

מבוא לניתוח נתונים - 096202
מבחן מועד א' – סמסטר חורף תש"ף
טור א'

תאריך הבחינה: 20.02.2020

מרצה: אורי פלונסקי

מתרגל: רפאל שללה

הוראות:

1. לרשותכם **שעתיים וחצי** לפתור את הבחינה.
2. מותר להשתמש במחשבון **פשוט** בלבד. אין להשתמש בכל חומר עזר אחר.
3. הבחינה כוללת 16 עמודים (כולל עמוד זה) ובהם 12 שאלות. בדקו בתחילת הבחינה שיש ברשותכם את כל העמודים.
4. ניקוד כל שאלה (ולרוב כל סעיף) מצוין לידה. סך כל הניקוד האפשרי בבחינה הוא **102** נקודות.
5. **את התשובות יש לכתוב רק על גבי טופס הבחינה, ובמקומות המיועדים לכך בלבד.** מחברות הטייטה יושמדו לאחר סיום הבחינה וממילא לא יבדקו.
6. חובה לכתוב מספר תעודת זהות על כל אחד מדפי טופס הבחינה.
7. אסור להפריד את דפי טופס הבחינה.
8. בשום שלב, במהלך הבחינה או לאחר סיומה, אסור להוציא מחדר הבחינה שום דבר שלא הבאתם איתכם לחדר הבחינה. בפרט, חובה להחזיר בסיום הבחינה את הטופס וכל מחברת בה השתמשתם.

בהצלחה!

שאלה 1 (15 נק')

ענו נכון או לא נכון. נמקו תשובתכם במשפט אחד.

א. כאשר השערת האפס נכונה, לעולם לא נדחה אותה אם נשתמש במבחנים סטטיסטיים

בצורה מתאימה. נכון / לא נכון.

נימוק (במשפט אחד):

ב. נתון משתנה המייצג הכנסה ומקבל את הערכים 1, 2, או 3 בלבד, כאשר: $1 = \text{הכנסה}$

מתחת ל-8000 ₪ בחודש; $2 = \text{הכנסה}$ בין 8001 ל-15000 ₪ בחודש; $3 = \text{הכנסה}$ של יותר

מ-15000 ₪ בחודש. משתנה זה הינו משתנה קטגורי. נכון / לא נכון.

נימוק (במשפט אחד):

ג. בתהליך K-fold cross validation על אלגוריתם kNN, מחלקים את הנתונים ל-K חלקים

ובכל חלק בודקים את הדיוק של שימוש במספר שכנים, k, אחר. לבסוף בוחרים את מספר

השכנים k שנותן את הדיוק הגבוה ביותר. נכון / לא נכון.

נימוק (במשפט אחד):

ד. נתון שבאמצעות שיטת בוטסטראפ חושב רווח סמך ברמת ביטחון של 99% להפרש בין

מספר הבאגים היומי הממוצע של מתכנת ממחלקה א' למספר הבאגים היומי הממוצע של

מתכנת ממחלקה ב' והתקבל [1.3, 2.7]. לכן, ברמת מובהקות של 5%, ניתן לדחות את

השערת האפס שמספר הבאגים היומי הממוצע של מתכנת ממחלקה א' זהה למספר

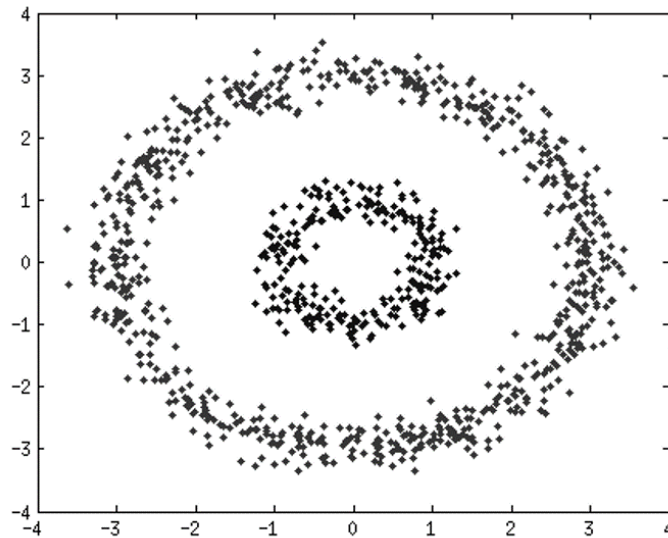
הבאגים היומי הממוצע של מתכנת ממחלקה ב'. נכון / לא נכון.

