

21 שנות בחינת אמיר"ם בקרת איכות ומחקר

טוני גוטנטג, אבי אללוף, מיכל באומר, מרינה פרונטון ונדב פודולר
מרכז ארצי לבחינות ולהערכה

כנס 2015 של האגודה הישראלית לפסיכומטריקה (אפ"י)

שני סוגי בקרת איכות

סוג שני

בקרת איכות לבחינות המועברות לקבוצות אוכלוסייה קטנות במספר תאריכים רב תוך שימוש במספר רב של נוסחי בחינה. למשל, אמיר"ם, מימ"ד, IELTS, TOEFL

סוג ראשון

בקרת איכות לבחינות המועברות לקבוצות אוכלוסייה גדולות במספר קטן של מועדים קבועים תוך שימוש במספר נוסחים קטן. למשל, הפסיכומטרי, SAT, ACT

מספר נבחנים		מספר נוסחים	תאריך
ממוצע לנוסח	סה"כ		
20	120	6	5 באפריל
30	210	7	7 באפריל
15	90	6	11 באפריל

מספר נבחנים		מספר נוסחים	תאריך
ממוצע לנוסח	סה"כ		
3000	3000	1	5 באפריל
4500	9000	2	7 ביוני
4000	8000	2	11 ביולי

בקרת איכות על שני סוגי המבחנים

- לסוג הבחינות הראשון והמסורתי פותחו שיטות בקרת איכות ידועות ומוכרות (Allalouf, 2007; Kolen & Brennan, 2004).
- שיטות אלה אינן שמישות בבחינות מהסוג השני, בעיקר בשל גודלי המדגמים הקטנים ובשל מספר רב של נוסחים המועברים במקביל.
- מחקר זה עוסק בסוג השני, שהינו סוג רלוונטי במיוחד של בקרת איכות, בהינתן שעולם הבחינות מתקדם בקצב מהיר לכיוון זה.
- ניתן לראות בבקרת איכות מתמשכת כעין "רפואה מונעת" תקלות עתידיות.

מטרת המחקר הנוכחי

■ בקרת איכות על ציוני בחינה היא הכרחית עבור הגופים העורכים מבחנים, הגופים המשתמשים בציונים והנבחנים, במיוחד כשמדובר במבחנים **עתירי סיכון** (high stake).

1. התבוננות "**היסטורית**" ארוכת טווח על הציונים ואיכותם, שכן בחינת אמיר"ם מועברת מזה שני עשורים.

2. עריכת **בקרת איכות** על בחינת אמיר"ם המועברת באופנות העברה מתמשכת, בהסתמך על שיטות סטטיסטיות שנחקרו ופותחו לאחרונה. (Lee & von Davier, 2013 ;Schafer et al., 2011)

שיטה

שיטת המחקר

■ אוכלוסייה

- 199,674 נבחנים
- בין השנים 1992 (אוגוסט) עד 2013 (דצמבר)
- ב-27 נוסחים (מתוכם ניתן לקבץ 10 נוסחים)

מבחני המרכז הארצי הבוחנים רמת אנגלית 2014

אמיר"ם	אמי"ר	בחינת מימ"ד	הבחינה הפסיכומטרית	
	מיון רמת אנגלית במוסדות לימוד אקדמיים.	בחינת ידע למועמדים למכינות הקדם-אקדמיות; מיון למסלולים.	כלי לחיזוי הצלחה אקדמית; מיון לחוגים; מיון רמת אנגלית.	מטרה
	אנגלית	עברית, מתמטיקה, אנגלית	מילולי, כמותי, אנגלית	תכנים
ממוחשב אדפטיבי	נייר ועפרון	ממוחשב	נייר ועיפרון	אופנות

