

## מבחן במערכות בסיסי נתונים

רובי בוים

סמסטר א' מועד א' - 07 בפברואר, 2010  
 סמסטר א' מועד ב' - 20 בספטמבר, 2010  
 חומר סגור, שלוש שעות  
 סה"כ 100 נקודות

נא לכתוב את כל התשובות על הדפים  
 (המבחן כולל 13 דפים כולל דף זה)

מס מחברת: 40 מס סטודנט:

חלק	ציון
1	40
2	20
3	15
4	25
סה"כ	100

## 1. [35 נקודות]

כיוון שאתם סטודנטים מצטיינים, מארק צוקרברג החליט לאפשר לכם הצצה חלקית וחד פעמית לבסיס הנתונים של FACEBOOK:

**Users**(uid, name, birthday, country)

**Events**(eid, name, uid, date)

**Guests** (eid, uid, status)

לכל משתמש יש מזהה ייחודי, שם, תאריך לידה ומדינה. עבור כל אירוע המערכת שומרת את מזהה ייחודי ובנוסף את שם האירוע, המשתמש שערך את האירוע והתאריך. בנוסף, אירועים יכולים להכיל רשימת מוזמנים שלכל אחד סטטוס הגעה (yes/no/maybe).

א. מארק (uid=123) רוצה לברר למה צלמי הפפראצי לא מגיעים לאירועים שהוא עורך. עזרו לו לכתוב שאילתה שתחזיר את רשימת כל האנשים (מזהה ושם) שערכו אירוע בזמן מקביל.

```
SELECT DISTINCT u.uid, u.name
FROM Users u, Events E1, Events E2
WHERE E1.uid = '123' AND E2.uid <> '123' AND
      E1.date = E2.date AND
      E2.uid = u.uid
```



ב. כתבו את אותה שאילתה באלגברה רלציונית.

$$\Pi_{uid, name} \left[ \left[ \rho_{E1. uid, date} \left( \Pi_{uid, date} \left( \sigma_{uid=123} (Events) \right) \right) \bowtie_{date=date} \Pi_{uid, date} \left( \sigma_{uid < 123} (Events) \right) \right] \bowtie_{uid=uid} \left[ \Pi_{uid, name} (Users) \right] \right]$$



ג. בעקבות הצלחתכם (בתקווה 😊) בסעיפים הקודמים, מארק מתעניין בהשקעה חדשה בשוק הישראלי. אך הוא עדיין מודאג שFACEBOOK אינה מספיק פופולארית בקרב המשתמשים המקומיים. עזרו לו לכתוב שאלתה שתחזיר את מספר האירועים שנערכו על ידי משתמשים ישראלים אשר הכילו לפחות 100 אנשים (כלשהם) שהסטטוס שלהם הוא yes או maybe.

```
SELECT COUNT(E.eid)
FROM Events E, Users U
WHERE U.country = 'Israel' AND U.uid = E.uid AND
      E.eid IN (SELECT E2.eid
                FROM Events E2, Guests G
                WHERE E2.eid = G.eid AND
                     (G.status = 'yes' OR
                      G.status = 'maybe'))
GROUP BY E2.eid
HAVING COUNT(G.*) >= 100)
```



ד. התבוננו בזוגות Q1, Q2 של השאלות הבאות. לכל זוג סמנו מי מהמשפטים הבאים מתקיים:

- (A) שתי השאלות מחזירות תמיד בדיוק אותה תשובה
- (B) התשובות המחוזרות ע"י Q1 מוחזרות גם ע"י Q2, אבל יתכן ש Q2 תחזיר ערכים נוספים או כפילויות של ערכים.
- (C) התשובות המחוזרות ע"י Q2 מוחזרות גם ע"י Q1, אבל יתכן ש Q1 תחזיר ערכים נוספים או כפילויות של ערכים.
- (D) אף אחד מהנ"ל

דוגמה 1:

Q1 = **SELECT \* FROM events WHERE date > 01-01-2010**

Q2 = **SELECT \* FROM events WHERE date > 01-01-2009**

עליכם לענות (B)

דוגמה 2:

Q1 = **SELECT name FROM Guests**

Q2 = **SELECT DISTINCT name FROM Guests**

עליכם לענות (C)

i.