נימוק (במשפט אחד):

[השאלה ממשיכה בעמוד הבא]

ה. בידינו תצפיות עם שני משתנים נומריים. תרשים פיזור של הקשר בין שני המשתנים מוצג להלן. על נתונים אלו הורץ אלגוריתם K-Means עם $K=2$. האלגוריתם סיווג את כל הנקודות שיוצרות מעגל קטן פנימי לאשכול (cluster) אחד ואת כל הנקודות שיוצרות מעגל גדול חיצוני לאשכול (cluster) שני. נכון / לא נכון.

נימוק (במשפט אחד):



שאלה 2 (6 נק')

לפני מבחן במתמטיקה לכיתות ח' בבית ספר "ברושים" נאספו נתונים מכלל תלמידי כיתות ח' בבית הספר בנוגע לסך כל זמן הלמידה שלהם למבחן (בדקות). לאחר מכן, נתונים אלו הושוו לציוני המבחן. התקבל כי הקורלציה (מתאם) הלינארית בין זמן הלמידה (בדקות) לבין ציון המבחן היא 0.74. לכל אחד מההיגדים הבאים כתבו נכון או לא נכון ונמקו במשפט אחד.

א. מאחר שהנתונים כוללים את כל האוכלוסייה (ולא רק מדגם מתוכה), ניתן להסיק שעבור אוכלוסייה זו, זמן למידה ארוך יותר למבחן במתמטיקה גורם לציונים גבוהים יותר במבחן.

ב. באמצעות הנתונים שמתואר שנאספו, ניתן להשתמש בשיטת בוטסטראפ על מנת לחשב רווח סמך ברמת ביטחון 95% לקורלציה הלינארית בין זמן למידה וציון המבחן ואם יתקבל שרווח הסמך לא כולל את הערך 0, נוכל לדחות את ההשערה שאין קורלציה לינארית בין זמן למידה וציון המבחן.

שאלה 3 (3 נק')

עבור איזה מבין השאלות הבאות הכי הגיוני להשתמש באלגוריתם kNN? (הקיפו את התשובה הנכונה)

- א. האם לקוחות רשומים מוציאים בממוצע יותר כסף באתר מאשר לקוחות שאינם רשומים?
- ב. בהתבסס על היסטורית החיפוש והגלישה של לקוח באתר, האם הוא צפוי לרכוש מוצר כלשהו בשבוע הקרוב?
- ג. מהו זמן הגלישה הממוצע של לקוחות רשומים באתר?
- ד. איזה מהלקוחות הרשומים באתר מתנהגים בצורה דומה זה לזה מבחינת היסטורית החיפוש והקנייה שלהם מהאתר?

שאלה 4 (7 נק')

בכד אטום מבחון יש 3 כדורים אדומים, 2 כדורים ירוקים וכדור אחד שחור. לבד מהצבע, כל הכדורים זהים לחלוטין.

- א. במשחק א', שחקן שולף מהכד שני כדורים בזה אחר זה **וללא החזרה**. אם אף אחד משניהם אינו אדום, הוא מנצח במשחק. מה ההסתברות של השחקן לנצח במשחק א'? (2 נק')

- ב. במשחק ב', החוקים זהים לאלו של משחק א', אבל כעת, לאחר שליפת כל אחד מהכדורים, הם **מוחזרים לכד** לפני שליפת הכדור הבא. שוב, על מנת לנצח, יש לשלוף 2 כדורים שאינם אדומים. מה ההסתברות לנצח במשחק ב'? (2 נק')

- ג. שחקן משחק ראשית את משחק א'. **אם** השחקן מנצח במשחק א', הכדורים ששולף מוחזרים לכד ואז השחקן זכאי לשחק את משחק ב'. **אם** השחקן משחק ומנצח גם במשחק ב', הוא זוכה בפרס. מה ההסתברות ששחקן לא יזכה בפרס? (3 נק')

שאלה 5 (16 נק')

נשיא הטכניון טוען שבוגרי טכניון שכסטודנטים למדו את הקורס מבוא לניתוח נתונים הם בעלי סיכוי גבוה יותר לעבוד כמדעני נתונים שנה לאחר סיום התואר מאשר בוגרי הטכניון שכסטודנטים לא למדו את הקורס מבוא לניתוח נתונים. לשם בדיקת הטענה, הוא מבקש מאלון לאסוף נתונים ולערוך מבחן סטטיסטי.