בחינת אמיר"ם

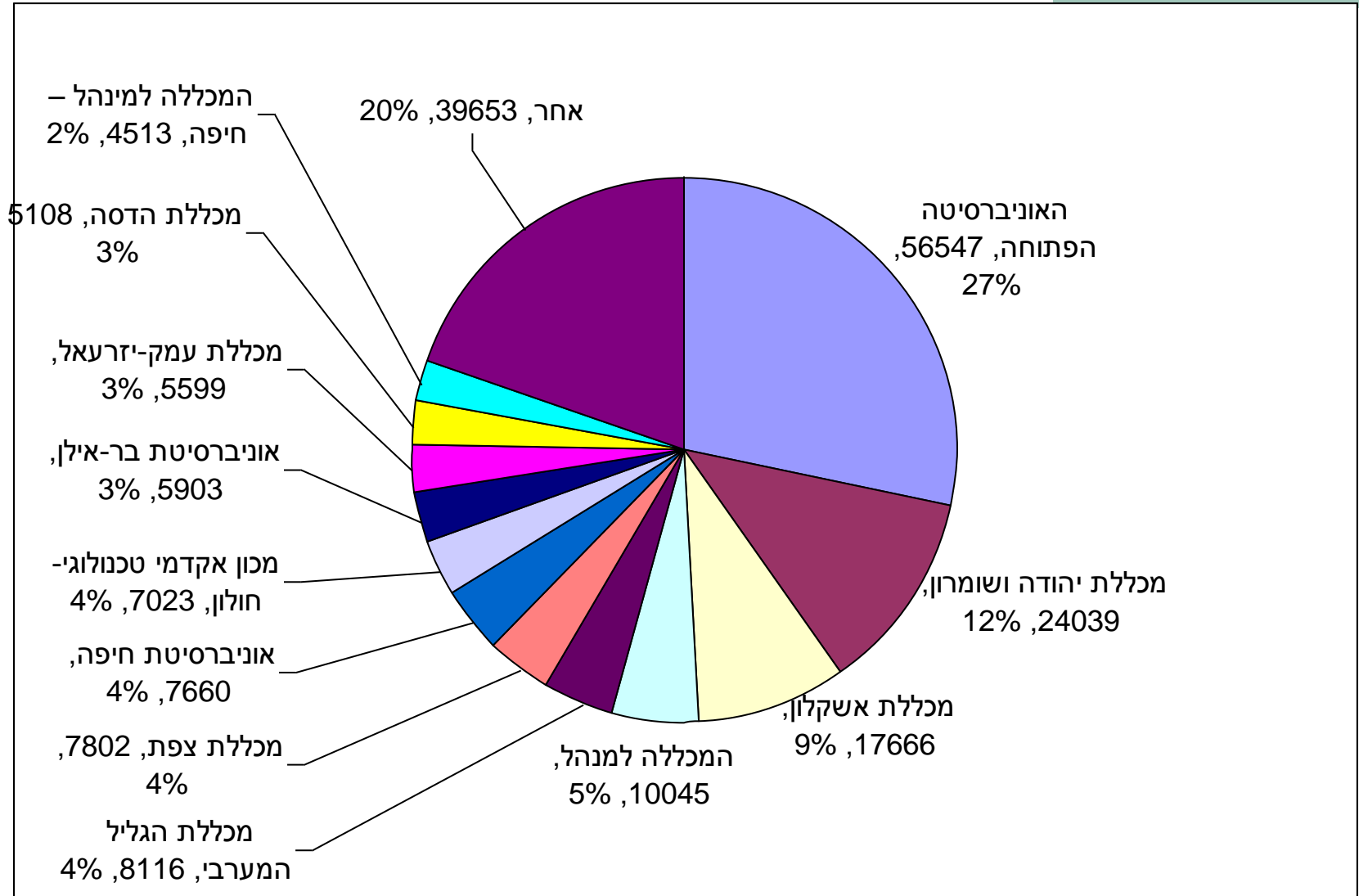
מבנה

- בבחינה שאלות השלמת משפטים, ניסוח מחדש והבנת הנקרא.
- שאלות ברירה עם ארבע תשובות מוצעות.
- אורך הבחינה הממוצע כשעה.
- בחינה **ממוחשבת**.
- בחינה **אדפטיבית**, לפי מודל IRT תלת פרמטרי*.
- רמת היכולת מומרת לציון סופי בבחינה באמצעות **טבלת המרה***.
- התאמת קושי השאלות ליכולת הנבחן מאפשרת לאמוד את רמת הידע שלו בעזרת **מספר שאלות קטן** (21,28 שאלות; ממוצע 23).

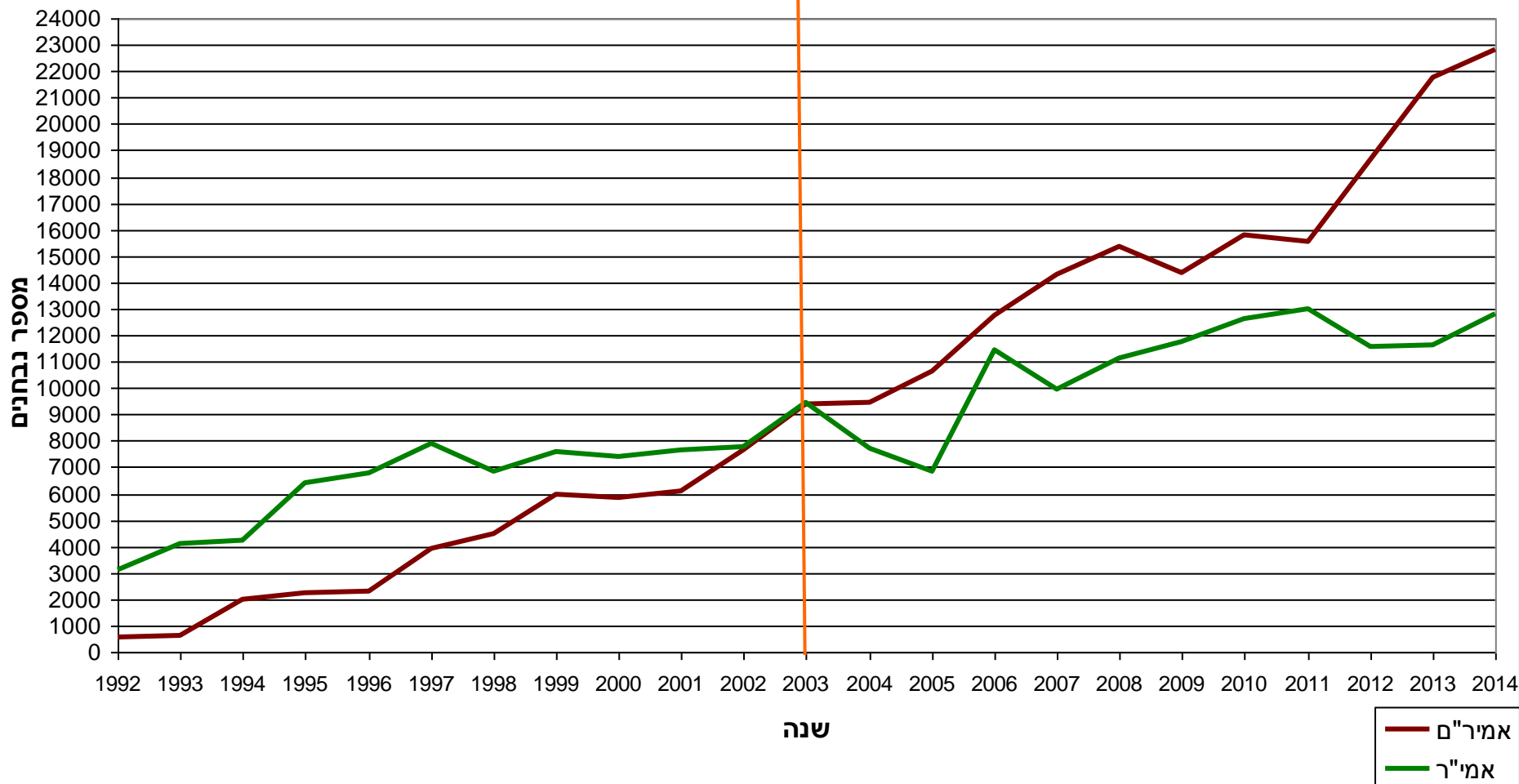
* ידובר בהמשך.

