Design Theory : 5 תרגיל

תאריך הגשה: 55:25, 25.06.23

הוראות הגשה:

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

- ex5.pdf עם התשובות מפורטות לשאלות.
- רמתאים לשאלה 3 סעיף ד.1. create.sql •
- מתאים לשאלה 3 סעיף ד.3. contradictions.sql
 - .4. המתאים לשאלה 3 סעיף ד.4. drop.sql •
- 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.

תזכורת: יש להגיש תרגיל מוקלד בלבד.

שאלה 1 (30 נקודות)

 ${
m R}$ מוסיפים ל R אטריבוט מבלי לבצע אף שינוי של F. מוסיפים ל פונקציונליות פונקציונליות פונקציונליות עם קבוצת הבאות. לכל טענה ענו דוגמה שמראה היתכנות של הטענה או הוכיחו שהטענה שגויה.

- .BCNF אי משאר ב R וגם לאחר השינוי B נשאר ב BCNF אי אי טענה: יתכן שלפני הוספת האטריבוט, R
- ולאחר (BCNF- או ב-3NF) אוב, היה לפחות האטריבוט אוב, R ב. טענה יתכן שלפני הוספת האטריבוט אוב היה לפחות ב-3NF ב. אינו ב-3NF ואינו ב-3NF ואינו ב-3NF השינוי אוב בר אינו ב-3NF ואינו ב-3NF השינוי אוב בר אינו ב-3NF ואינו ב-3NF השינוי אוב בר אינו ב-3NF היה לפחות ב-3NF היה לפח
- BCNFב R ולאחר השינוי, BCNF אבל אבל אבל אבל אבל אבל הוספת האטריבוט, R. היה ב-3NF אבל אבל אבל אונה BCNF

<u>שאלה 2 (35 נקודות)</u>

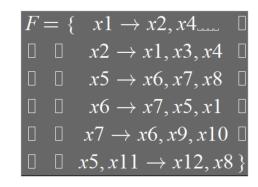
נחזור וניזכר במידול מידע על ההרשמה לאוניברסיטאות מהתרגיל בית הראשון. הפעם, במקום למדל בעזרת דיאגרמת ישויות קשרים, נשתמש בגישת תיאוריית התכנון על מנת להבין איך יש להפריד טבלה אחת גדולה לתתי טבלאות.

<u>הערה :</u> בטבלה המקורית של מידע ההרשמה היו גם ערכי null. מכיוון שלא דיברנו על טיפול ב null בתיאוריית התכנון, ניתן להניח שכל השדות תמיד מקבלות ערך שאינו null.

נתון היחס enrollment עם הסכמה הבאה: (היחס דומה לטבלה מתרגיל 1 אך הורדנו כמה שדות לשם הפשטות של התרגיל*)*

enrollment (country, countrycode, region, incomegroup, iau_id1, eng_name, orig name, foundedyr, latitude, longitude, year, students5 estimated)

נתונה קבוצת התלויות הפונקציונליות הבאה מעל הסכמה של enrollment:



```
F = { country → countrycode, incomegroup
countrycode → country, region, incomegroup
iau_id1 → eng_name, orig_name, foundedyr
eng_name → orig_name, iau_id1, country
orig_name → eng_name, latitude, longitude
iau_id1, year → students5_estimated, foundedyr }
```

: ענו על הסעיפים הבאים

- א. מצאו את כל המפתחות של הטבלה enrollment. (אין צורך לפרט את החישובים)
 - ב. מה הצורה הנורמלית של הטבלה enrollment?
 - ג. נתון פירוק של enrollment לתתי סכמות:

 $R_1 = (iau_id1, eng_name, orig_name, country)$

R₂ = (iau_id1, year, students5_estimated, foundedyr)

 R_3 = (country, countrycode, region, incomegroup)

 $R_4 = (eng_name, latitude, longitude)$

האם הפירוק הוא ללא אובדן ? נמקו בעזרת תוצאת האלגוריתם שנלמד בכיתה.

- ד. מצאו כיסוי מינימאלי ל-F. (אין צורך לפרט את החישובים)
- ה. מצאו פירוק של enrollment ל-3NF על פי האלגוריתם הנלמד בכיתה. לכל אחד מתת הסכמות בפירוק, כתבו מה הצורה הנורמלית.
- . מצאו פירוק של enrollment ל-*BCNF* על פי האלגוריתם הנלמד בכיתה. פרטו את השלבים של האלגוריתם, את התלות על פיה פירקתם בכל שלב, והדגישו את התשובה הסופית.
 - ו. האם הפירוק שמצאתם בסעיף הקודם (ו) משמר תלויות? אם לא, פרטו אילו תלויות אינן נשמרות.

שאלה 3 (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.

 $\frac{Amazon\ Top\ 50\ Bestselling\ Books\ 2009\ -\ 2019\ I}{Amazon\ Large}$ במידע שנמצא בקישור: על ספרים שהיו רבי מכר בי 2009 עד 2019. ניתן גם להוריד באות: מכילה מידע על ספרים שהיו רבי מכר בי $\frac{Kaggle}{Amazon}$ את המידע הזה מאתר הקורס. הטבלה מכילה את העמודות הבאות:

- .Name •
- .Author Author
- Amazon, דירוג הקוראים בUser_Rating ●
- Reviews, מספר הביקורות של הספר בReviews
 - Price , מחיר הספר בדולרים
 - א רב מכר להיות רב מכר, Year •
- (nonfiction או fiction) הזיאנר של הספר, Genre •

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 לזיאנר אחד.
- .13/10/20 פעם אחת, בתאריך User_Rating, Reviews, Price לקחו מuzon.
 - 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.
 - ב. מה המפתח של הטבלה? אם יש מספר מפתחות, ציינו את כולם.
 - :. מה הצורה הנורמלית של הטבלה! נמקו.
- 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:
- טם כל העמודות הנתונות וללא אילוצים create.sql שמייצר טבלה בשם create.sql שמייצר טבלה בשם בכלל.
- 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.
 - 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?
- ו. תנו פירוק מומלץ של הטבלה לתתי יחסים והסבירו איך שמירת הנתונים בפירוק היתה מונעת הכנסת שורות לא קונסיסטנטיות.

בהצלחה!