב	ש	ה	א	ר	צ				

חץ אחד מסמן את המרחק הקצר ביותר בין המד"ש המינימלי של המלה "הא/להים" למד"ש המינימלי של המלה "בוראכם". חץ שני מסמן את המרחק בין שתי אותיות עוקבות של המד"ש "בוראכם". חץ שלישי מסמן את המרחק בין שתי אותיות עוקבות של המד"ש "הא/להים". כל המרחקים נמדדים על פני הטבלה. למשל, המרחק בין שתי אותיות עוקבות של המלה "הא/להים", הוא אות אחת בדיוק. הכימות יעדיף מפגשים שבהם שלושת המרחקים הללו הנם קטנים. כלומר, המפגש – "דחוס".

נספח א2 מתאר כיצד נבנתה "מידת הקרבה" על פי עקרונות אלה בדרך פשוטה וטבעית. באמצעות "מידת הקרבה" של זוג מד"שים אנו מקבלים ערך מספרי בשביל כל מפגש בין צמד מד"שים. ככל שהמד"שים מינימליים "יותר", קרובים "יותר" ומכונסים "יותר" – כלומר, המפגש מוצלח "יותר" – יהיה הערך המספרי של "מידת הקרבה" גבוה יותר.

כל שתיארתי עד כה, אינו אלא הצעד הראשון בדרך לחישוב הסתברות המפגשים. עד כאן תיארתי בצורה איכותית את בניית "מידת הקרבה" עבור מפגש מסויים בין שני מד"שים, כי היא

¹ מינימלי בקטע שאורכו 3/4 מספר בראשית.
² מינימלי על כל ספר בראשית.

המשקפת את ההתרשמות האינטואיטיבית שלנו מן המפגשים שכבר ראינו, ואף עתידים אנו לראות, בספר זה. אין בדעתי להאריך כאן בעניינים טכניים, שמקומם בנספח. נספח א2 עונה על השאלות הבאות:

- כיצד מגדירים "מידת קרבה" בין ביטויים?
- כיצד מעריכים את הסתברות המפגש בין שני ביטויים?

כפתרון לשתי השאלות הללו פותחה "מידת הקרבה המכוילת" (המוגדרת בנספח א2). זו היא פונקציה המושתתת על "מידת הקרבה" הנזכרת לעיל. "מידת הקרבה המכוילת" בין שני ביטויים זה מספר בין 0 ל-1. ככל שמופעי הביטויים כצפני ELS מינימליים בקטעים גדולים יותר, וככל שהם קרובים יותר זה לזה ופחות מפוזרים על פני הטבלאות, אזי "מידת הקרבה המכוילת" קטנה וקרובה יותר ל-0. לעומת זאת "מידת הקרבה המכוילת" של זוג ביטויים קרובה ל-1 כאשר צפני ה- ELS שלהם רחוקים ומפוזרים במיוחד. מידה זו היא הסתברות המפגשים בין שני ביטויים.

לא נותר אלא להפוך פתרונות אלה לתוכנת מחשב. יואב רוזנברג התנדב להתמודד עם משימה קשה זו. היה עליו להקדיש זמן רב כדי לבנות תוכנה יעילה, המתחשבת במגבלות החמורות של אמצעי החישוב שעמדו לרשותנו – והוא עשה זאת בהצלחה. בסוף שנת התשמ"ה (סתיו 1985 למניינם) היתה בידינו תוכנה, המסוגלת למדוד את המפגשים השייכים לתכונה א: מפגשים בין צפני ELS לבין עצמם. אתה מקליד למחשב צמד ביטויים, וכפלט מתקבל מספר שהוא ערכה של "מידת הקרבה המכוילת". כאמור, מספר זה מודד את ההסתברות למפגשים של זוג הביטויים כצפני ELS.

הערה: יש להבהיר, בהנחה שהביטויים אכן מופיעים בספר בראשית בדילוג שווה של אותיות, תמדוד התוכנה מהו הסיכוי שהמפגש יהיה גם "מינימלי" וגם "דחוס". אך התוכנה איננה מודדת מהו הסיכוי, שהביטויים יופיעו בדילוג שווה בראשית. ואם יתמה הקורא מדוע איננו מתחשבים בכך, אזכיר עקרון ידוע: הגישה הנכונה לחקר תופעה נעלמה היא דווקא לבדד תכונות מסוימות. הנטייה למפגשים היא תכונה אחת. הופעה לא-מסתברת של ביטוי ארוך ו/או מורכב מאותיות נדירות כמד"ש – היא תכונה אחרת. יתכן מאד, שכדאי לחקור גם תכונה זאת – אך בשלב זה, חשוב יותר להפריד בין התכונות.