

# תרגיל 5 : Design Theory

תאריך הגשה : 55 : 23, 03.01.21.

## הוראות הגשה:

בתרגיל זה אתם נדרשים להגיש קובץ zip בודד שיכלול את הקבצים הבאים :

- ex5.pdf עם התשובות מפורטות לשאלות. יש לפרט חישובים לא רק תשובה סופית!
- create.sql המתאים לשאלה 2 סעיף ד.1.
- contradictions.sql המתאים לשאלה 2 סעיף ד.3.
- drop.sql המתאים לשאלה 2 סעיף ד.4.
- README שמכיל שורה בודדת ובו ה-login של הסטודנט שמגיש את התרגיל. אם התרגיל מוגש בזוגות, על שורה זאת להכיל את שני ה-login מופרדים בפסיק.

## שימו לב:

- נא לקרוא על הדרישות המנהליות של הקורס בלינק באתר הקורס כדי למלא אחר ההוראות להגשה של קבצים סרוקים!
- תרגיל מוקלד יזכה ב- 2 נקודות בונוס!

## שאלה 1 (35 נקודות)

נחזור וניזכר במידול מידע על האולימפיאדה מהתרגיל בית הראשון. הפעם, במקום למדל בעזרת דיאגרמת ישויות קשרים, נשתמש בגישת תיאוריית התכנון על מנת להבין איך יש להפריד טבלה אחת גדולה לתתי טבלאות.

הערה: בטבלה המקורית של מידע אולימפי היו גם ערכי null. מכיוון שלא דיברנו על טיפול ב null בתיאוריית התכנון, ניתן להניח שכל השדות תמיד מקבלות ערך שאינו null. כמו כן, שינינו מעט את הסכמה על מנת לפשט את השאלה. (הטבלה המקורית אפילו לא הייתה בצורה נורמלית ראשונה.)

נתון היחס athlete\_events עם הסכמה הבאה :

athlete\_events (ID, Name, Sex, Age, Height, Weight, Team,  
NOC, Year, Season, City, Sport, Event, Medal)

נתונה קבוצת התלויות הפונקציונליות הבאה מעל הסכמה של athlete\_events :

$F = \{$  ID  $\rightarrow$  Name, Sex  
Year, Season  $\rightarrow$  City  
ID, Year, Season, City  $\rightarrow$  Name, Sex, Age, Height, Weight, Team, NOC  
Event  $\rightarrow$  Sport  
Team  $\rightarrow$  NOC  
NOC, Year  $\rightarrow$  Team  
ID, Year, Season, Team, NOC, Event  $\rightarrow$  Sport, Medal  $\}$

ענו על הסעיפים הבאים :

1. מצאו את כל המפתחות של הטבלה athlete\_events. יש רק מפתח אחד ID, Year, Season, Event.

2. מה הצורה הנורמלית של הטבלה athlete\_events? לא 3NF ולא BCNF.

3. נתון פירוק של athlete\_events לתתי סכמות:

$R_1 = (ID, Year, Season, Name, Sex, Age, Height, Weight)$

$R_2 = (ID, Year, Season, City, Team, Event, Sport, Medal)$

$R_3 = (Team, NOC)$ .

האם הפירוק הוא ללא אובדן? נמקו! כן

4. מצאו כיסוי מינימאלי ל-F.

$F = \{$   
ID  $\rightarrow$  Name  
ID  $\rightarrow$  Sex  
Year, Season  $\rightarrow$  City  
ID, Year, Season  $\rightarrow$  Age  
ID, Year, Season  $\rightarrow$  Height  
ID, Year, Season  $\rightarrow$  Weight  
ID, Year, Season  $\rightarrow$  Team  
Event  $\rightarrow$  Sport  
Team  $\rightarrow$  NOC  
NOC, Year  $\rightarrow$  Team  
ID, Year, Season, Event  $\rightarrow$  Medal  
 $\}$

5. מצאו פירוק של athlete\_events ל-3NF על פי האלגוריתם הנלמד בכיתה.  
לכל אחד מתת הסכמות בפירוק, כיתבו מה הצורה הנורמלית.