תוצאות

שכיחויות נבחנים על-פי מוסד אקדמי מעבר לשנים

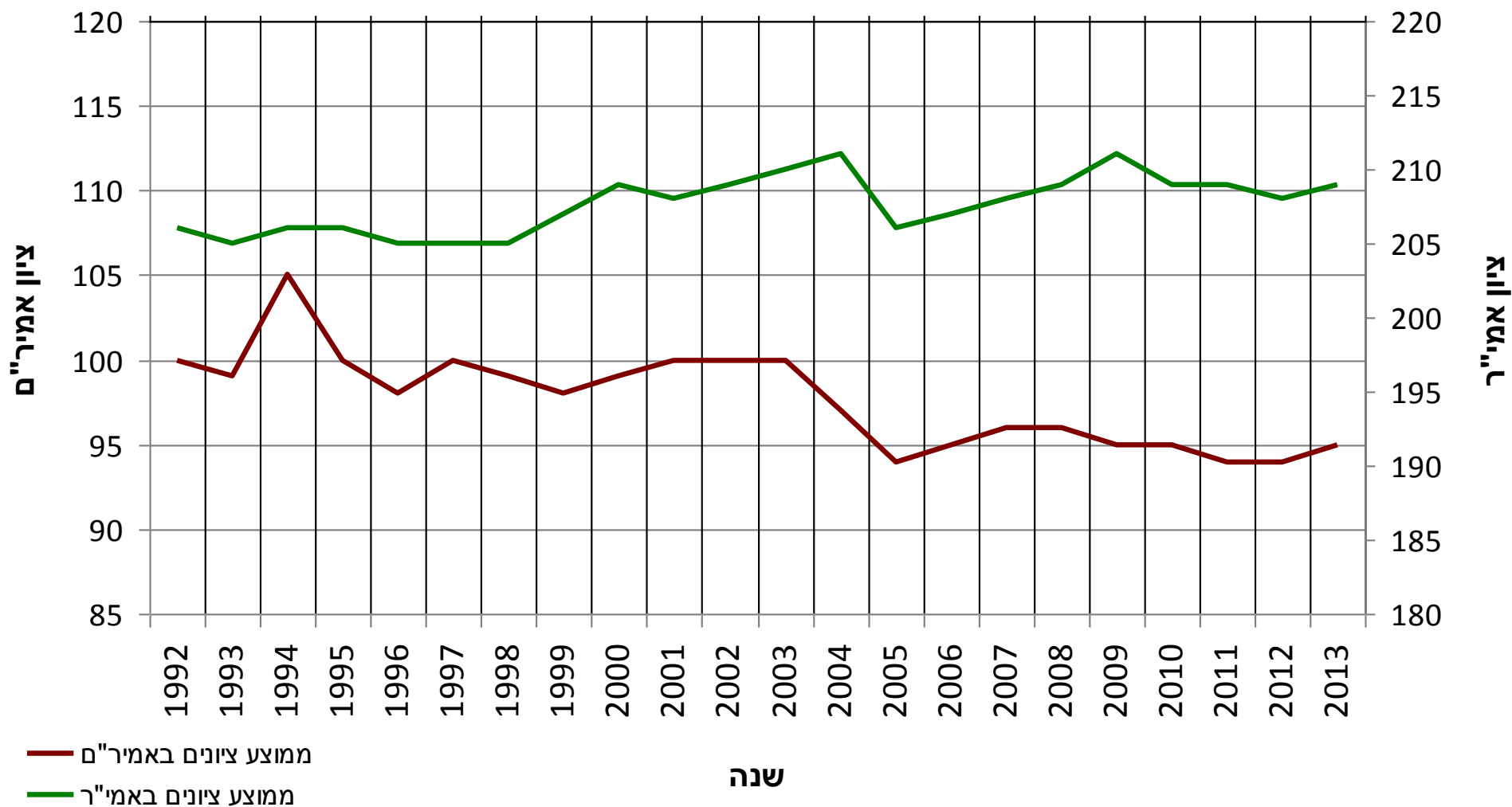


שכיחויות נבחנים באמיר"ם מול אמיר"ר לפי שנה



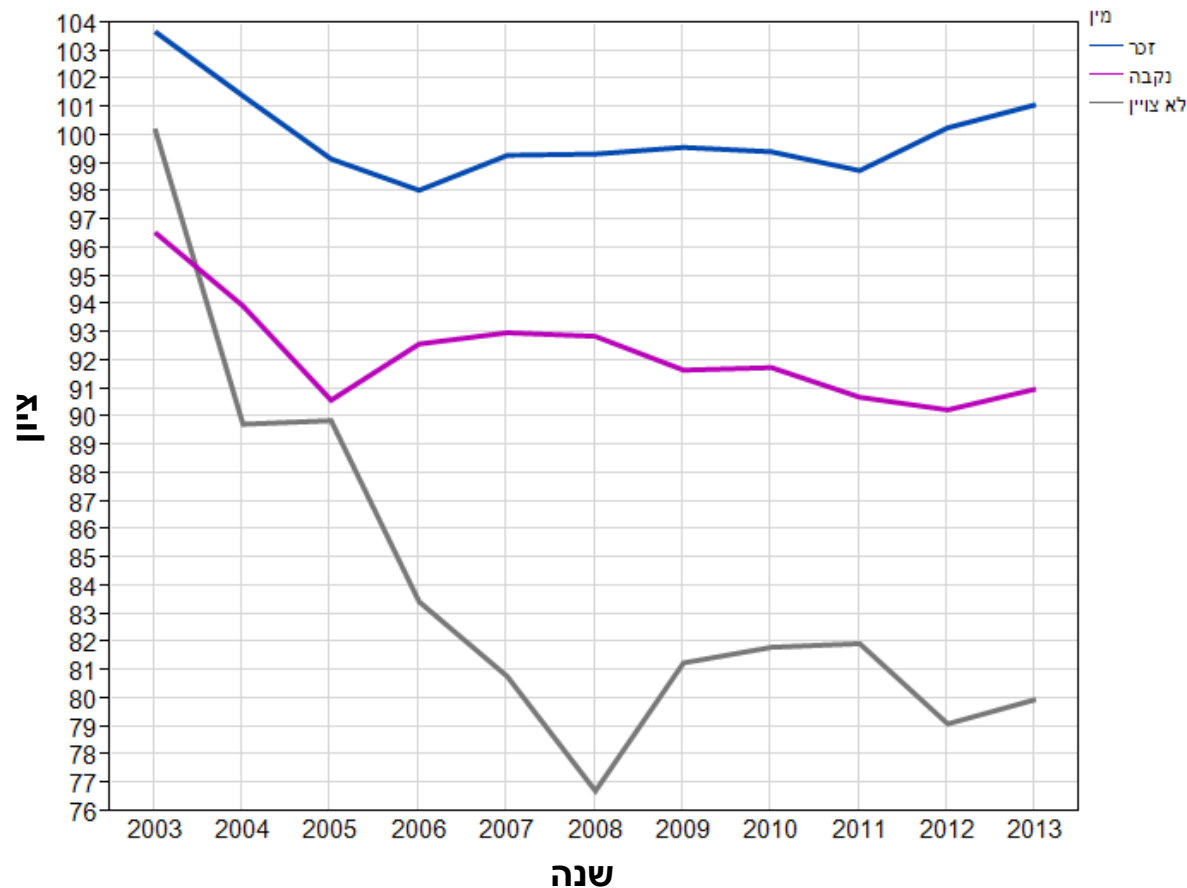
ציוני אמיר"ם לעומת אמיר"ר לפי שנה

[הבדל של 100 בסולם]



