

נספח א5

טעות פשוטה

בנספח א2 ("מידת מפגשים") הגדרנו את "מידת הקרבה". זו מקבלת ערך גדול יותר אם המד"שים הנפגשים מינימליים יותר, קרובים יותר זה לזה ו"לא מפוזרים" במידה רבה יותר. שם גם הגדרנו את "מידת הקרבה המכילת" המודדת מה הסיכוי, של "מידת הקרבה" ערך כה גדול. כפי שהוסבר שם באריכות, "מידת הקרבה המכילת" היא הדרוג ב"מרוץ" בין המד"שים למדכ"שים (מלים בדילוגים כמעט שווים). ככל שערכה נמוך יותר, הסיכוי שהמפגשים אכן נוצרו במקרה – קטן יותר. להפתעתנו, הסטטיסטיקאי הנודע, פרופסור פרסי דיאקוניס כתב במכתבו הראשון (ראו בפרק טז), כי לא הצליח להבין מדוע נזקקנו להגדיר את "מידת הקרבה המכילת". דיאקוניס חשב בטעות, כי ניתן להשוות ישירות את "מידת הקרבה" של זוג ביטויים אחד ל"מידת הקרבה" של זוג ביטויים אחר. מכאן הסיק, כי אפשר ליישם שיטות סטטיסטיות סטנדרטיות המבוססות על השוואה, ישירות על קבוצת ערכי "מידת הקרבה" שנתקבלה בניסוי הגדול. לכן, הוא הציע באותו מכתב מבחן מסוים המבוסס על השוואה ישירה של "מידת הקרבה" של זוגות ביטויים. כפי שנראה להלן, זו טעות חמורה. למרבה האירוניה, העובדה שאנחנו לא עשינו כך – עוררה בו חשד, ששיטת ההשוואה לדילוגים כמעט שווים, "נתפרה" כדי לייצר מובהקות שאינה אמיתית!...

מבקר אחר של מחקרנו, פרופסור אברהם הסופר (אף הוא סטטיסטיקאי), חזר על אותה שגיאה כעשר שנים מאוחר יותר, באופן בלתי תלוי, ואף הוסיף נופך משלו. הוא הצביע, כי לכאורה הגדרת "מידת הקרבה המכילת" מניבה מוזרויות ביחס להגדרת "מידת הקרבה". לדוגמה הוא מביא "פרדוקס":

"מידת הקרבה המכילת" של זוג הביטויים "הגאון" – "ט"ו ניסך" שווה ל-0.076. פירושו של דבר, כי קיים סיכוי קטן למדי לקבל במקרה מפגש כל כך מוצלח. לעומת זאת, "מידת הקרבה המכילת" של זוג הביטויים "רבי משה" – "ב"ח איר" היא 0.4. כלומר, היא גרועה פי 5.3 מזו של הזוג הראשון, ומצביעה שהמפגש של הזוג הראשון טוב מזה של השני. אבל, אם בוחנים את "מידת הקרבה" מקבלים תמונה הפוכה: "מידת הקרבה" של הזוג השני, היא טובה פי 4.5 מזו של הזוג הראשון – ומצביעה שהמפגש של הזוג השני טוב מזה של הראשון!

ננסה להבהיר את שגיאתו של דיאקוניס ואת טעות הסופר באמצעות משל. היוונים הקדמונים רצו להוכיח את עליונותם בריצה. לטענתם, עליונות הרץ היווני היא תכונה מולדת, ולכן היא קיימת בכל הגילים. הם יזמו סדרה של 99 תחרויות ריצה.

- בכל מרוץ היו אמורים להשתתף 125 רצים: אחד מהם יווני ושאר 124 הרצים הם בני עמים ובני שבטים אחרים.