$R_1 = (ID, Name)$  BCNF  
 $R_2 = (ID, Sex)$  BCNF  
 $R_3 = (Year, Season, City)$  BCNF  
 $R_4 = (ID, Year, Season, Age)$  BCNF  
 $R_5 = (ID, Year, Season, Height)$  BCNF  
 $R_6 = (ID, Year, Season, Weight)$  BCNF  
 $R_7 = (ID, Year, Season, Team)$  BCNF  
 $R_8 = (Event, Sport)$  BCNF  
 $R_9 = (NOC, Year, Team)$  3NF  
 $R_{10} = (ID, Year, Season, Event, Medal)$  BCNF

6. מצאו פירוק של athlete\_events ל-BCNF על פי האלגוריתם הנלמד בכיתה.

$R_1 = (ID, Name, Sex)$   
 $R_2 = (Year, Season, City)$

$R_3 = (ID, Year, Season, Age, Height, Weight, Team)$

$R_4 = (Event, Sport)$

$R_5 = (NOC, Team)$

$R_6 = (ID, Year, Season, Event, Medal)$

7. האם הפירוק שמצאתם בסעיף הקודם (ח) משמר תלויות? נמקו.

לא, התלות  $NOC, Year \rightarrow Team$  לא נשמרת.

## שאלה 2 (35 נקודות)

בשיעור למדנו ששמירת נתונים בטבלה בצורה נורמלית גבוהה (BCNF או 3NF) הוא חשוב, על מנת למנוע הכנסה לטבלה של נתונים שאינם עקביים. בשאלה זו אתם תתנסו בהתמודדות עם מידע אמיתי שלא נשמר בצורה נורמלית טובה. כאשר מעוניינים לבצע אנליזה על מאגר מידע נתון, שלב חשוב בתחילת התהליך הוא ניקוי המידע מהשגיאות שנמצאות בו.

לצורך התרגיל, אנחנו נשתמש במידע שנמצא בקישור: [Amazon Top 50 Bestselling Books 2009 - 2019](#) | [Kaggle](#). טבלה זו מכילה מידע על ספרים שהיו רבי מכר ב Amazon בשנים 2009 עד 2019. ניתן גם להוריד את המידע הזה מאתר הקורס. הטבלה מכילה את העמודות הבאים:

- Name, שם הספר.
- Author, שם המחבר.
- User Rating, דירוג הקוראים ב Amazon
- Reviews, מספר הביקורות של הספר ב Amazon
- Price, מחיר הספר בדולרים
- Year, השנה בו דורג הספר להיות רב מכר
- Genre, הז'אנר של הספר (fiction או nonfiction)

לפי התיאור של המאגר, אמורים להתקיים ההנחות הבאות:

1. לספרים שונים יש שמות שונים.
2. לכל ספר יש בדיוק מחבר אחד.
3. כל ספר שייך לז'אנר אחד.
4. הנתונים Rating, Num\_Reviews, Price נלקחו מ Amazon פעם אחת, בתאריך 13/10/20.
5. ספר יכול להופיע ברשימת רבי המכר ביותר משנה אחת.

ענו על השאלות הבאות:

א. כתבו את קבוצת התלויות הפונקציונליות שאמורות להתקיים בטבלה לפי כל ההנחות הנ"ל. כתבו את התלויות בצורה אטומית, כלומר שבצד ימין של כל תלות יופיע רק שדה אחד. אין לציין תלויות טריוויאליים.