מגדר, ציון ממוצע לפי שנה 2013-2003

(לפני שנת 2003 הנבחנים
לא נתבקשו לציין מגדר)



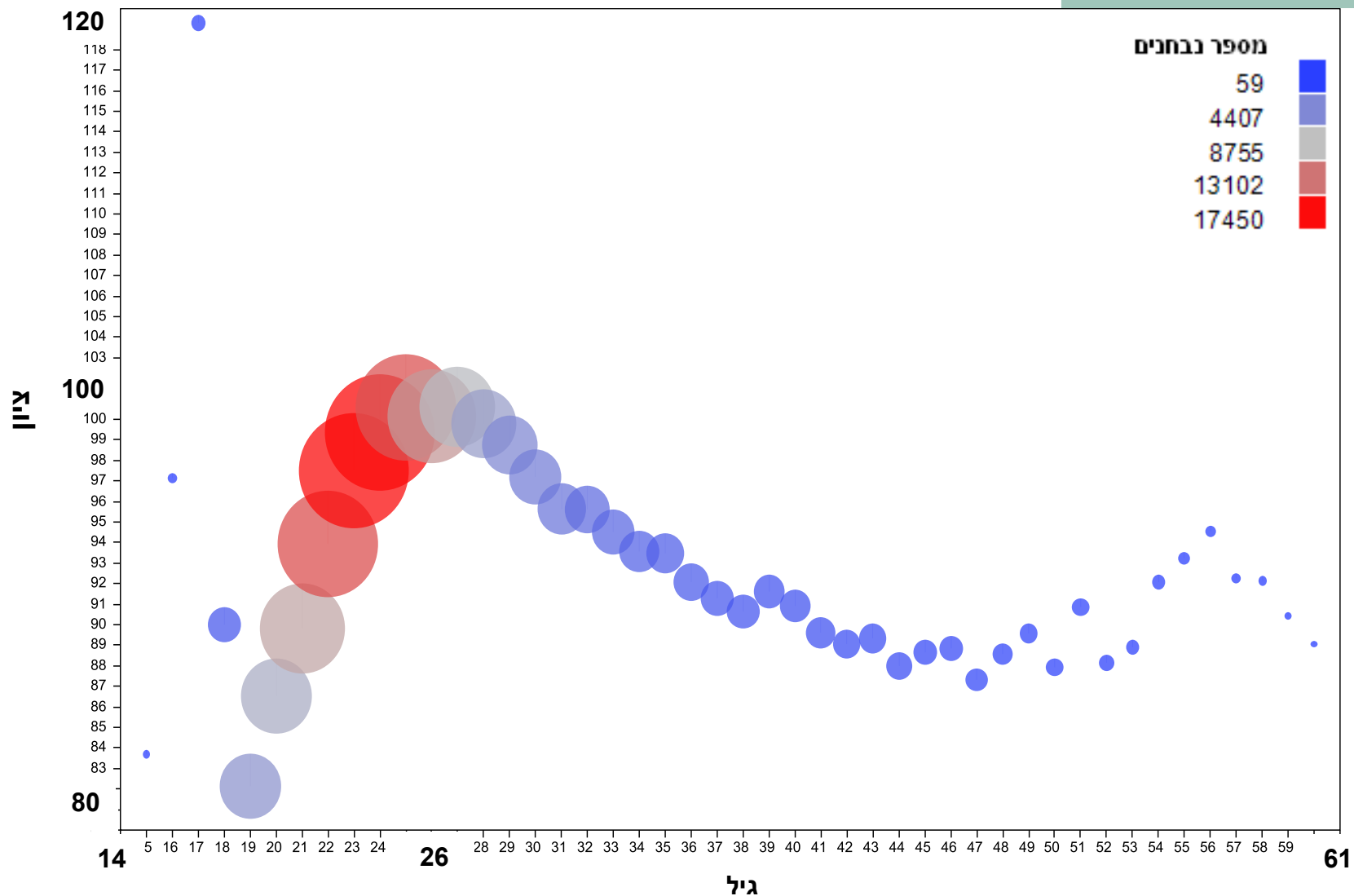
רוב הנבחנים הן נבחנות (נשים: 89,391, גברים: 66,779).

ציון של הנשים (M=92.0, SD=24.3) נמוך במובהק

מציון הגברים (M=99.9, SD=25.9) בכ-0.4 ס.ת.

(לפני שנת 2003 הנבחרים
לא נתבקשו לציין גיל)

גיל, שכיחות וציון ממוצע מעבר לשנים 2003-2013



שכיחות, ממוצע וסטיית תקן של הציון לפי נוסח מקובץ מעבר לשנים

ניתן לראות כי ציון הבחינה
הממוצע וסטיית התקן
יציבים מעבר לנוסחים
השונים, החל מהנוסח ה-3.