- המרוצים התנהלו לפי שנתונים: במרוץ הראשון השתתפו רצים בני שנתיים, במרוץ השני – בני שלוש שנים, בשלישי – בני ארבע, וכן הלאה, עד כי במרוץ ה-99 השתתפו רצים בני 100 שנה.
- לכל אחד מן הרצים המשתתפים נרשמה "מידת המהירות" באמצעות שעון חול משוכלל.
- בכל מרוץ דרגו את הרצים לפי סדר הגעתם למטרה. ואז, הגדירו את "מידת המהירות המכילת" של המשתתף היווני כך: הדרוג שלו, מחולק במספר המשתתפים בפועל במרוץ. למשל, אם היווני הגיע למקום השני מתוך 109 משתתפים בפועל (16 רצים לא התייצבו לתחרות או מתו בדרך), אזי "מידת המהירות המכילת" שלו היא $2/109$.
- תוצאות האולימפידה הזאת היו, אם כן, 99 מספרים הקטנים או שווים ל-1. מארגני האולימפידה צריכים היו לקבוע, אם התוצאות אכן מאשרות את עליונות היוונים כטענתם – ובאיזו רמת מובהקות סטטיסטית.
- ניתוח התוצאות יכול להיעשות על ידי "המידות הכוללות לקרבה", שהוגדרו בנספח א3 ("קרבה כוללת"). אך לא נעסוק בכך עכשיו, רק נדווח כי באותו מעמד אכן הוכחה טענת היוונים במובהקות ניכרת, ובטקס נעילת האולימפידה הוכרז על כך בחגיגות. והנה, קם הנציג המצרי – בכיר הנציגים הזרים – וערער על התוצאות ועל המסקנות. לטענתו היה פגם בסיסי בהגדרת "מידת המהירות המכילת" – פגם הגורם לפרדוקסים משונים. לדוגמה, המצרי הצביע על כך, שבתחרות לבני 33 שנים, הגיע הרץ היווני רק למקום ה-50 מתוך 125 רצים, כך ש"מידת המהירות המכילת" שלו היא $50/125$. לעומת זאת, בתחרות לבני 99 הגיע הרץ היווני למקום הראשון מתוך 101 רצים, כך ש"מידת המהירות המכילת" שלו היא $1/101$, כלומר, הישג שהוא טוב יותר פי 40 מזה של הרץ בן ה-33. "הרי זה אבסורד!" – צעק המצרי בהתרגשות – "כל אחד יודע שהמצב הפוך: 'מידת המהירות' של הרץ בן ה-33 טובה פי 50 מזו של בן ה-99!"
- מארגני האולימפידה הסבירו למצרי הנרגש, כי טעות בידו. לכל מרוץ דרגת קושי ייחודית, התלויה בכושר הגופני של המתחרים. יש משמעות ל"מידת המהירות" של הרצים ביחס לרצים "דומים": היא נועדה לבדוק, אם לרץ היווני יכולת ריצה טובה במיוחד בהשוואה לרצים בני אותו שנתון, שלהם תכונות גופניות דומות. היא קובעת את הדרוג, הדרוש להגדרת "מידת המהירות המכילת", שהיא המדד ל"יכולת המיוחדת". לעומת זאת, אין טעם להשוואה בין "מידת המהירות" של רץ בן 33, לזו של רץ בן 99, או לזו של רץ בן שנתיים: אי אפשר ללמוד מהשוואה כזו על יכולת הריצה המיוחדת של הרץ בן 33, אלא רק על השפעות הזקנה המופלגת על בן ה-99, או על ההשפעה של רגלים קצרצרות ומוטוריקה שלא בשלה בבן השנתיים. איננו יודעים אם הנציג המצרי השתכנע או לא. נניח לו ונחזור לנמשל.
- הנמשל הוא - 152 "המרוצים" שנערכו במסגרת הניסוי הגדול בין המד"שים לשאר המדכ"שים. גם כאן, לכל "מרוץ" היתה דרגת קושי ייחודית. כאן לא גיל המתחרים קבע את דרגת הקושי, אלא זוג הביטויים הנמדדים קבע זאת, כפי שנסביר מיד.
- למשל, אם "מלה א" בזוג מלים היא בת 5 אותיות שכיחות, צפוי כי המד"שים שלה יופיעו בדילוג קטן מ-10. אם "מלה ב" בזוג היא מלה בת 7 אותיות, צפוי כי המד"שים שלה יופיעו בדילוג

של כמה אלפים. במקרה כזה, קל יחסית למד"ש מינימלי בעל דילוג גדול (השייך ל"מלה ב") "ללכוד" מד"ש מינימלי בדילוג קטן (השייך ל"מלה א") במפגש קרוב. המד"ש בעל הדילוג הגדול קובע את הטבלה הדו-ממדית, ולכן יחסית הוא אינו מפוזר. המד"ש בעל הדילוג הקטן אף הוא אינו מפוזר (כי הדילוג קטן ולאורך השורה בטבלה). כך מקבלים בקלות יחסית מפגש קרוב ולא מפוזר.

