

72+2

מועד א'
סמסטר א', תשס"ט
12.2.2009

הפקולטה למדעים מדויקים
החוג למדעי המחשב
אוניברסיטת תל אביב

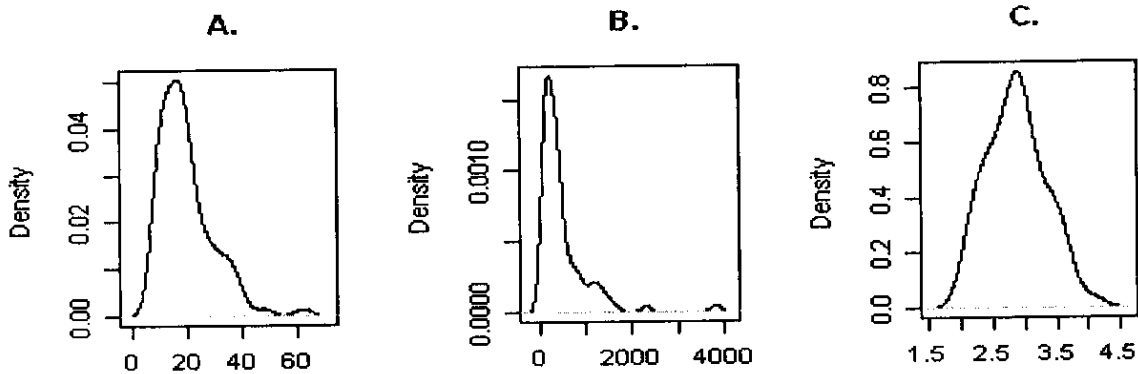
סטטיסטיקה למדעי המחשב
מרצה: ד"ר סהרון רוסט

16
25
333
415
513
*2

משך הבחינה 3 שעות. הינכם רשאים להביא עמכם מחשבון בלבד. דפי נוסחאות וטבלאות יסופקו עם הבחינה. עליכם לענות על כל השאלות בגוף השאלון. הקפידו על תשובות קצרות וברורות. המחברת תשמש לחישובים בלבד. סעיפים המסומנים ב(*) הם סעיפי בונוס והמבחן מסתכם ל100 נקודות בלעדיהם. בהצלחה!

מס' מחברת: 47

שאלה 1 (10): נתונות שלוש הצגות של מחיריהן של 93 דגמי מכוניות בארה"ב (בדולרים, נכון לשנת 1993). בחלון A מוצגת החלקת חלון של המחירים המקוריים, ב B של המחירים אחרי העלאה בריבוע, וב C של המחירים אחרי טרנספורמציה log.



א. אם נתון שערכי תמצית פירסון לנטייה עבור שלושת הגרפים הם 0.25, 1.48, 3.29, שייכו כל אחד מהערכים לאחד מהגרפים והסבירו בקצרה.

ידיד זה תגיד ד זכירה ג קיימת: סכום $T(x) \leq 3$ אורך ימניה, לפונקציה $f(x)$ ימניה. זהו 0.25 - נראה לפי $f(x) \leq 3$ שכן זהו האדק הנמוך ביותר להיות הפונקציה, ולפי C היא הקרובה ביותר לסימטריה. 1.48 - נראה לפי A שכן הנטייה בו ימניה גדולה מאוד ב-C, אך לפי הפקולטה נראה לפי B שכן זהו האדק הנמוך ביותר להיות הפונקציה, ולפי C היא הקרובה ביותר לסימטריה.

ב. מה ניתן ללמוד על התפלגות מחירי מכוניות מכך שהגרף ב C (בסקאלת לוג) נראה כמעט סימטרי, בניגוד לגרף ב A? הדרכה: נסו להימנע משימוש במונח לוג ולדבר במונחים שיהיו מובנים לאדם ברחוב.

נראה לפי A שיש להימנע משימוש במונח לוג ולדבר במונחים שיהיו מובנים לאדם ברחוב. נראה לפי B שיש להימנע משימוש במונח לוג ולדבר במונחים שיהיו מובנים לאדם ברחוב. נראה לפי C שיש להימנע משימוש במונח לוג ולדבר במונחים שיהיו מובנים לאדם ברחוב.

שאלה 2 (15): בדיקה חדשה וזולה ל-HIV מפורסמת כ"מדויקת ב-99%". העלון המצורף מסביר שהכוונה ש-99% מהנשאים מזוהים ככאלה ו-99% מהבריאים מזוהים ככאלה.