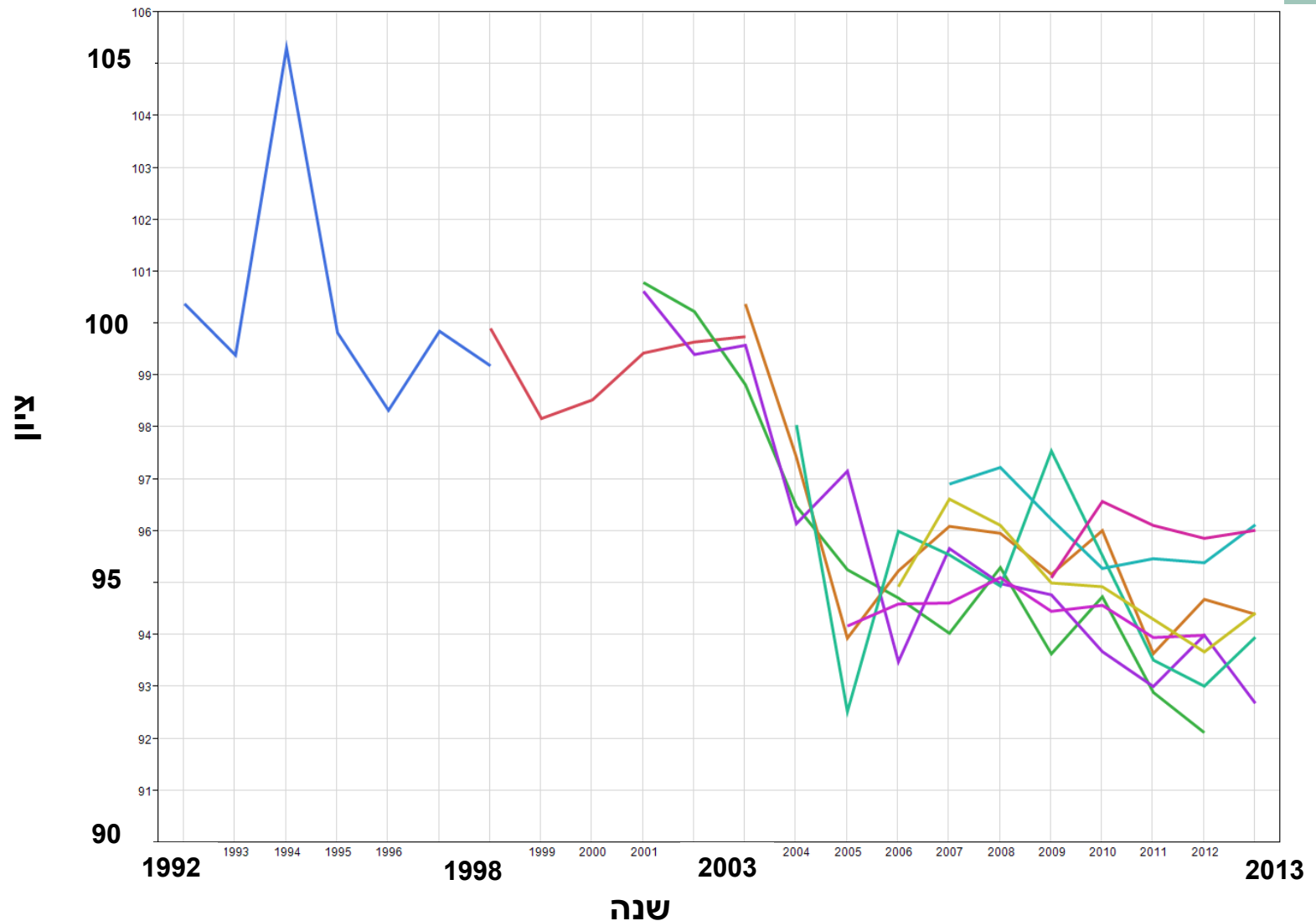
ציון		N	נוסח
ממוצע	ס.ת.		
100.1	20.4	16,072	1
99.0	20.8	24,805	2
95.8	24.4	18,049	3
95.5	24.6	20,274	4
95.6	25.0	27,397	5
94.7	25.7	22,165	6
94.5	25.5	17,570	7
95.0	25.4	19,668	8
96.0	25.6	20,149	9
96.0	26.6	13,529	10

נוסח ישן



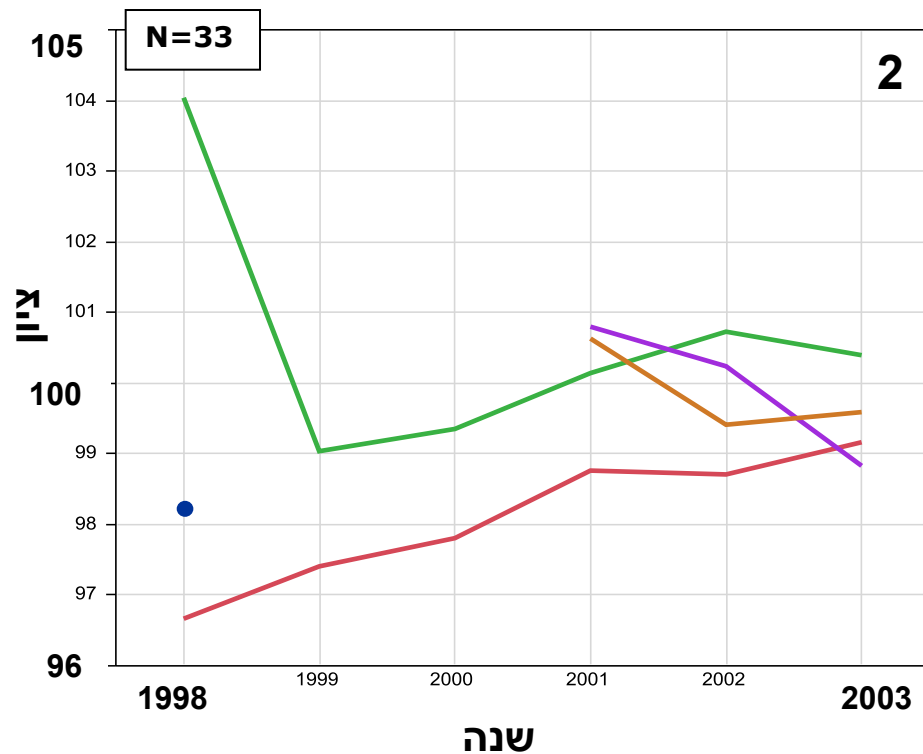
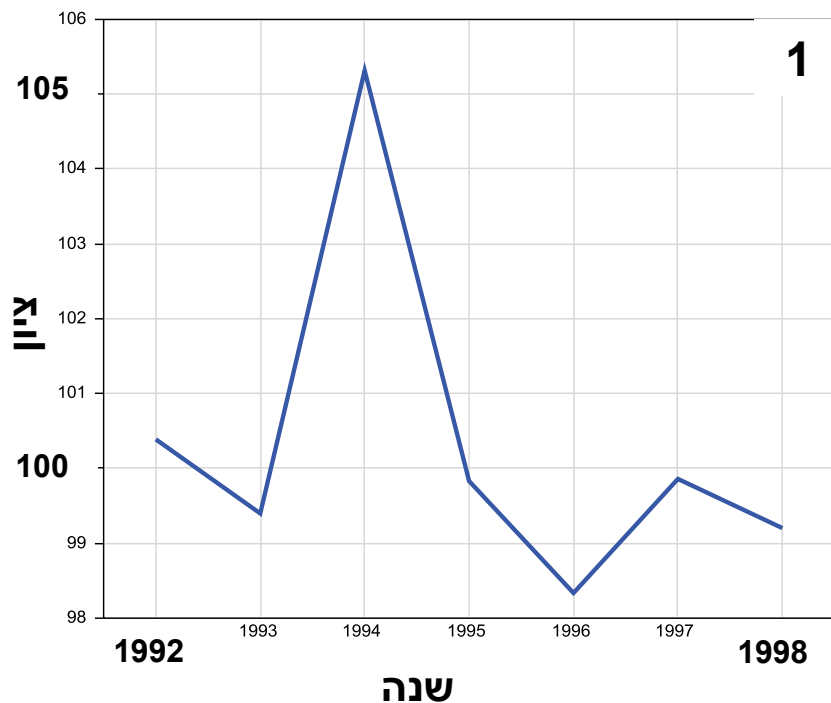
נוסח חדש