Q1 = **SELECT** eid  
**FROM** Events

Q2 = **SELECT** eid  
**FROM** Guests

הערה: אפשר להניח שמוגדר foreign key בטבלת Guests על שדה eid

(A)

(B)

(C)

(D)

ii.

Q1 = **SELECT** U.uid  
**FROM** Users U  
**WHERE** 1 >=  
( **SELECT** count(\*)  
**FROM** Events E, Guests G  
**WHERE** E.eid = G.eid **AND**  
E.uid = U.uid **AND**  
G.uid = U.uid )

Q2 = **SELECT** E.uid  
**FROM** Events E, Guests G  
**WHERE** E.uid = G.uid **AND**  
E.eid = G.eid

(A)

(B)

(C)

(D)

iii.

Q1 = **SELECT**   **DISTINCT** U.uid, U.name  
     **FROM**     Users U, Guests G1, Guest G2  
     **WHERE**    G1.uid     =    123        **AND**  
                  G1.eid     =    G2.eid    **AND**  
                  G2.uid     =    U.uid

Q2 = **SELECT**   uid, name  
     **FROM**     Users  
     **WHERE**    uid        **IN**  
                  (SELECT uid FROM Guests WHERE eid **IN**  
                          (SELECT eid FROM Guests WHERE uid = 123)

(A)

(B)

(C)

(D)



2. [20 נקודות]

נתונה הטבלה  $R(A, B, C, D, E)$  עם התלויות הפונקציונאליות הבאות:

- $A \rightarrow C$
- $B, C \rightarrow E$
- $B \rightarrow D$

א. מצא את המפתח המינימאלי של טבלה R

המפתח המינימאלי הוא  $\{A, B\}$  ✓

ב. פרקו את הטבלה לטבלאות BCNF ונמקו את תשובתכם. התשובה:

צריכה להכיל:

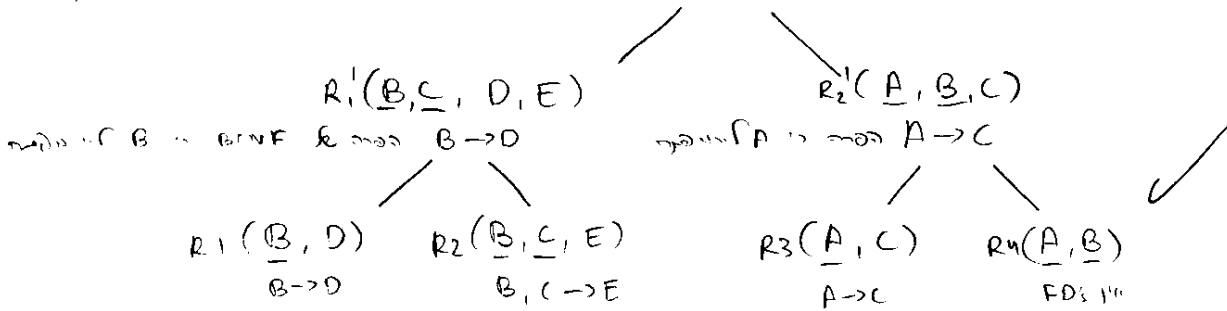
- רשימת טבלאות והשדות בכל טבלה
- המפתחות בכל אחת מהטבלאות

תלויות פונקציונאליות:

$B, C \rightarrow D, E$

$R(\underline{A}, \underline{B}, C, D, E)$

$B, C \rightarrow D, E$  - הפכה ל BCNF כי  $B, C$  אינם מפתח



הגשימה יסודית:  $R_1(B, D)$  עם המפתח  $\{B\}$ ,  $R_2(B, C, E)$  עם המפתח  $\{B, C\}$

$R_3(A, C)$  עם המפתח  $\{A\}$ ,  $R_4(A, B)$  עם המפתח  $\{A, B\}$  בהם אלה נחלקו לפי

FD שמתחילים מ- BCNF, וצריכים להבטיח את ה- FD's המקוריים (ה- non-trivial) שמאזינים  $\{B, C\}$  ו-  $\{A, B\}$  חלשי צינור. BCNF. המפתחות שלהם (שהם זהה ל- FD's אחרים).