א. מה הם שני סוגי הטעויות שהמבחן יכול לעשות? כתבו אותן כהסתברויות מותנות.

נניח ש- X מוכר כבריא ולא אם הוא אדם נגיף. H ו- S הם אירועי המבחן: I : המבחן עלה לזרוע ולא S : המבחן עלה לזרוע, כשהוא אמיתי

על אולם, המבחן ההסתברותי המאמין (נסמן את X כ- X התקבל S לבריא ו- S לחולה) נוסמן תוצאת הבדיקה כ- Y התקבל H לנצטוב בריא ו- S לנצטוב חולה):

הסיכוי שאדם בריא ינצטוב חולה עם הבדיקה: $P(Y=S | X=H)$

(2) למצא מסומן II: המבחן עלה לזרוע ולא S : המבחן עלה לזרוע, כשהוא אמיתי

המבחן. הקולנו ההסתברותי המאמין אדם נגיף ו- S : המבחן עלה לזרוע, כשהוא אמיתי $P(Y=H | X=S)$

ב. אנחנו בוחרים מדגם מקרי של 50 אנשים שהוכרוזו נשאים על ידי הבדיקה ובודקים אותם שוב בבדיקה יקרה שדיוקה 100%. אנו מגלים ש-45 מתוכם למעשה בריאים. כתבו את האמירה הזו כהסתברות מותנית אמפירית.

נסמן את מספר החולים האמיתי ב- X , התקבל ערכים $0 \leq k \leq 50$ (ב- X).

נסמן את מספר החולים שהבדיקה עלה ב- Y התקבל ערכים $0 \leq l \leq 50$ (ב- Y).

האמירה ההסתברותית היא: $P(Y=l | X=k)$

הישר הבדיקה $k=5$ (5 חולים בריאים) ו- $l=50$ (תוצאת הבדיקה) -

ג. הניחו ששתי הסתברויות הטעות בסעיף א שוות ל-1% בדיוק. האם התוצאה בסעיף ב אפשרית? נמקו בקצרה.

התוצאה בסעיף ב אינה יכולה להתקיים, ייתכן שגם התוצאה הזו קרה. ייתכן גודלן. למעשה רק בצורה אחת בריא כחולה, גודלן אינו, ייתכן שגם חולים בריאים ו- S : המבחן עלה לזרוע, כשהוא אמיתי. המבחן עלה לזרוע, כשהוא אמיתי. המבחן עלה לזרוע, כשהוא אמיתי.

50 איש ← 45 בריאים: המבחן עלה לזרוע, כשהוא אמיתי, ונניח לו שמספר

נסמן לחולה ונניח ~~שהמבחן עלה לזרוע, כשהוא אמיתי~~ המבחן עלה לזרוע, כשהוא אמיתי.

S : המבחן עלה לזרוע, כשהוא אמיתי (נניח ונניח) המבחן עלה לזרוע, כשהוא אמיתי.

ייתכן שגם המבחן עלה לזרוע, כשהוא אמיתי. המבחן עלה לזרוע, כשהוא אמיתי.

על אולם, המבחן ההסתברותי המאמין (נסמן את X כ- X התקבל S לבריא ו- S לחולה) נוסמן תוצאת הבדיקה כ- Y התקבל H לנצטוב בריא ו- S לנצטוב חולה):

הסיכוי שאדם בריא ינצטוב חולה עם הבדיקה: $P(Y=S | X=H)$

7. (*5 נקודות) בהינתן שני מבחנים א, ב, צרו דוגמא לפרדוקס של סימפסון שבה התניה על תוצאות מבחן א מייצרת אשליה כאילו מבחן ב אינו עוזר לזהות נשאים, בעוד שהמסקנה ללא ההתניה היא שונה.

המשקל $\frac{1}{6}$:
 חצי היגד 6 אנשים 'לבו כחול' קיבו הנמינא תניכתיס, וזינו
 \Rightarrow אם יבין אם יבן שזכו 50 חולס, $\frac{1}{6}$ אגישור האמיני אל התולס בטוליסיה
 הוא כפי שנוצק בבדיקה היצויקד ח-1000 (5 יגוק 50). אין גוצט-סאז
 ב' לא יאשרה.