ציון ממוצע על-פני ציר הזמן לפי נוסח מקובץ 1992-2013



ציון ממוצע על-פני ציר הזמן לפי נוסח העשור הראשון

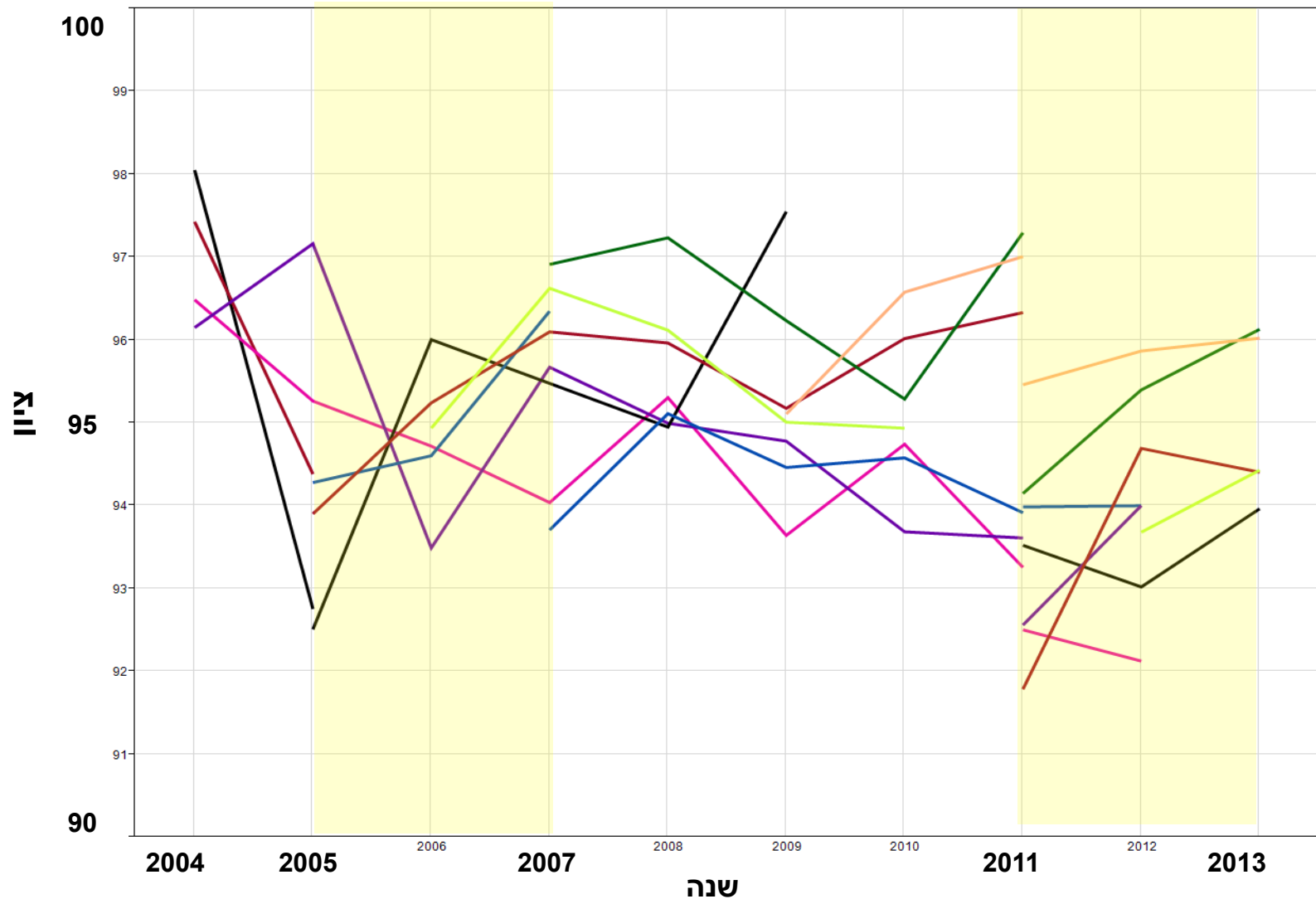
בעשור הראשון ניכרה
תנודה משמעותית בציון הממוצע
משנה לשנה.