נדגים זאת: הטבלה הבאה, ובה $1700=4/6801$ טורים, נקבעה על ידי המד"ש המינימלי של המלה "הסתברות" בכל ספר בראשית:

טבלה 1

ברואתו יצחקו ישמעאל בני יואל מערתה מכ
אמראה נהאשתכה ואוי כאמרת אחת יהוא
קליל אשוראני מצוה את כל נא אלה צאנו ק
ובלבו יקרבו ימי אבלאבי ואהר גהאת יעק
למקמהו י אמר לה מי עקב אחי מאי נאתמו יא
ליה יעקב ותהרבלהה ותלד לי עקב בני ותאמ
שבימה פרידי יעקבו יתנפן יהצאנא לעקדו
ספתה לביתאבי כלמה גנתא תאלה יועני
מהלכלק י אתכו ארב עמאותא ישעמו וירא
שפחותנהו ילדיהו ותשתחונו ותגשגמל
אתבנתמנן קחלנו ולנשימו אתבנתינו נתנל
ההפילגשאבי וישמע ישראל יהיו בני י
אדומבלעב נבעורו שמע ירודנה בהו ימתב
תואתכתנתהפסימ אשרעליו ויקחהו וישל
אלאתאנשי מקמהלאמראיההקדשהו אבי
נחנו ובענין ישרבית הסהרו יתנשרבית הסה
יהי בבקרו תפעמרו חו וישלחו יקרא את כל
צמצרימו יסרפרעה את טבעתו מעל ידו ית
משמר שלשתי מימו י אמר אלהי ספבי ומה
קשנו ואמלאהב יאתי ואל יכוה צגתי ולפני
הרחיקו ויוס פאמרל אשרעל לביתו קמרדפ
יולא יכלואח יולענו ותאתוכינבה לו מפנ
ימה יעקב וכל זרעו את בניו ובני בניו ו
עמדה ולפני פרעהו יברכי יעקב את פרעהו י
לו ישבע על המטה ויאמר יעקב אליו ספאלשד

האמור, הסתברות

על פני הטבלה מופיע מפגש קרוב ולא מפוזר עם מד"ש מינימלי של המלה המלאכותית "האמור", בת 5 אותיות שכיחות (להלן יוסבר איך הגענו אליה). למרות שהמפגש נראה "פוטוגני" – ואכן "מידת הקרבה" שלו גדולה – הוא איננו מובהק. בתחרות בין המד"שים למדכ"שים מתברר, כי הדרוג של מפגש זה הוא 19 מתוך 55 מתחרים, ולכן "מידת הקרבה המכוילת" של מפגש נאה זה (ראו בסוף נספח 2א) היא רק 19/55. הסיבה לכך, כי גם למתחרים – המדכ"שים – אותה דרגת קושי (במקרה זה – אותה קלות) ליצור מפגש כזה או טוב ממנו. מקרה כזה מקביל למרוץ של בני 33 במשל הנ"ל, ואולי אפילו למרוץ של בני 25.

לעומתו, נציג עתה "מרוץ" בדרגת קושי המקבילה (אולי) למרוץ של בני 70. הדבר יקרה כאשר שני הביטויים הם בני 7 אותיות.

נדגים זאת: הטבלה הבאה, ובה $680=10/6801$ טורים, נקבעה על ידי המד"ש המינימלי של המלה "הסתברות" בכל ספר בראשית:

עד כאן טיפלנו רק בדרגת הקושי של האפשרות למפגשים קרובים ולא מפוזרים בין מד"שים, המושפעת מ"מגושות" המד"שים. כדאי לציין, כי קיים גורם חשוב נוסף, המשנה את דרגת הקושי במפגשים בין ביטויים ארוכים: הכוונה למספר המד"שים המייצגים את הביטויים. לדוגמא, ביטוי בן 8 אותיות צפוי להופיע בספר בראשית מספר פעמים מועט בדילוג שווה של אותיות, בדרך כלל רק פעם אחת בלבד (במקום כ-10 מד"שים המייצגים בדרך כלל את הביטוי בסכמת המדידה שלנו – ראו נספח א2). ולכן, במרוץ הנערך עבור זוג ביטויים, שאחד מהם בן 8 אותיות, ישתתפו מד"שים מועטים. נובע מכך, כי "מידת הקרבה" – שהיא סכום התרומות של מפגשי המד"שים (ראו נספח א2) – תקטן באופן משמעותי.