א. מהי האוכלוסייה שלגביה או האוכלוסיות שלגביהן רוצים לבדוק את ההשערות? (2 נק')

ב. הגדירו את השערת האפס ואת ההשערה האלטרנטיבית (3 נק')

ג. מהו סטטיסטי המבחן שבאמצעותו תבדקו את ההשערות? (2 נק')

ד. אלון מצליח להשיג גישה לרשימת מספרי הטלפון של כלל בוגרי הטכניון (שנה לאחר סיום התואר) שלקחו את הקורס מבוא לניתוח נתונים. הוא בוחר באקראי 300 מספרי טלפון ומתקשר אליהם במטרה לשאול האם הם עובדים כיום כמדעני נתונים. מאחר שחלקם לא זמינים וחלקם לא מעוניינים לענות, הוא בסופו של דבר מקבל נתונים של 237 מהם. הציעו הטיה אחת אפשרית בתהליך איסוף הנתונים המתואר אשר עלולה להשפיע על התוצאות של הניתוח שאלון צפוי לעשות. הסבירו מה הבעיה שעלולה להיווצר. (3 נק')

[השאלה ממשיכה בעמוד הבא]

ה. בנוסף, אלון מקבל גישה לנתונים מתוך מדגם מקרי פשוט של מדעני נתונים בארץ. נתונים אלו כוללים בין היתר את שנת סיום התואר ומוסד אקדמי בו למדו. האם נתונים האלה, יחד עם הנתונים המתוארים בסעיף ד', מספיקים על מנת לבדוק את הטענות של נשיא הטכניון? אם כן, הסבירו בקצרה כיצד תבדקו את הטענה. אם לא, הסבירו מדוע הנתונים לא מספיקים ואילו נתונים נוספים דרושים כדי לבדוק אותה (הניחו שבנתונים שנאספו בסעיף ד' אין הטיות). (4 נק')

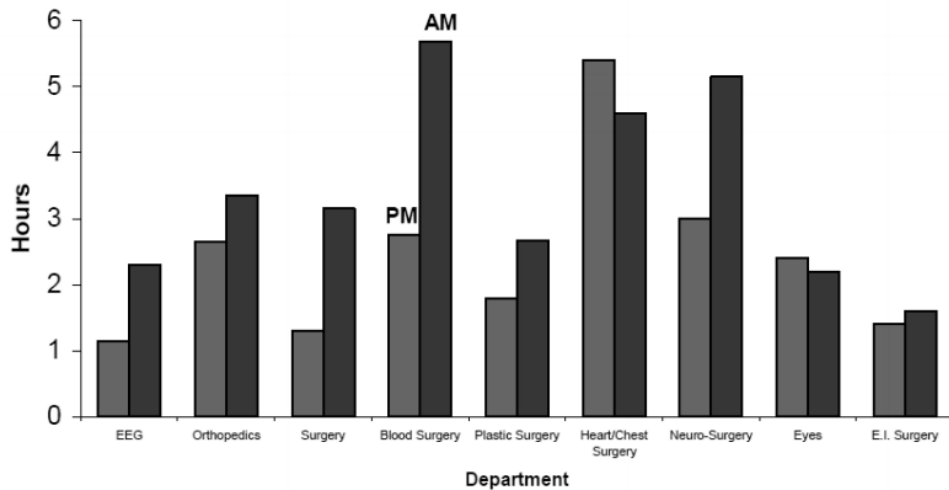
ו. בסופו של דבר אלון מצליח להשיג מספיק נתונים (אלו מסעיפים קודמים ו/או נתונים אחרים) על מנת לבחון כיאות את הטענה של נשיא הטכניון. הוא עורך את המבחן הסטטיסטי (בצורה נכונה) ומסיק שיש לדחות את השערת האפס. מכאן ניתן להסיק שלמידת הקורס מבוא לניתוח נתונים גורמת לעלייה בסיכוי לעבוד כמדען נתונים. נכון או לא נכון? נמקו במשפט אחד. (2 נק')

[illegible]

שאלה 7 (7 נק')

