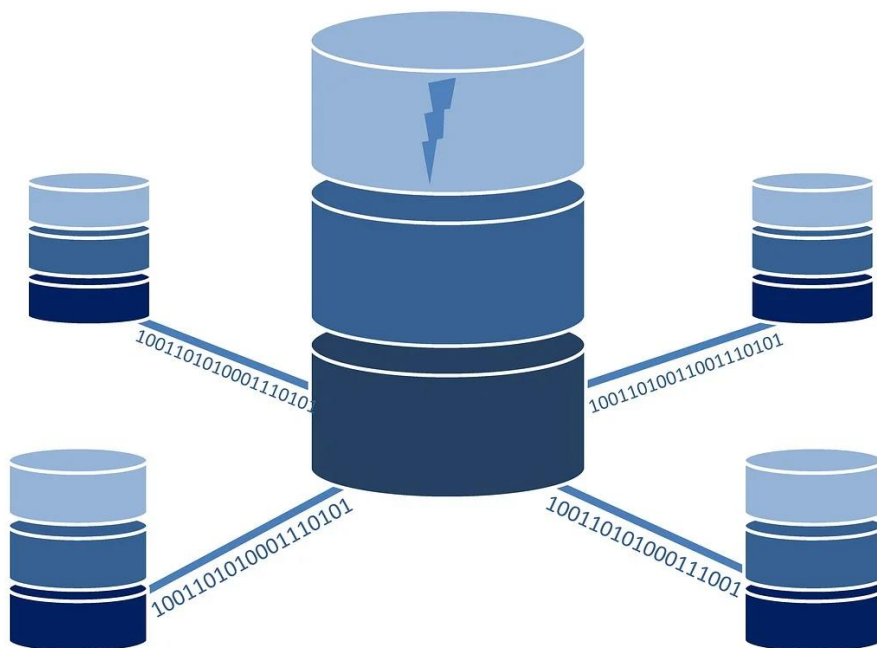


# פרויקט במדידת חוסר עקביות במסדי נתונים

דו"ח סיום

## Inconsistency Measures



מגישה: רינה קוצ'ירגן  
תעודת זהות: 316465707

3.....	מבוא.....
5.....	כלים וסביבת פיתוח.....
7.....	הקלט לתוכנית.....
9.....	אלגוריתם יצירת ההפרות.....
10 .....	מימוש המדדים.....
13 .....	זמני הריצה ומגבלות התוכנית.....
14 .....	גרפים ותוצאות.....
14 .....	<b>Voters</b>
15 .....	<b>Adult</b>
16 .....	<b>Food</b>
17 .....	<b>Tax</b>
18 .....	<b>Airport</b>
19 .....	<b>Flight</b>
20 .....	<b>Hospital</b>
21 .....	<b>Stock</b>
22 .....	ניתוח התוצאות.....
23 .....	ביבליוגרפיה.....

הפרויקט נערך במסגרת מחקר<sup>1</sup> של בני קימלפלד, אסתר ליבשיץ, איהב אליאס, סודיפה רוי ושגב צור העוסק במדידת אי עקביות במסדי נתונים. במסגרת המחקר נחקרו 6 מדדים המשמשים למדידת חוסר עקביות במסדי נתונים, מטרת החקירה היא לאפיין כל מדד ולבחון את התנהגותו תוך כדי הכנסת הפרות מלאכותיות אל תוך מסד הנתונים.

ששת המדדים שמומשו ונחקרו הם :

- **Drastic inconsistency value  $I_D$** .1

מדד בינארי שערכו 1 במידה והמסד איננו עקבי (קיימות הפרות במסד נתונים) ו-0 אחרת.

- **Set of all minimal inconsistent subsets of  $D$   $I_{MI}$** .2

מספר זוגות הרשומות המופיעות יחד בהפרה כלשהי במסד הנתונים.

- **Problematic facts  $I_P$** .3

מספר הרשומות המשתתפות בהפרה כלשהי במסד הנתונים.