אם כן טעותם הבסיסית של דיאקוניס והסופר היתה, שלא הבחינו כי "מרוצים" הנקבעים על ידי זוגות ביטויים שונים, אינם עומדים כלל באותה דרגת קושי, ולכן "מידת הקרבה" במרוץ אחד אינה בת-השוואה לזו של מרוץ אחר. לפי טעותם, ערך התוחלת של "מידת הקרבה" אינו תלוי בזוג הביטויים הנמדד. הרי זה כאילו טענו שערך התוחלת של "מידת המהירות" במרוץ מסוים באולימפידה (במשל שלנו), אינו תלוי בגיל הרצים!...

לפי דברינו עד כה, ניתן לחזות מראש, כי ערך "מידת הקרבה" במפגשים עם מלה כמו "הסתברות" יהיה גדול יחסית כשהמלה השניה קצרה, אך קטן יחסית – כשהמלה השניה ארוכה. אפשר להעמיד זאת במבחן הניסוי.

ננסה להפגיש ביטויים משלושה סוגים - A, B ו-C - עם המלה "הסתברות".

סוג A: ביטויים בני 5 אותיות השכיחות ביותר בספר בראשית. כך מובטח כי צפוי שהם יופיעו בדילוג קצר. כדי לקבל קבוצה גדולה כזאת יצרנו את הביטוי "אהוימ", המורכב מחמש האותיות השכיחות ביותר בספר בראשית, שסודרו כאן לפי מיקומן בסדר האלפבית. ניתן לסדר אותיות אלה ב-120 אופנים שונים. כיוון שאין אנו מבדילים בין דילוג קדימה לדילוג אחורה, הרי לעניין הדילוגים הביטוי "אהוימ" והביטוי "מיוהא" מהווים אותו הביטוי. לכן, סך כל הביטויים השונים הוא $120/2=60$.

סוג B: לכל אחד מן הביטויים מסוג A, נוסיף "לר" בסופו. כך שבמקום "אהוימ" – יתקבל "אהוימלר". האותיות "לר" נבחרו משום ש"ל" היא האות הששית מבחינת השכיחות בספר בראשית, ו"ר" – השביעית. כך קבלנו 60 ביטויים שונים בני 7 אותיות.

סוג C: לכל אחד מן הביטויים מסוג B, נוסיף "ב" בסופו, במקום "אהוימלר" – "אהוימלרב". האות "ב" נבחרה משום שהיא האות השמינית מבחינת השכיחות בספר בראשית. כך קבלנו 60 ביטויים שונים בני 8 אותיות.

שימו לב, כי בנית קבוצות הביטויים מאותן האותיות מבטיח, כי דרגת הקושי דומה היא עבור כל הביטויים השייכים לאותו הסוג.

בטבלה הבאה ציינתי בצד כל אחד מן הביטויים את ערכה של "מידת הקרבה" במפגש עם המלה "הסתברות". לנוחות ההצגה הוכפל כל מספר פי 100. בטור האחרון סימנתי ב"-" את המקרים שבהם הביטוי אינו מופיע בדילוג שווה.