על מנת לשפר את היעילות, מנהלת בית חולים "רפואה שלמה" החליטה על תכנית לפיה חדרי הניתוח יהיו פעילים גם בשעות אחר הצהריים. על מנת לייצר תמריצים לרופאים שיישארו לנתח גם בשעות אחה"צ, הוחלט כי בשעות אלו הם יקבלו שכר לפי מספר ניתוחים ולא שכר שעתי כמקובל. כלומר, בעוד שבשעות הבוקר, התשלום לרופאים נעשה לפי שעת עבודה, בשעות אחה"צ, התשלום נעשה לפי מספר ניתוחים שהם מבצעים. לאחר תקופת הרצה של תכנית זו שכללה מאות רבות של ניתוחים, נמדד משך ניתוח ממוצע, בכל אחת מ-9 מחלקות, עבור ניתוחים שבוצעו בבוקר וניתוחים שבוצעו אחה"צ בנפרד. להלן גרף המציג את התוצאות (עמודות כהות מייצגות זמני ניתוח ממוצעים בשעות הבוקר, עמודות בהירות מייצגות זמני ניתוח ממוצעים בשעות אחה"צ):

Operations Time - Morning (by Hour) vs. Afternoon (by Case):



א. ישנה אסוציאציה (association) בין המשתנה "משך ניתוח ממוצע" לבין המשתנה "זמן ביום" (המקבל את הערכים "בוקר" או "אחר הצהריים"). נכון או לא נכון? נמקו במשפט אחד (2 נק')

[השאלה ממשיכה בעמוד הבא]

ב. לאחר שראה את הגרף, סגן מנהל בית החולים הסיק שהרופאים ממהרים לסיים יותר ניתוחים כאשר הניתוחים מתבצעים בשעות אחה"צ וזאת על מנת להספיק לבצע יותר ניתוחים וכך לקבל שכר גבוה יותר. מנהל המחלקה הכירורגית טען שזה לא יתכן שהרופאים מסכנים את חיי המטופלים במודע ע"י כך שמזדרזים לקצר את הניתוחים. עם זאת, לטענתו רופאים מאריכים באופן מלאכותי את זמני הניתוח בבקרים על מנת לבצע פחות ניתוחים בבקרים וכך לעבוד פחות קשה, בעוד שהתשלום עבור זמן העבודה שלהם בבוקר קבוע. לפי הנתונים שבידיכם, מי מהם צודק (אם בכלל)? (2 נק')

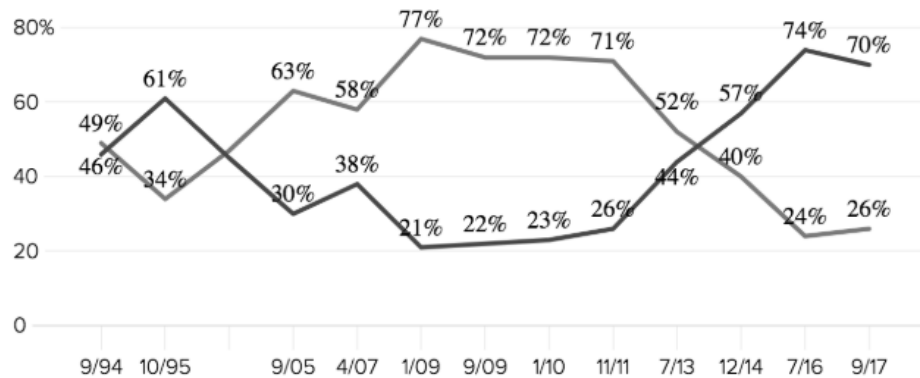
ג. מנהלת בית החולים תשאלה מספר מנתחים וטענה כי ישנו משתנה מתערב (confounder) אשר מסביר את התוצאות בגרף טוב יותר מההסברים של סגן מנהל בית החולים ושל מנהל המחלקה הכירורגית. הציעו משתנה מתערב אשר סביר שגורם לתוצאות הנצפות. (3 נק')

שאלה 8 (7 נק')

נתונה הויזואליזציה הבאה:

In general, do you think race relations in the U.S. are ____ ?