- **Minimal cost of a sequence of operations that repairs the database  $I_R$** .4

המספר המינימאלי של רשומות שיש להסיר בכדי שמסד הנתונים יהיה עקבי. מדד זה מחושב באמצעות תכנון לינארי המוגדר באופן הבא :

$$\begin{aligned} \text{minimize } \left\{ \sum_{i=1}^n x_i \right\} \quad & \text{s.t. } i \text{ represents a tuple in the database} \\ & x_i \in \{0,1\} \end{aligned}$$

$$x_i + x_j \geq 1 \quad \text{for each two tuples } i, j \text{ that violate a constraint}$$

- **Linear relaxation of  $I_R$   $I_R^{lin}$** .5

מדד זה כמעט זהה למדד הקודם מלבד הקלה על התחום האפשרי עבור כל  $x_i$ .

$$\begin{aligned} \text{minimize } \left\{ \sum_{i=1}^n x_i \right\} \quad & \text{s.t. } i \text{ represents a tuple in the database} \\ & 0 \leq x_i \leq 1 \end{aligned}$$

$$x_i + x_j \geq 1 \quad \text{for each two tuples } i, j \text{ that violate a constraint}$$

<sup>1</sup> [Principles of Progress Indicators for Database Repairing](#), Ester Livshits, Ihab F. Ilyas, Benny Kimelfeld, Sudeepa Roy, 13 Apr 2019

– The set of all maximal consistent subsets of  $D$   $I_{MC}$ .6

מדד זה מתבסס על גרף הקונפליקטים של מסד הנתונים. גרף הקונפליקטים הוא גרף בו הצמתים מהווים רשומות וישנה קשת בין שני צמתים המייצגים רשומות שמשתתפות בהפרה כלשהי זו עם זו.

עבור הגרף המשלים לגרף הקונפליקטים – כלומר הגרף בו ישנה קשת בין שני צמתים המייצגים רשומות שאינן משתתפות בהפרה זו עם זו, המדד מחשב את המספר המקסימלי של קליקות בגרף.



- סביבת הפיתוח בה השתמשתי לכתיבת הקוד והרצתו היא **jupyter lab** גרסה 1.2.6

- הקוד נכתב בשפת פייתון גרסה 3.7



- לצורך מימוש המדדים  $I_R, I_R^{lin}$  השתמשתי בכלי המאפשר פתירת בעיות בתכנות לינארי הנקרא **Gurobi optimizer** גרסה 9.0

### מהו תכנות לינארי?

תכנות לינארי (Linear programming) הוא למעשה מציאת פתרון אופטימלי אשר מביא למינימום/ למקסימום את פונקציית המטרה של התוכנית.

בעיית תכנות לינארי מורכבת ממספר מאפיינים חשובים המאפשרים את פתרונה:

1. קבוצת משתנים (Variants) :  $x_1, \dots, x_n$  כאשר לכל משתנה יש טווח (למשל: בינארי או רציף בתחום כלשהו או קבוצת ערכים סופית)

2. קבוצת אי שוויוניים (Constraints) מהצורה :  $x_i + x_j + x_m + \dots \leq Y$

3. פונקציית מטרה : הערך אותו אנו מעוניינים להביא למינימום/למקסימום, למשל :

$$\text{minimize } \left\{ \sum_{i=1}^n x_i \right\}$$

בכדי לפתור את בעיית התכנות הלינארי יש להגדיר כל אחד מהמאפיינים הללו באמצעות פונקציות ייעודיות של Gurobi ולבסוף ניתן לבצע את החישוב באמצעות פונקציית `optimize`.

- לצורך מימוש המדד  $I_{MC}$  המחשב את מספר הקליקות המקסימלי בגרף המשלים לגרף הקונפליקטים השתמשתי באלגוריתם `parallel enum`<sup>2</sup> האלגוריתם מקבל כקלט גרף בפורמט NDE ומחשב את מספר הקליקות המקסימלי בגרף זה.

- ספריות נוספות שהשתמשתי בהן :

**Pandas** – ספרייה המשמשת לאנליזה ומניפולציות על מידע. מבנה הנתונים שמאחסן את המידע מכונה `dataframe` והוא מאפשר מגוון של פעולות על מידע. באמצעות `dataframes` טענתי את מסד הנתונים מה-CSV לקוד ולאחר מכן השתמשתי בפונקציות הייעודיות להוצאת רשומות, עדכון רשומות, מניה ועוד.



<sup>2</sup> [Parallel enum algorithm GitHub](#)

**Pandasql** - ספרייה המאפשרת הרצת שאילתות בשפת SQL על dataframes של Pandas. באמצעות כלי זה הרצתי שאילתות המוצאות את כל זוגות הרשומות המופיעות בהפרה זו עם זו וכל הרשומות המופיעות בהפרה כלשהי.