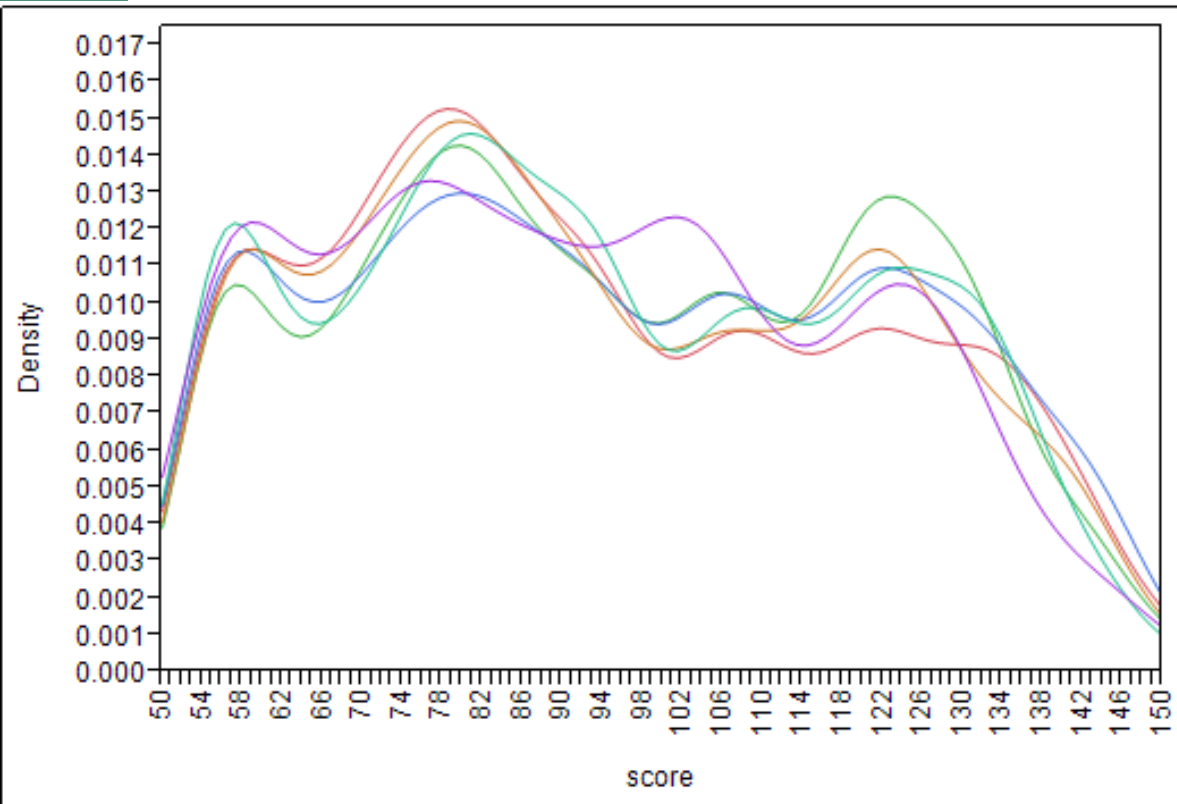
לקראת העשור השני,
כל הנוסחים המוצגים כאן
הורדו מהמחזור.

ציון ממוצע על-פני ציר הזמן לפי נוסח

העשור השני



שקילות נוסחים 2013



Level		Mean
א'	A	96.1
ב'	A B	96.0
ג'	B C	94.4
ד'	B C	94.4
ה'	C	94.0
ו'	C	92.7

נמצא הבדל מובהק בין נוסחי הבחינה, אך האפקט זעום:
 $F(5, 21707) = 7.06, p < .001, adj. r^2 = .01$

טבלאות המרה: מתטא לציון הבחינה

- הפרמטר תטא (θ) מייצג את יכולת הנבחן, במקרה זה בתחום האנגלית.
- הפרמטר מחושב על סמך מודל IRT לוגיסטי תלת פרמטרי (3PL) בשיטת הנראות המרבית.
- הטבלאות ממירות רמת יכולת (ציון תטא, לרוב בטווח של 3- עד 3) לציון סופי בבחינה (על סולם של 50 עד 150)
- טבלת ההמרה נקבעה על סמך מחקרים אמפיריים (רפ, 1995, 1996) והחלטות לגבי ערכים ספציפיים: "כללי עצירה".

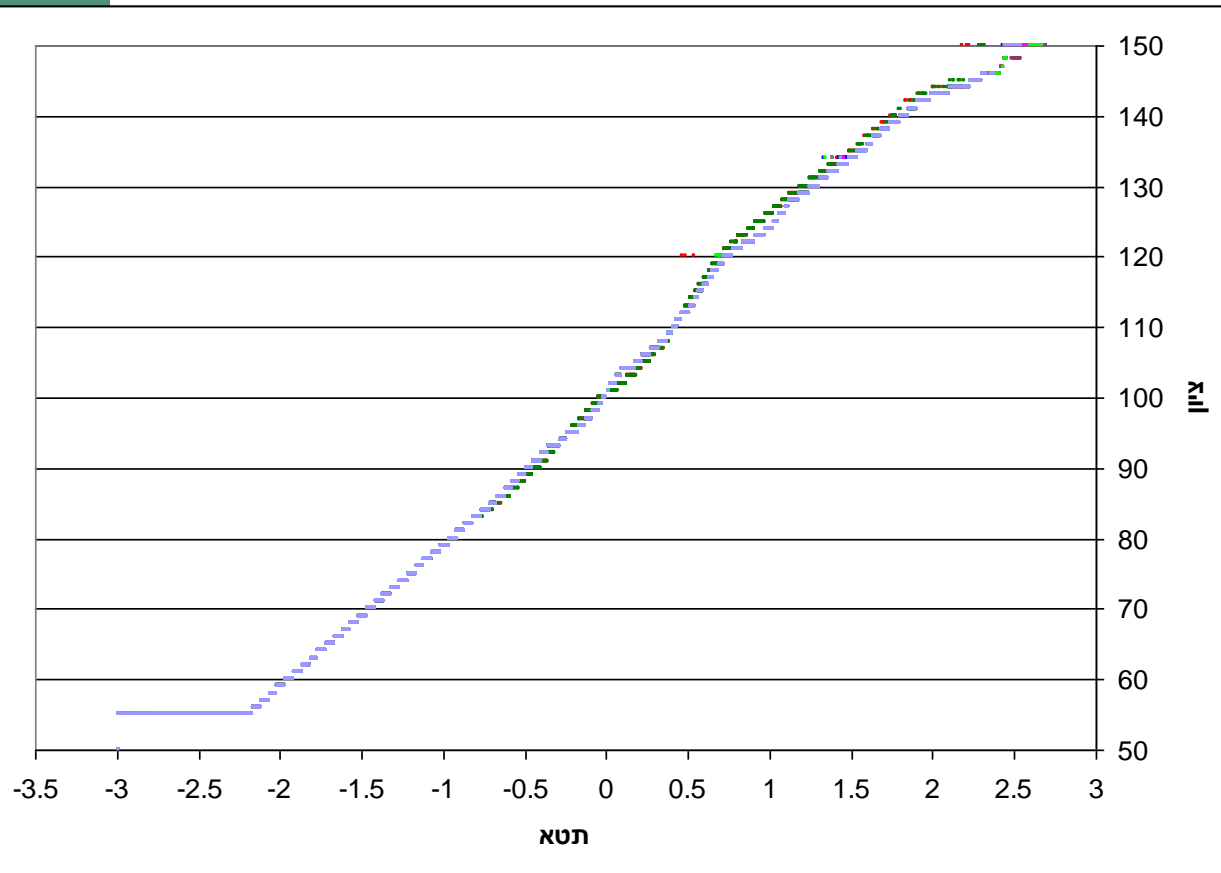
אחוזי התכנסות לפי נוסח 2012-2013

■ P-var היא טעות המדידה, נאמדת לאחר כל פריט במבחן.