טבלה 3

$\Omega \times 100$	C	$\Omega \times 100$	B	$\Omega \times 100$	A	
0.325	אהוילרב	1.118	אהוילרב	1.403	אהוי	1
0.453	המאילרב	1.060	המאילרב	1.810	המאוי	2
-	המאילרב	1.556	המאילרב	2.483	המאוי	3
-	המאילרב	0.726	המאילרב	6.842	המאוי	4
-	המאילרב	1.001	המאילרב	3.517	המאוי	5
0.231	המאילרב	0.663	המאילרב	3.100	המאוי	6
-	המאילרב	1.600	המאילרב	5.124	המאוי	7
0.490	המאילרב	0.842	המאילרב	6.678	המאוי	8
0.626	המאילרב	1.023	המאילרב	9.158	המאוי	9
0.161	המאילרב	1.144	המאילרב	4.038	המאוי	10
0.447	המאילרב	0.989	המאילרב	1.379	המאוי	11
0.588	המאילרב	0.702	המאילרב	1.732	המאוי	12
0.400	המאילרב	0.731	המאילרב	3.998	המאוי	13
0.423	המאילרב	1.635	המאילרב	0.601	המאוי	14
0.386	המאילרב	1.035	המאילרב	3.506	המאוי	15
0.255	המאילרב	1.077	המאילרב	4.193	המאוי	16
0.295	המאילרב	1.023	המאילרב	0.548	המאוי	17
0.217	המאילרב	1.215	המאילרב	3.260	המאוי	18
-	המאילרב	1.006	המאילרב	1.987	המאוי	19
0.351	המאילרב	0.850	המאילרב	1.608	המאוי	20
0.562	המאילרב	1.303	המאילרב	2.353	המאוי	21
-	המאילרב	1.461	המאילרב	1.488	המאוי	22
0.402	המאילרב	0.833	המאילרב	2.934	המאוי	23
0.399	המאילרב	1.392	המאילרב	10.600	המאוי	24
0.106	המאילרב	1.088	המאילרב	1.605	המאוי	25
0.192	המאילרב	0.701	המאילרב	1.898	המאוי	26
0.632	המאילרב	1.029	המאילרב	10.518	המאוי	27
-	המאילרב	1.019	המאילרב	2.775	המאוי	28
0.233	המאילרב	1.015	המאילרב	1.403	המאוי	29
1.304	המאילרב	0.969	המאילרב	4.816	המאוי	30
0.469	המאילרב	0.541	המאילרב	24.714	המאוי	31
0.399	המאילרב	1.164	המאילרב	1.523	המאוי	32
0.318	המאילרב	1.215	המאילרב	2.781	המאוי	33
0.750	המאילרב	1.625	המאילרב	1.877	המאוי	34
0.360	המאילרב	0.580	המאילרב	1.071	המאוי	35
0.381	המאילרב	1.091	המאילרב	0.906	המאוי	36
0.474	המאילרב	1.216	המאילרב	2.049	המאוי	37
0.687	המאילרב	0.700	המאילרב	3.658	המאוי	38
0.176	המאילרב	1.674	המאילרב	1.736	המאוי	39
0.123	המאילרב	0.853	המאילרב	1.982	המאוי	40
0.437	המאילרב	1.280	המאילרב	1.628	המאוי	41
0.263	המאילרב	0.774	המאילרב	4.215	המאוי	42
-	המאילרב	0.290	המאילרב	3.127	המאוי	43
-	המאילרב	1.049	המאילרב	3.962	המאוי	44
0.448	המאילרב	0.962	המאילרב	1.458	המאוי	45
0.184	המאילרב	1.335	המאילרב	2.187	המאוי	46
-	המאילרב	1.158	המאילרב	3.868	המאוי	47
0.187	המאילרב	0.919	המאילרב	3.339	המאוי	48
-	המאילרב	1.288	המאילרב	1.419	המאוי	49
-	המאילרב	1.544	המאילרב	8.261	המאוי	50

$\Omega \times 100$	C	$\Omega \times 100$	B	$\Omega \times 100$	A	
1.415	ואהימלרב	0.959	ואהימלר	6.225	ואהימ	51
-	ואיהמלרב	0.736	ואיהמלר	5.457	ואיהמ	52
0.748	ויהאמלרב	0.736	ויהאמלר	15.151	ויהאמ	53
-	ויהאמלרב	0.741	ויהאמלר	2.739	ויהאמ	54
-	יהאומלרב	0.682	יהאומלר	2.876	יהאומ	55
0.112	יהואמלרב	0.946	יהואמלר	3.518	יהואמ	56
0.277	יהאומלרב	1.076	יהאומלר	2.440	יהאומ	57
0.253	יאהמלרב	1.309	יאהמלר	2.301	יאהמ	58
0.125	יאהאמלרב	0.526	יאהאמלר	6.011	יאהאמ	59
0.361	יאהמלרב	1.306	יאהמלר	13.772	יאהמ	60
0.409		1.03		4.06		ממוצע

הטבלה ממחישה הבדל רציני וקבוע בין "מידת הקרבה" של כל סוג וסוג:

- למעט 3 מקרים (שהובלטו ברקע אפור) מתוך 60, "מידת הקרבה" של סוג A גדולה מזו של סוג B.
- למעט 3 מקרים (שהובלטו ברקע אפור) מתוך 45, בהם הוגדרה "מידת הקרבה" של סוג C, "מידת הקרבה" של סוג B גדולה יותר.
- בכל 45 המקרים, בהם הוגדרה "מידת הקרבה" של סוג C, "מידת הקרבה" של סוג A גדולה יותר.
- הממוצעים עבור כל סוג משקפים הבדל מהותי בין הסוגים.

כל "דרך סטטיסטית סטנדרטית" (כרצונו של דיאקוניס) תגלה, כי שלוש ההתפלגויות של ערכי "מידת הקרבה", שנקבעו על ידי שלושת הסוגים הנ"ל, הן זרות זו לזו, ולכן ההנחה שהניח – שגויה לחלוטין.