**matplotlib** – ספרייה המשמשת בתור כלי לוויזואליזציה של נתונים למשל: גרפים, תמונות, צורות ועוד. השתמשתי בכלי זה בפרויקט לצורך הצגת פלטי המדדים, הקוד מייצר הפרה במסד הנתונים ומעדכן את הרשומה עם ההפרה, לאחר כל שינוי המדדים מחושבים מחדש ולכל מדד שחושב התוכנית מייצרת גרף בו ציר ה- x הוא מספר השינויים במסד נתונים וציר ה- y הוא ערך המדד. ניתן לצפות בדוגמאות [בגרפים ותוצאות](#)

**ipywidgets** – היא ספרייה של Jupyter המאפשרת הוספה של ווידג'טים למשל: תיבת טקסט, בחירה מרובה, כפתורי רדיו, אנימציות, רכיבי HTML ועוד. באמצעות ספרייה זו יצרתי חווית קלט אינטראקטיבית למשתמש המעוניין להריץ את הקוד דרך jupyter lab כפי שאציג [בקלט לתוכנית](#)

**random** – ספרייה של פייתון שבאמצעותה שגב ואני מימשנו את האלגוריתם שמייצר את ההפרות, כחלק מיצירת ההפרה התבקשנו להגריל ערכים לרשומות כדי ליצור את ההפרות.

**Numpy** - ספרייה המאפשר חישובים מדעיים באמצעות פייתון.  NumPy

**Re** – ספרייה המאפשר פעולות על ביטויים רגולריים. השתמשתי בספרייה בכדי לבצע מניפולציות על מחרוזות ביתר קלות.

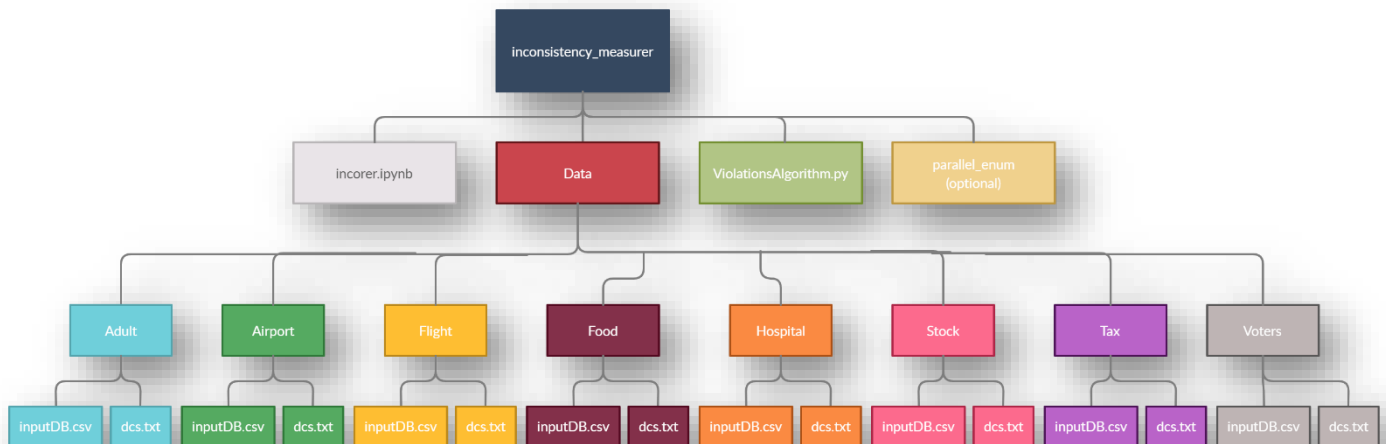
**Os** – ספרייה בפייתון המאפשרת מניפולציות הקשורות למערכת ההפעלה. השתמשתי בספרייה בכדי למצוא כתובות של קבצים בתיקיית הפרויקט במהלך הקוד ובנוסף כדי לקרוא, לכתוב, ליצור ולסגור קבצים.

**Subprocess** – ספרייה המאפשרת יצירת תהליכים חדשים והרצת תוכניות. השתמשתי בספרייה כדי להריץ את האלגוריתם <sup>3</sup>parallel\_enum המחשב את המדד  $I_{MC}$  (מקסימום קליקות בגרף המשלים לגרף הקונפליקטים) לאחר ההרצה חילצתי את הפלט והחזרתי את ערך המדד.

**Time** – ספרייה בפייתון המאפשרת פעולות הקשורות בזמן. השתמשתי בספרייה כדי לחשב את זמני הריצה של כל אחד מהמדדים ובפרט זמן הריצה של חישוב המדדים הכולל על המסד הנתונים המסוים.