$F = \{$  Name  $\rightarrow$  Author  
Name  $\rightarrow$  Rating  
Name  $\rightarrow$  Num\_Reviews

Name → Price  
Name → Genre }

ב. מה המפתח של הטבלה? אם יש מספר מפתחות, ציינו את כולם. המפתח הוא Name, Year.  
ג. מה הצורה הנורמלית של הטבלה? נמקו. המפתח הוא Name, Year ולכן הצורה הנורמלית היא לא 3NF

- ד. בסעיף זה נבחן אילו תלויות פונקציונליות מתקיימות בפועל במופע של הטבלה באתר kaggle ואילו לא מתקיימות. כלומר, אנחנו נגלה את בעיות העקביות של הנתונים. כדי לעשות זאת:
1. כתבו קובץ create.sql שמייצר טבלה בשם bestsellers עם כל העמודות הנתונות וללא אילוצים בכלל.
  2. טענו את הנתונים מהקובץ לתוך הטבלה (בצורה הרגילה, שתוארה בתרגילים קודמים). שימו לב: ההכנסה תהיה פשוטה יותר אם לפני כן תורידו את כל הפסיקים שמופיעים בשמות הספרים. בגרסה באתר הקורס הורדנו את הפסיקים עבורכם.
  3. כתבו שאילתת SQL בקובץ contradictions.sql שמחזירה את כל השורות המעורבות בסתירה של תלות פונקציונלית. על השאילתה להחזיר רק את העמודות שם ספר, מחבר הספר ושנה, ממוינים לפי שם ספר, ואח"כ שנה, בסדר עולה וכן, להחזיר כל שורה רק פעם אחת.
  4. כתבו קובץ drop.sql שמוחק את הטבלה.
- ה. אלו תלויות פונקציונליות שכתבתם בסעיף א מתקיימות בנתונים, ואילו תלויות מופרות? תנו פירוק מומלץ של הטבלה לתתי יחסים והסבירו איך שמירת הנתונים בפירוק היה מונע הכנסת שורות לא קונסיסטנטיות.

התלויות המופרות הן:

Name → Rating  
Name → Num\_Reviews  
Name → Price

הפירוק המוצע:  $R1 = (\text{Name}, \text{Author}, \text{Rating}, \text{Reviews}, \text{Price}, \text{Genre})$ ,  $R2 = (\text{Name}, \text{Year})$ .  
בפירוק זה, מכיוון שName הוא מפתח בR1, לא יהיה ניתן להכניס מידע שונה לאותו ספר כי לכל ספר תהיה שורה אחת בלבד.

שאלה 3 (30 נקודות)

- א. תנו דוגמה פשוטה ליחס, תלויות פונקציונליים ופירוק כך שהפירוק משמר תלויות אך עם אובדן. אם לא קיימת דוגמה כזאת, הסבירו.

$R = (A, B, C)$ ,  $F = A \rightarrow B$   
 $R1 = (A, B)$ ,  $R2 = (C)$

ב. תנו דוגמה פשוטה ליחס, תלויות פונקציונליים ופירוק כך שהפירוק אינו משמר תלויות אך ללא אובדן. אם לא קיימת דוגמה כזאת, הסבירו.

$R = (A, B, C)$ ,  $F = A \rightarrow B, CB \rightarrow A$   
 $R_1 = (A, B)$ ,  $R_2 = (A, C)$

ג. נתון יחס  $R = (A, B, C, D, E)$  וקבוצה  $F$  של תלויות פונקציונליים  $F$ . ידוע שצד ימין של כל אחד מהתלויות ב- $F$  הוא בדיוק האטריבט  $A$ . הוכח או תן דוגמה נגדית:  $R$  ב-BCNF אם ורק אם  $R$  ב-3NF.

Correct.

Case 1: All the dependencies in  $F$  are trivial. Then  $R$  is both in 3NF and in BCNF

Case 2: There is some non-trivial dependency in  $F$ . Since  $B, C, D, E$  do not appear on the right side of any dependency, they must be in every key. Clearly,  $A$  is in the closure of  $BCDE$ , due to the non-trivial dependency. Therefore, the only key is  $BCDE$ . We have shown already that any relation with a single key is in BCNF if and only if it is in 3NF.