התכנסות מוגדרת	אחוז הנבחנים שציונם התכנס	נוסח
כ- $P\text{-var} < .08$.	81	א'
לעיתים, המבחן יסתיים ללא	80	ג'
התכנסות בנקודת החתך	80	ה'
(אך במרחק סביר ממנה).	80	ו'
	80	ז'
	80	ח'
	79	ב'
	77	ד'

טבלאות המרה

נוסחים שהועברו ב-2013



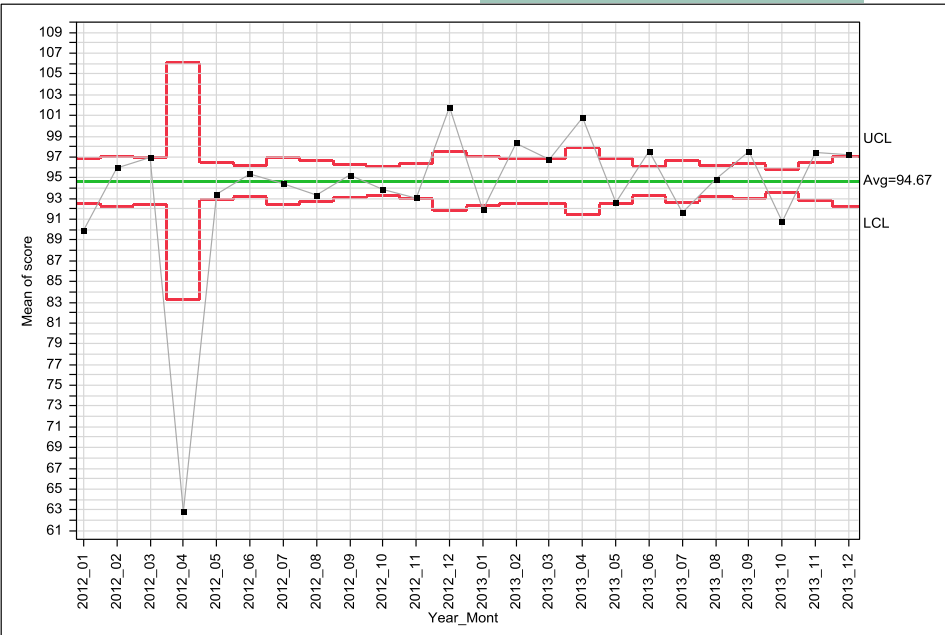
הקשר בין תטא לציון
כמעט מושלם ($r = .99$).

בתרשים מוצגים כל
הנוסחים שבשימוש כיום,
כל נוסח בצבע קו שונה.

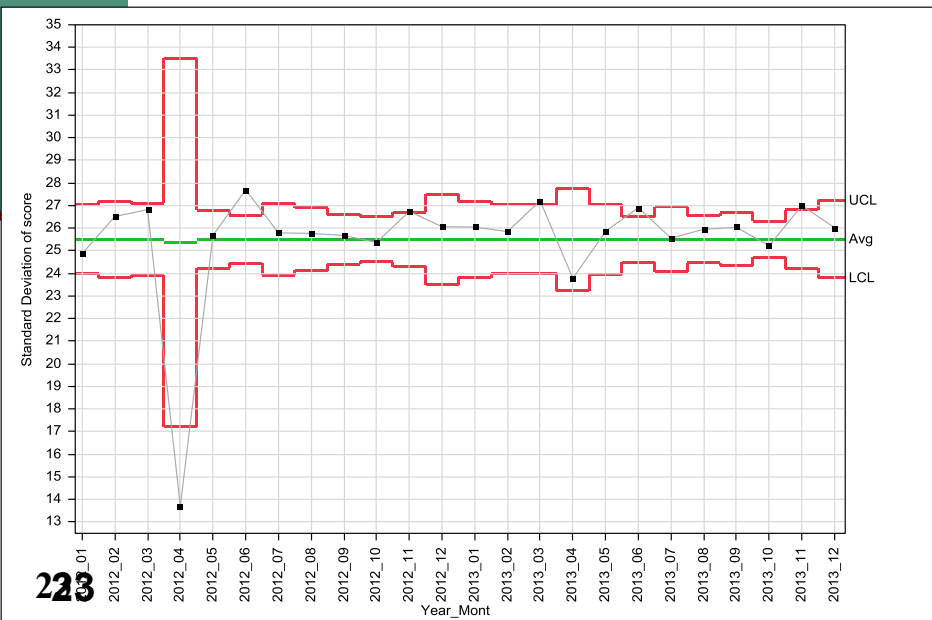
ניתן לראות את הפעלתם
של "כללי עצירה".

תרשימי בקרה על הציונים הגולמיים בשנים 2012-2013

ממוצע



סטיית תקן



דיין

מסקנות

סטטיסטיקה תיאורית

- ישנו שימוש גובר בבחינה במוסדות חדשים לצד הוותיקים (למשל, נבחני האו"פ).
- לצד העלייה העקבית במספר הנבחנים, ניכרת ירידה עקבית בציון הממוצע בבחינה.
- שכיחות הנבחנים באמיר"ם עולה כיום על שכיחות הנבחנים באמי"ר; הציון הממוצע באמיר"ם נמוך מזה של אמי"ר.
- גיל הנבחנים הממוצע הינו 26.6 (ס.ת. 7.6); וגברים מראים בממוצע הישגים גבוהים מנשים (בכ-0.4 ס.ת.).

מסקנות (2)

בקרת איכות

- הציון הממוצע וסטיית התקן בנוסחים השונים די יציבים.
- התכנסות הבחינה עומדת על כ-80 אחוז בממוצע.
- תרשימי בקרת איכות מאפשרים מעקב ויזואלי נוח אחר ממוצעי הציונים.

המחקר הנוכחי מרחיב את היריעה בדרך לתפיסה כוללת של בקרת איכות של מבחנים המועברים תכופות לקבוצות נבחנים קטנות, תוך נתינת מענה גם למבחנים אדפטיביים בנוסף למבחנים לינאריים.

פלט המערכת הממוכנת לבקרה מתמשכת על בחינת המימ"ד

אחד מתוצרי המחקר המתוכננים הינו מערכת ממוכנת אוטומטית לבקרת הציונים טרם דיווחם, שתפעל באופן שבועי.