---

<sup>3</sup> [Parallel\\_enum algorithm GitHub](#)



התוכנית מצפה להכיל בתיקיית הפרויקט בתיקיית Data/database\_name שני קבצים עיקריים לכל מסד נתונים :

## 1. מסד נתונים בפורמט CSV

## 2. קובץ תלויות על המסד

כל תלות היא מסוג Denial constraint כלומר :  $\forall t \neg (P_1, \dots, P_m)$   
 התלויות יופיעו בשורות נפרדות באופן הבא :  
 $not(t1.City \neq t2.City \& t1.Location = t2.Location)$   
 $not(t1.Longitude \neq t2.Longitude \& t1.Location = t2.Location)$

הקלט לתוכנית נוצר באמצעות HTML וסיפריית ipywidgets בכדי לספק למשתמש חוויה אינטראקטיבית ולבחור דרכים שונות לחשב את המדדים השונים .

כאן המשתמש בוחר את שם המסד שברצונו להריץ – אחד או יותר מהמסדים המסופקים עם הפרויקט או מסד אחר.

השלב הבא הוא לבחור האם להריץ את אלגוריתם ההפרות מספר פעמים ולחשב את המדדים לאחר כל הרצה.

```

HelloNewUser()

Welcome to the inconsistency measurer

Please specify the databases you wish to compute the measures on (seperated by ',')
The databases included are: Adult,Airport,Flight,Food,Hospital,Stock,Tax,Voters


Welcome to the inconsistency measurer

Please specify the databases you wish to compute the measures on (seperated by ',')
The databases included are: Adult,Airport,Flight,Food,Hospital,Stock,Tax,Voters


Do you wish to run a simulation that introduces random violations in the database? [y/n]
Default value is: y

  
```

במידה והמשתמש בחר להריץ את האלגוריתם עליו לספק את מספר הפעמים שהאלגוריתם ייצר הפרות ויעדכן את מסד הנתונים

כעת, המשתמש יבחר את המדדים שברצונו לחשב. ההודעה המוצגת בהמשך מיועדת לאלו הרוצים להריץ את המדד  $I_{MC}$  שמריץ אלגוריתם לחישוב מספר קלילות בגרף.

לאחר מכן יש ללחוץ על המשך בכדי להריץ את המדדים

כאשר התוכנית מתחילה לרוץ מודפסת הודעה למסך יחד עם מד התקדמות

לאחר סיום הרצת המדדים על כל מסד, הפלט מופיע בתיקייה עם שם המסד תחת השם מצוין. בתיקייה ניתן למצוא את הגרפים שנוצרו יחד עם קבצים המתארים את זמני הריצה ותוצאות המדדים.

Welcome to the inconsistency measurer

Please specify the databases you wish to compute the measures on (seperated by ',')  
The databases included are: Adult,Airport,Flight,Food,Hospital,Stock,Tax,Voters  
Stock,Voters,Tax

Do you wish to run a simulation that introduces random violations in the database? [y/n]

Default value is: y

y

Please specify the number of iterations of the simulation

Default value is: 100

Please choose the measures you wish to compute:

☐ I\_D

☐ I\_MI

☐ I\_P

☐ I\_R

☐ I\_lin\_R

☐ I\_MC

This message only applies in case you wish to compute the I\_MC measure

Please make sure :

1. To build the parallel\_enum project inside the current folder
  2. that the file text\_ui created by parallel\_enum is located in /parallel\_enum/build/
- For additional instructions regarding the parallel\_enum algorithm, please visit:

[parallel\\_enum github](#)

✓ Proceed

```
---
Starting tests ['Stock', 'Voters', 'Tax'] from database inputDB.csv
with the following measurers: ['I_D', 'I_MI']; iterationsNum = 100
---
Test Stock : running 100 iterations; startTime:1601661709.264454
Computing...
```

```
---
Starting tests ['Stock', 'Voters', 'Tax'] from database inputDB.csv
with the following measurers: ['I_D', 'I_MI']; iterationsNum = 100
---
Test Stock : running 100 iterations; startTime:1601661709.264454
Computing...
Test Stock : runTime = 1601661729.3201776
Test Stock finished, preparing the results.
End of test Stock; total time = 20.420215368270874
Computation finished, outputs can be found in Data/Stock/1601661709.264454_results

Test Voters : running 100 iterations; startTime:1601661729.6866636
Computing...
Test Voters : runTime = 1601662057.3623223
Test Voters finished, preparing the results.
End of test Voters; total time = 328.203813791275
Computation finished, outputs can be found in Data/Voters/1601661729.6876616_results

Test Tax : running 100 iterations; startTime:1601662057.8924754
Computing...
Test Tax : runTime = 1601662095.2310407
Test Tax finished, preparing the results.
End of test Tax; total time = 37.607834815979004
Computation finished, outputs can be found in Data/Tax/1601662057.8934703_results
```