■ Total Good ■ Total Bad



NBC NEWS

Data: NBC News/Wall Street Journal poll. September 14-18, 2017

א. פרטו בנוגע לשתי בעיות לפחות בויזואליזציה המטעות או מבלבלות את הקורא. (הצבעים דומים רק עקב אילוצי הדפסת המבחן: צבעים דומים זו אינה טעות של הויזואליזציה) (4 נק')

ב. ציירו מחדש את הויזואליזציה כך שלא תהיה מטעה (אין צורך לדייק, רק שנבין מה הרעיון)

(3 נק') [מקום לאיור]

שאלה 9 (8 נק')

במסגרת הקורס "מבוא לניתוח נתונים", דני התבקש לייצר אלגוריתם סיווג מסוג kNN שישתמש בשני פיצ'רים (predictor_1, predictor_2) על מנת לסווג תצפיות לאחת משתי מחלקות (Class). להלן סיכום הנתונים שדני קיבל לשם ביצוע המשימה:

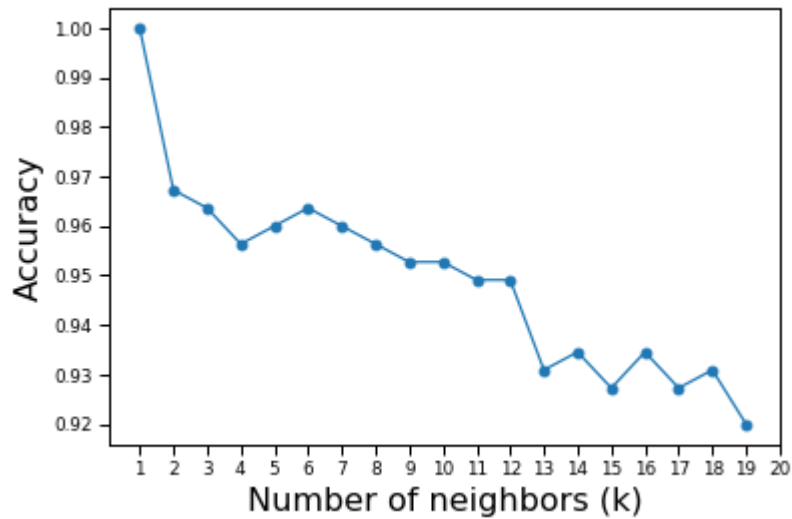
```
1 knn_df.describe()
```

	predictor_1	predictor_2	Class
count	1372.000000	1372.000000	1372.000000
mean	0.539114	0.587301	0.444606
std	0.205003	0.219611	0.497103
min	0.000000	0.000000	0.000000
25%	0.379977	0.451451	0.000000
50%	0.543617	0.602168	0.000000
75%	0.711304	0.770363	1.000000
max	1.000000	1.000000	1.000000

דני החליט לאמן את האלגוריתם עם ערכים שונים של k כדי לבדוק איזה מהם הוא הטוב ביותר. להלן הקוד שדני כתב והפלט שקיבל:

```
1 knn_df = knn_df.sample(frac=1)
2
3 X = knn_df.loc[:, knn_df.columns != 'Class'].values
4 Y = knn_df.loc[:, 'Class'].values
5
6 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.20)
7
8 test_scores = []
9 k_list = range(1, 20)
10 for nn in k_list:
11     knn_classifier = KNeighborsClassifier(n_neighbors=nn)
12     knn_classifier.fit(X, Y)
13     test_scores.append(knn_classifier.score(X_test, Y_test))
14
15 # output results
16 plt.plot(k_list, test_scores, '-o')
17 plt.xlabel('Number of neighbors (k)')
18 plt.xticks(np.arange(1,21))
19 plt.ylabel('Accuracy');
```

[השאלה ממשיכה בעמוד הבא]



א. לפי הגרף שהתקבל (בלבד), באיזה ערך של k כדאי לדני לבחור? נמקו במשפט אחד. 2)
(נק')

ב. בקוד שדני כתב ישנה טעות מהותית שגורמת לו לקבל תוצאות שגויות. מהי מספר השורה שבה מופיעה הטעות? מהי הטעות ומדוע זו טעות? כתבו שורת קוד מתוקנת. (שימו לב: כשדני כתב את הקוד, הוא עדיין לא למד את תהליך ה-cross validation ולכן העובדה שלא השתמש בכך אינה טעות) (6 נק')

שאלה 10 (13 נק')