שאלה 3 (35): חוק האבל (Hubble) באסטרונומיה קובע שמהירות התרחקותו v של גלקסיות זו מזו פרופורציונלי למרחקן d . פרמטר היחס ידוע בשם קבוע האבל והוא מסומן ב- H . בשאלה הזו אנו מבקשים להעריך אותו. להלן נתונים על שתי גלקסיות (ליתר דיוק, צבירי גלקסיות): מרחקן משביל החלב (הגלקסיה שלנו) ומהירות ההתרחקות שלהן מאיתנו.

צביר	מרחק d (מיליוני שנות אור)	מהירות v (אלפי ק"מ/שניה)
וירגו	22	1.13
ג'מיני	405	23.2

נניח שהמרחקים d_i ידועים בדיוק מראש (דהיינו הם מספרים) ואילו המהירויות v_i הן תצפיות מקריות הנמדדות עם רעש: $v_i = Hd_i + \varepsilon_i$, עבור $i=1,2$, כאשר ε_i הן שגיאות בלתי תלויות עם תוחלת אפס.

אנו שוקלים שני אומדים אפשריים עבור הפרמטר H :

$$\hat{H}_1 = \frac{d_1 v_1 + d_2 v_2}{d_1^2 + d_2^2} \quad \hat{H}_2 = 0.5 \cdot \left(\frac{v_1}{d_1} + \frac{v_2}{d_2} \right)$$

א. בהנחה שהשגיאות ε_i מתפלגות $N(0, \sigma^2)$ (ולכן $(v_i | H) \sim N(Hd_i, \sigma^2)$), כתבו את לוג פונקציית הנראות של H במונחים של d_1, d_2, v_1, v_2 .

$$L(H; d_1, d_2, v_1, v_2) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left\{-\frac{(v_1 - Hd_1)^2}{2\sigma^2}\right\} \cdot \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left\{-\frac{(v_2 - Hd_2)^2}{2\sigma^2}\right\} = \frac{1}{2\pi\sigma^2} \cdot \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^2 (v_i - Hd_i)^2\right\}$$

$$\Rightarrow \ell(H; d_1, d_2, v_1, v_2) = \log\left(\frac{1}{2\pi\sigma^2}\right) - \frac{1}{2\sigma^2} \sum_{i=1}^2 (v_i - Hd_i)^2 = -\log(2\pi\sigma^2) - \frac{1}{2\sigma^2} [v_1^2 - 2Hv_1d_1 + H^2d_1^2 + v_2^2 - 2Hv_2d_2 + H^2d_2^2]$$

$$= -\log(2\pi\sigma^2) - \frac{1}{2\sigma^2} [v_1^2 + v_2^2] + H \left[\frac{v_1d_1 + v_2d_2}{\sigma^2} \right] - H^2 \left[\frac{d_1^2 + d_2^2}{2\sigma^2} \right]$$

ב. גזרו אותו והסיקו שאחד משני האומדים לעיל הוא אומד נראות מקסימלי במצב זה (איזה מהם?).

$$\frac{\partial \ell(H)}{\partial H} = 0 + 0 + \frac{v_1d_1 + v_2d_2}{\sigma^2} - \frac{2H(d_1^2 + d_2^2)}{2\sigma^2} = 0$$

$$\frac{\partial \ell(H)}{\partial H} = 0 \Leftrightarrow v_1d_1 + v_2d_2 = H(d_1^2 + d_2^2) \Leftrightarrow \boxed{H = \frac{v_1d_1 + v_2d_2}{d_1^2 + d_2^2}}$$

$$\frac{\partial^2 \ell(H)}{\partial H^2} = -\frac{d_1^2 + d_2^2}{\sigma^2} < 0$$

לכן \hat{H}_1 הוא האומד הנראות מקסימלי.

ג. האם האומדים מוטים?

$$(1) E(\hat{H}_1) = E\left(\frac{d_1v_1 + d_2v_2}{d_1^2 + d_2^2}\right) = \frac{1}{d_1^2 + d_2^2} \cdot (d_1E(v_1) + d_2E(v_2)) = \frac{d_1 \cdot Hd_1 + d_2 \cdot Hd_2}{d_1^2 + d_2^2} = \frac{d_1^2 + d_2^2}{d_1^2 + d_2^2} H = H$$

לכן \hat{H}_1 אינו מוטה.

$$(2) E(\hat{H}_2) = E\left(\frac{1}{2} \left[\frac{v_1}{d_1} + \frac{v_2}{d_2} \right]\right) = \frac{1}{2} \left(\frac{1}{d_1} E(v_1) + \frac{1}{d_2} E(v_2) \right) = \frac{1}{2} \left(\frac{H}{1} + \frac{H}{1} \right) = \frac{2H}{2} = H$$

לכן \hat{H}_2 אינו מוטה.