ג. עבור כל טבלה שקיבלתם (בפרוק הסופי) ציינו האם היא במקרה גם ב

3NF כיין  $3NF \subset BCNF$  ויש זכרון של טבלה שקיבלנו

בפירוק הסופי, טבלה א- ב- BCNF, הן

הפכה גם ב- 3NF ✓



### 3. [20 נקודות]

נתון מסמך ה XML הבא המתאר חנות משחקי ווידאו. המוצרים נחלקים לפי סוג קונסולת המשחק (XBOX, WII, PLAYSTATION..) אשר לכל אחת יש שם, יצרן ורשימת משחקים. עבור כל משחק שומרים את השם ואת מספר המשתתפים המקסימלי האפשרי.

#### Input.xml

```
<rubi-store>
  <console>
    <name>Xbox360</name>
    <manufacture>Microsoft</manufacture>
    <games>
      <game players=1>halo</game>
      <game players=4>guitar hero 5</game>
      ...
    </games>
  </console>
  <console>
    <name>Wii</name>
    ...
  </console>
  ...
</rubi-store>
```

א. כתוב ביטוי XPath שמחזיר את כל המשחקים בעלי אפשרות משחק ל 4 שחקנים בדיוק עבור קונסולת ה Xbox360. התשובה צריכה להראות בסגנון הבא:

guitar hero 5  
call of duty 4

...

//console [name/text() = 'xbox360']/games/game [@players = '4']/text()



ב. כתוב ביטוי XQuery שמחזיר את רשימת הקונסולות ואת כמות המשחקים המתאימים שלהם. סננו קונסולות עם פחות מ 50 משחקים. התשובה צריכה להראות בסגנון הבא:

```
<result>
  <ConsoleName>Xbox360</ConsoleName >
  <ConsoleCount>322</ConsoleCount>
```

</result>

```
<result>
  <ConsoleName>Wii</ConsoleName >
  <ConsoleCount>65</ConsoleCount>
```



</result>

...

(הנחיה: יזון כפוליה קונסולות ומחזיר קונסולות כ"ן הפוליה מסוג משחקים)

```
FOR $c IN document('input.xml')//console
LET $g := count(document('input.xml')//console [name/text() = $c/name/text()])
WHERE $g >= 50
RETURN $ <result>
  <ConsoleName> $c/name/text() </ConsoleName>
  <ConsoleCount> $g </ConsoleCount>
</result> }
```

4. [25 נקודות]

נתונות שתי הטבלאות  $R(\underline{A}, B, C)$   $S(\underline{D}, E, F)$  עם האינדקסים הבאים:

א. על השדה R.A (מסוג unclustered)

ב. על השדה R.B (מסוג clustered)

ג. על השדה S.D (מסוג clustered)

בניח ששתי הטבלאות גדולות ולא ניתן להחזיק אותן בזיכרון, אבל האינדקסים קטנים וכן יושבים בזיכרון בשלמותם.

א. נתבונן בשאלתה הבאה:

**SELECT \* FROM R,S WHERE R.C = S.D**

בניח שמימוש ה JOIN בשאלתה הוא ע"י אלגוריתם Index nested join.

האם ישפיע שינוי האינדקס על שדה S.D ל unclustered על ביצועי

המערכת בחישוב השאלתה?

(במקו בקצרה – שתי שורות לכל היותר!)

השינוי לא ישפיע כיוון ש-D הוא אינדיקס ב-S, ולכן הוא מקרה יחיד צמוד ל-B ואגב  
המקרה של הדיק המאגים (כיוון שכל האינדקסים בזיכרון, אלו 2/0 לא ישפיעו על הביצועים).