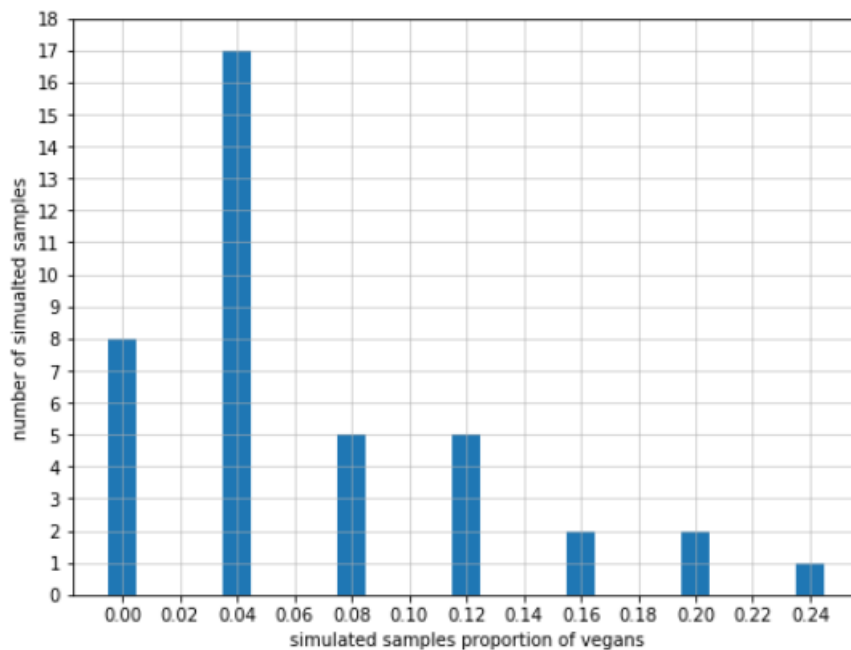
תמר קראה מאמר הטוען כי 6% מבני הנוער בישראל הם טבעונים. היא חושדת כי אצלה בבית הספר אחוז הטבעונים גדול יותר. כדי לבדוק זאת, היא לקחה מדגם מקרי פשוט של 25 תלמידים בבית הספר שלה והתקבל כי 5 מהם טבעונים. לאחר מכן, תמר הריצה סימולציה שבה היא לקחה 40 מדגמים מאוכלוסייה שבה 6% טבעונים וחישבה את אחוז הטבעונים שהתקבלו בכל מדגם בסימולציה.

א. מהי השערת האפס אותה בודקת תמר? (2 נק')

ב. בהנחה שתמר ביצעה את התהליך בצורה נכונה, כמה תלמידים היו בכל מדגם אותו תמר

סימלצה? (2 נק')

ג. להלן ההיסטוגרמה שקיבלה תמר עבור התפלגות אחוז הטבעונים תחת השערת האפס:



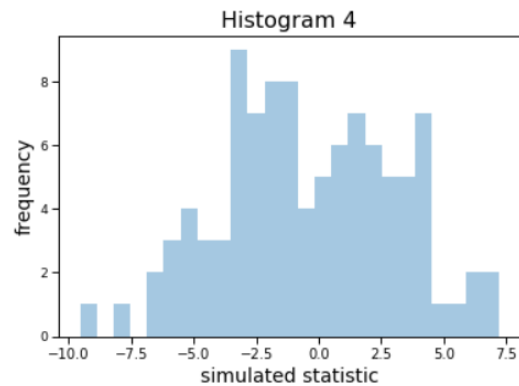
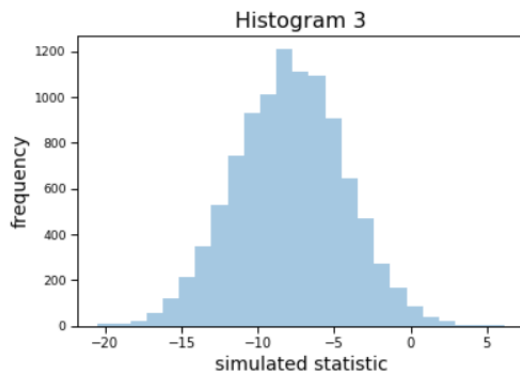
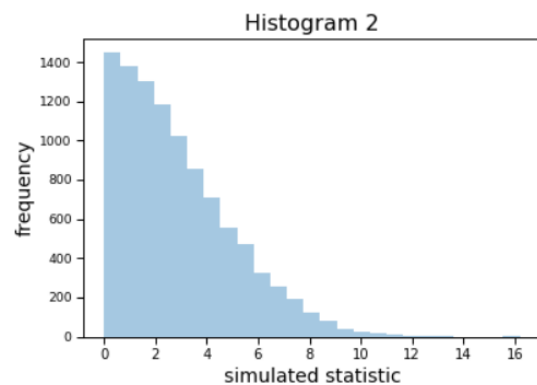
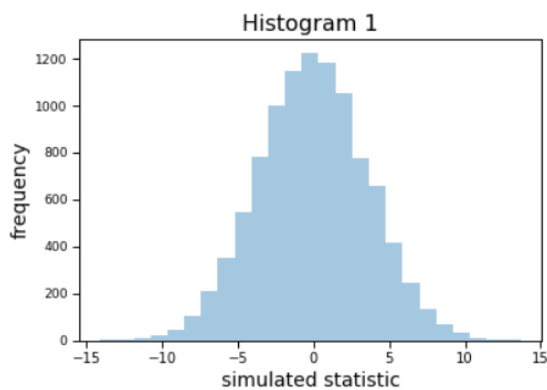
לפי תוצאות אלו, מהו ה-P-value של המבחן? נמקו בקצרה (5 נק')