הן v_1, v_2

7. מהן השונות שלהם? איזו גדולה יותר?

$$(1) \text{Var}(\hat{H}_1) = \text{Var}\left(\frac{d_1 v_1 + d_2 v_2}{d_1^2 + d_2^2}\right) = \frac{1}{(d_1^2 + d_2^2)^2} (d_1^2 \text{Var}(v_1) + d_2^2 \text{Var}(v_2)) = \frac{d_1^2 + d_2^2}{(d_1^2 + d_2^2)^2} \sigma^2 = \frac{\sigma^2}{d_1^2 + d_2^2}$$

$$(2) \text{var}(\hat{H}_2) = \text{var}\left(\frac{1}{2} \left[\frac{v_1}{d_1} + \frac{v_2}{d_2} \right]\right) = \frac{1}{4} \left[\frac{1}{d_1^2} \text{var}(v_1) + \frac{1}{d_2^2} \text{var}(v_2) \right] = \frac{d_1^2 + d_2^2}{4 d_1^2 d_2^2} \sigma^2$$

הן v_1, v_2

נראה כי ההפרש: המכנה של $\text{var}(\hat{H}_1)$ הוא המכנה של $\text{var}(\hat{H}_2)$ כפי שהתקבלה. $d_1^4 - 2 d_1^2 d_2^2 + d_2^4 \stackrel{?}{>} 0$
 עבור $d_1 = d_2$ נקבל שני המכנים זהים וזה נכון. $d_1 \neq d_2$ נקבל גישות שונות. $d_1^4 - 2 d_1^2 d_2^2 + d_2^4 = (d_1^2 - d_2^2)^2 \geq 0$
 אז $\text{var}(\hat{H}_1) \leq \text{var}(\hat{H}_2)$ כלומר המכנה של \hat{H}_1 הוא יותר מזה של \hat{H}_2 .
 ה. על סמך סעיפים א-ד, איזה משני האומדים עדיף? ציינו שתי סיבות לבחירתכם.

לשני האומדים התעניין נעזרתי או \hat{H}_1 בעיקר, משתי סיבות:

- הישגה שלו קטנה יותר מזה של \hat{H}_2 , כלומר הוא יותר מדויק (פחות כיוונו).

- (הוא אוניברסלי יותר) מקסימום H .

8. חשבו את הערך של שני האומדים על הנתונים שלנו.

הערה: האומדים הם ביחידות מוזרות שכן המרחק נמדד בשנות אור והמהירות בק"מ לשנייה. כדי להגיע למדידה ביחידות מקובלות יש להכפיל בקבועי המרה עצומים, ואנו חוסכים זאת מכאן.

$$\hat{H}_1 = \frac{22 \cdot 1.13 + 405 \cdot 23.2}{22^2 + 405^2} \approx 0.0572$$

$$\hat{H}_2 = \frac{1}{2} \left(\frac{1.13}{22} + \frac{23.2}{405} \right) \approx 0.0543$$

9. בהנחה שהקצב אכן קבוע, $1/H$ הוא הזמן שעבר מאז המפץ הגדול. האם $1/\hat{H}_2$ או $1/\hat{H}_1$ הוא אומדן

חסר הטיה לזמן זה? האם אחד מהם הוא אומדן נראות מקסימלי?

\hat{H}_1, \hat{H}_2 שני $1/H$ חסרי הטיה \Rightarrow $\frac{1}{\hat{H}_1}$ ו $\frac{1}{\hat{H}_2}$ הם אומדנים חסרי הטיה \Rightarrow $\frac{1}{\hat{H}_1}$ ו $\frac{1}{\hat{H}_2}$ הם אומדנים חסרי הטיה ל H .