ב. נתון המידע הסטטיסטי הבא:

(מספר הדפים בטבלה R)	B(R) = 100
(מספר הרשומות בטבלה R)	T(R) = 2000
(מספר הערכים השונים בשדה B בטבלה R)	V(R,B) = 20
(מספר הדפים בטבלה S)	B(S) = 20
(מספר הרשומות בטבלה S)	T(S) = 1000
(גודל הזכרון בדפים)	M=80

למקרה B.C  
hash join

נתבונן בשאלתה הבאה:

**SELECT \* FROM R,S WHERE R.C = S.D AND R.B = 123**

עבורה, נתונות תוכניות הביצוע הלוגיות הבאות:

$P1 = S \text{ join}_{D=C} (\text{select}_{B=123}(R))$

$P2 = \text{select}_{B=123}(S \text{ join}_{D=C} R)$

נניח שכל אחת מתוכניות הביצוע  $\{P1, P2\}$  מחושבת ב pipeline, כלומר תוצאת אופרטור אחד ניתנת ישירות לאופרטור הבא (בלי להיכתב לדיסק בדרך).

יש שני אופרטורים בסיסים:

~~Sort-Merge~~ JOIN - נניח שממומש על ידי <sup>hash</sup> Sort-Merge Join

SELECT - יכול להיות ממומש על ידי האינדקס או ע"י מעבר סדרתי על הטבלה. כאשר הוא לא האופרטור הראשון, הוא ניתן לביצוע ישירות על ה tuple שהוא מקבל כקלט (בלי תשלום של I/O).

חשבו את מחיר כל אחת מתוכניות הביצוע במונחים של פעולות I/O (לא כולל כתיבת התוצאה הסופית לדיסק). בטאו כל אחת מהתשובות כפונקציה של הסטטיסטיקות הנ"ל. לדוגמה:

$COST(P1) = B(R)B(S) / V(R,A)$  (לא התשובה האמיתית!)

הסבירו כל תשובה בקצרה (2 שורות מקסימום)

$$\begin{aligned} \text{A. } COST(P1, \text{select ע"י אינדקס}) &= \frac{B(R)}{V(R,B)} + 3 \cdot B(S) + 2 \cdot \frac{B(R)}{V(R,B)} = \\ &= \frac{100}{20} + 3 \cdot 20 + 2 \cdot \frac{100}{20} = 5 + 60 + 10 = 75 \end{aligned}$$

האינדיקס יהיה מהלפי ה-select עם אינדקס clustered, ושני האינדיקס הלא-ממוינים הם  
לפי ה-hash join, הרישוי הוא אינדיקס סלקט/אינדיקס (R) הוא ק  $\frac{B(R)}{V(R,B)}$ .

B.  $\text{COST}(P1, \text{מעבר סידרתי select}) = B(R) + 3B(S) + 2 \cdot \frac{B(R)}{\sqrt{R \cdot B}} =$   
 $= 100 + 3 \cdot 20 + 2 \cdot \frac{100}{20} = 100 + 60 + 10 = 170$

היזכרתי שהיה עלול להיות select עם סינון היינקה - למה לא? כי הולך ר. חסר היינקה  
 הם עלול להיות hash, וזה עדיין R קטנות יותר מ  $\frac{B(R)}{\sqrt{R \cdot B}}$  יחדיו ה select מנסה.



C.  $\text{COST}(P2, \text{ע"י אינדקס select}) = 3 \cdot B(S) + 3 \cdot B(R) = 3 \cdot 20 + 3 \cdot 100 = 360$

זהו אולי יגידו כי ה hash הוא יותר יעיל מסינון אינדקסים. כן, זה select יחסית  
 כי למה עדיין, לפי גודל הולך יותר או אולי I/O יאלו.



D.  $\text{COST}(P2, \text{מעבר סידרתי select}) = 3 \cdot B(S) + 3 \cdot B(R) = 360$

חשוב לא יזהר לקבל, שוב, נדע שה select היא אולי אל פאזה I/O.  
 חשב ה hash הוא יותר יעיל.



ג. ציינו מהי התוכנית הזולה ביותר לחישוב השאילתה, ואת מחירה (מבוטא כמספר)

לפי הנגזרת שהקבלו נראה כי היעיל ביותר היא ה select  
 השאילתה היא P1 הוצע ה select נראה עדיין יעיל יותר.  
 ופאזה יק היא : 75 פאזה I/O.