ד. ברמת מובהקות 5%, האם תמר יכולה לדחות את השערת האפס שלה? (2 נק')

ה. הציעו לתמר דרך פשוטה (שלא כוללת איסוף נתונים נוספים) שתגדיל את הדיוק של הסימולציה שלה. (2 נק')

שאלה 11 (4 נק')

חוקרים רצו לבחון את האפקטיביות של חיסון השפעת השנתי. הם לקחו מדגם מקרי פשוט גדול מאוכלוסיית המתחסנים לשפעת בשנת 2017 וגילו שמתוך האנשים שנדגמו 68% לא חלו בשפעת ב-2017. בנוסף, הם לקחו מדגם מקרי פשוט גדול מאוכלוסיית המתחסנים לשפעת בשנת 2018 וגילו שמתוך האנשים שנדגמו 60% לא חלו בשפעת ב-2018. מכאן, הם חשדו כי אפקטיביות החיסון לשפעת היתה נמוכה יותר בשנת 2018 מאשר בשנת 2017. לשם בדיקת הטענה, הם הריצו 10,000 סימולציות של סטטיסטי המבחן תחת השערת האפס והציגו את ההתפלגות האמפירית של הסטטיסטי בהיסטוגרמה. אחת מארבע ההיסטוגרמות להלן היא ההיסטוגרמה אותה הם הציגו. איזו? (הקיפו את התשובה הנכונה)



- א. היסטוגרמה 1 (Histogram 1)
 ב. היסטוגרמה 2 (Histogram 2)
 ג. היסטוגרמה 3 (Histogram 3)
 ד. היסטוגרמה 4 (Histogram 4)

שאלה 12 (7 נק')

בידיכם נתונים אשר מסוכמים בפלט הבא:

```
1 df.describe(include='all')
```

	x_1	x_2	group
count	100.000000	100.000000	100
unique	NaN	NaN	2
top	NaN	NaN	B
freq	NaN	NaN	92
mean	0.070712	-0.026916	NaN
std	0.760040	0.856362	NaN
min	-2.367725	-2.297804	NaN
25%	-0.405385	-0.689449	NaN
50%	0.093605	0.022972	NaN
75%	0.621761	0.617290	NaN
max	2.100413	1.996156	NaN

מעוניינים להשתמש בנתונים אלו כדי לאמן אלגוריתם kNN שישוו את התצפיות לקבוצות A או B (לפי משתנה group) באמצעות המשתנים x_1 ו-x_2.

א. הוחלט להשתמש בכל הנתונים כסט אימון, ללא ולידציה. כאשר יבחן דיוק האלגוריתם על נתונים חדשים הבאים מאותה אוכלוסייה (סט מבחן), מה צפוי להיות (בערך) הדיוק שיתקבל עבור $k = 17$? נמקו. (4 נק')

ב. הוחלט לבצע cross-validation על מנת לבחור את k. באיזה מדד ביצועים כדאי להשתמש? נמקו בקצרה. (3 נק')