- $\frac{1}{\hat{H}_1}$ הוא אומדן חסר הטיה ל H כי $\frac{1}{\hat{H}_1} = \frac{1}{\frac{d_1 v_1 + d_2 v_2}{d_1^2 + d_2^2}} = \frac{d_1^2 + d_2^2}{d_1 v_1 + d_2 v_2}$ הוא אומדן חסר הטיה ל H כי $H = \frac{1}{\frac{d_1 v_1 + d_2 v_2}{d_1^2 + d_2^2}}$

ולכן סיוק צבאי אינו סביר (הסבירה), והיא אומדן חסר הטיה ל H כי $H = \frac{1}{\frac{1}{2} \left(\frac{v_1}{d_1} + \frac{v_2}{d_2} \right)}$

- $\frac{1}{\hat{H}_2}$ הוא אומדן חסר הטיה ל H כי $H = \frac{1}{\frac{1}{2} \left(\frac{v_1}{d_1} + \frac{v_2}{d_2} \right)}$

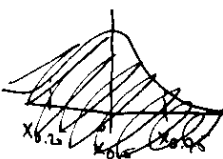
שאלה 4 (20): ציינו נכון/לא נכון ונמקו בקצרה (תשובות ללא נימוק תקף לא יקבלו נקודות):

א. מבחן א ומבחן ב בודקים את אותה השערת אפס. אזי אם מבחן א דוחה בכל מקרה שמבחן ב דוחה, אזי מבחן א הוא בעל עצמה גבוהה יותר. לא נכון

ב. ניח $X \sim N(\mu, \sigma^2)$, ואנו עושים הסקה בייזאנית על התוחלת הנורמלית μ , תוך שימוש בהתפלגות מקדימה על μ של $N(0, \sigma^2)$. אזי לשונות ההתפלגות המקדימה σ^2 אין השפעה על תוחלת האומדן הבייזאני הא-פוסטריורי. לא נכון

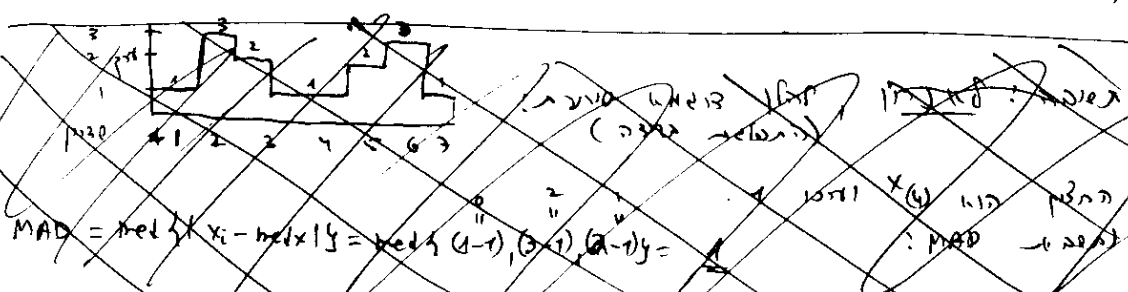
ג. בונים רווח סמך ברמה 95% לתוחלת נורמלית, פעם בהנחת שונות ידועה σ^2 ושימוש ב Z, ופעם בהנחת שונות לא ידועה --- אומדים את σ^2 ובונים רווח סמך תוך שימוש בערך t המתאים. אזי רווח הסמך הראשון יהיה קטן יותר. נכון

ד. אם התפלגות היא סימטרית אזי MAD שווה למחצית הטווח הבין רבעוני IQR (בדיוק על ההתפלגות התיאורטית ובערך על מדגם). נכון



~~שני רגעי הסמך א סוגים (משנים) והוויכוח ראה בעמוד הבא... Q_1 ו Q_3~~

ה. פירסון המבוסס על הקירוב הנורמלי הוא גם מבחן ניימן-פירסון לבעיה המקורית. נכון



שאלה 5 (20): בארץ גיבורי הילדות התמודדו בבחירות המועמדים קיפי וביילבי. בסקר ראשון שנערך בקרב מדגם מקרי של 600 מצביעים, תמכו בקיפי 40% וביילבי 60% מתוכם. אז הכריזה קיפי מלחמה על ארץ חלאס השכנה, שיושביה הטרידו את גיבורי הילדות. לאחר המלחמה נערך מדגם מקרי נוסף של 600 מצביעים, וכעת תמכו בקיפי 45% מתוכם, והשאר בביילבי.

א. ברצוננו לבחון את השערת האפס שהמלחמה לא השפיעה על התמיכה בקיפי (מול האלטרנטיבה שהעלתה התמיכה). נסחו את ההשערות ובצעו את המבחן הבינומי המתאים ברמה 0.05. מה המסקנה? יהי p שיעור התמיכה בקיפי לאחר המלחמה:

$$H_0: p = 0.4 \text{ (פנימה)} - \text{אין שינוי במעמד היחסים בין המועמדים}$$

$$H_1: p > 0.4 \text{ (צדקה)} - \text{התמיכה בקיפי גדולה יותר (שינוי 0.4 דברגסס 6 בקר 1)}$$

המבחן הבינומי (המשמ):

$$C_{\alpha} = \left[p_0 + z_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{p_0(1-p_0)}{n}}, \infty \right) = \left[0.4 + 1.645 \cdot \sqrt{\frac{0.4 \cdot 0.6}{600}}, \infty \right) = [0.4329, \infty)$$

כיצד $\beta = 0.45$ הוא בעצם הדחייה, ולכן צדקה H_1 .

לפני המסקנה היא:

ברמה מוחלטת של 95%, המלחמה הגבירה את שיעור התמיכה בקיפי בטווחים מסוימים.

ב. מכון הסקרים "מין הצמח" העדיף לבדוק את השערת האפס לעיל באמצעות מבחן חי-בריבוע לטבלת 2×2 המתאימה. בצעו את המבחן והראו שמסקנתו ברמה 0.05 שונה מזו של סעיף א.

נמצא המבחן χ^2 בין גמיכה בקיפי/בילבי לבין לפני/אחרי המלחמה:

obs	גמיכה בקיפי	גמיכה בבילבי	סה"כ
לפני	240	360	600
אחרי	270	330	600
סה"כ	510	690	1200

נמצא המבחן χ^2 לפי ההתפלגות המשותפת:

כל מה שהיה צריך
הצדקה ולכן ניתן לבדוק
למבחן χ^2 , במקרה זה
סה"כ 1 ד"ח.

סה"כ	בילבי	קיפי	סה"כ
600	345	255	לפני
600	345	255	אחרי
1200	690	510	סה"כ

$$\sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}} = 0.832 + 0.652 + 0.182 + 0.652 = 2.318 \approx 2.32$$

כבר לפני הדחייה במקרה זה יהיה: $\chi^2 > 3.84$ (עם $\alpha = 0.05$)

(המקרה, כלומר שיעור התמיכה
לפני המלחמה)

עבור המבחן χ^2 נראה שיש H_0 (כלומר אין שינוי)

הגמיכה בקיפי לא גדלה (95%).

ג. הסבירו בקצרה את שני ההבדלים העיקריים בין המבחנים בסעיפים א וב. אילו מהמבחנים מתאים יותר לסיטואציה המתוארת?
 רמז: אחד מההבדלים נוגע להנחות על המודל הסטטיסטי (אילו גדלים ידועים מראש ואילו מקריים) ואחד להשערות הנבדקות.

לג' מ'ג' א'ק

(1) - הרפיה (הנבדק) המבחן הביומי היוו אב בייני - אפס נוצמה מקסימלית

✓ גוף ממון חי במוצד ~~המבחן~~ (בדק נטל) העב י' גוף או אין גוף

לזו המבחן גוף גוף, לזו המבחן גוף גוף או $p < 0.05$.

(2) - הרפיה א המבחן הביומי המבחן בייני הוא - בייני ... וממון

"- לסיטואציה, במצב ממון חי במוצד המבחן בייני הוא ממון בייני

X $(\sum_{i=1}^n x_i^2 \in \sum_{i=1}^n x_i^2)$, וזו המבחן גוף גוף ~~המבחן~~ (מקרה לה -

מ'כ'י'ק נס'ק יי המבחן המבחן יגר הוא המבחן הפוסורציה הביומי

5/8

ד. (5+*) נקודות) נניח כעת ששני המדגמים בסעיף א הם על אותם 600 מצביעים. דונו בקצרה איך ניתן כעת לבחון את השערת האפס במדגם מזווג כזה.
 רמז והצעה: כדי לענות נכון על שאלה זו תצטרכו לפתח הרחבות מעבר לדברים שלמדנו בכיתה. לכן מומלץ לנסות אותה רק אם זמנכם בידכם.