## אלגוריתם יצירת הפרות

לצורך חישוב המדדים נדרשנו לייצר הפרות במסד הנתונים, לצורך כך שגב צור ואני כתבנו אלגוריתם המייצר הפרות ומכניס אותן למסד הנתונים.  
(האלגוריתם נמצא בקובץ ViolationsAlgorithm.py)

הגרל שתי רשומות מהמסד באופן רנדומלי  
 $t, t'$

הגרל תלות מרשימת התלויות שסופקו לתוכנית, התלות היא מהצורה:  
 $not(t_i.att_k < operand > t_j.att_m \& t_i.att_p < operand > t_j.att_l \& \dots)$

לכל תנאי מהצורה:  $t_i.att_k < operand > t_j.att_m$   
האם התנאי כבר מתקיים עבור  $t$  ו- $t'$ ?

true

המשך לתנאי הבא

false

לכל תנאי מהצורה:  $t_i.att_k < operand > t_j.att_m$  קרא לפונקציה המתאימה ע"פ האופרנד

Operand is </>

Operand is !=

Operand is = or  
<= or >=

בחר אקראית את  
 $attribute_m$  או  $attribute_k$ ,  
בחר אקראית רשומה אחת מתוך  
 $t$  או  $t'$  וחפש עבורה ערך גדול  
או קטן בהתאמה במסד נתונים,  
אם לא נמצא תייצר ערך שונה  
באמצעות הגרלת מספר גדול או  
קטן בהתאמה לאופרנד

בחר אקראית את  $attribute_m$  או  
 $attribute_k$ , בחר אקראית  
רשומה אחת מתוך  $t$  או  $t'$  וחפש  
עבורה ערך שונה במסד נתונים,  
אם לא נמצא תייצר ערך שונה  
כתלות בסוג העמודה

בחר אקראית את  
 $attribute_m$  או  $attribute_k$   
ועדכן את הערכים עבור שתי  
הרשומות לשוויון

עדכן את שתי הרשומות שהוגרלו לערכים החדשים לאחר יצירת ההפרה

בחלק זה אסביר את האופן שבו מימשי את המדדים השונים

ראשית, אתאר שתי פונקציות מרכזיות המאפשרות את מציאת ההפרות במסד הנתונים ומשמשות לחישוב המדדים בהמשך :

#### 1. build\_dynamic\_queries -

הפונקציה מקבלת קבוצה של מחרוזות כאשר כל מחרוזת מהווה תלות כלשהי (constraint) מקובץ התלויות. כל תלות היא מהצורה :

$$not(t_i.att_k < operand > t_j.att_m \& t_i.att_p < operand > t_j.att_l \& \dots)$$

$$i, j \in \{1, 2\}, operand \in \{=, \neq, <, >\}$$

ובנוסף הפונקציה מקבלת את ה - dataframe של מסד הנתונים.

הפונקציה "תנקה" את המחרוזות של התלויות מתווים מיותרים ותבנה שתי שאליות דינאמיות.

1. unionOfAllPairs – מחרוזת של שאליות SQL שמבצעת איחוד בין כל זוגות הרשומות

המשתתפות בהפרה כלשהי זו עם זו.

השאלית תהיה מהצורה :

```
SELECT t1.rowid as t1ctid , t2.rowid as t2ctid
FROM df t1, df t2
WHERE t_i.att_k < operand > t_j.att_m AND t_i.att_p < operand > t_j.att_l AND
...AND t1.rowid != t2.rowid
UNION SELECT t1.rowid as t1ctid , t2.rowid as t2ctid FROM df t1, df t2
WHERE t_i.att_x < operand > t_j.att_y AND t_i.att_z < operand > t_j.att_t AND
...AND t1.rowid != t2.rowid
...
```

כלומר לכל תלות השאלית מחשבת את כל הרשומות המשתתפות זו עם זו בהפרה ולבסוף עושה איחוד עם כל הרשומות שחושבו עד כה עבור שאר התלויות.

2. **unionOfAllTuples** - מחרוזת של שאילתת SQL שמבצעת איחוד בין כל הרשומות המשתתפות בהפרה כלשהי. השאילתה תהיה מהצורה:

```
SELECT DISTINCT t1.* FROM (SELECT *
FROM df t1, df t2
WHERE ti.attk < operand > tj.attm AND ti.attp < operand > tj.attl AND
...AND t1.rowid! = t2.rowid
UNION SELECT * FROM df t1, df t2
WHERE ti.attx < operand > tj.atty AND ti.attz < operand > tj.attt AND
...AND t1.rowid! = t2.rowid ...) AS A
```

כלומר לכל תלות השאילתה מחשבת את כל הרשומות המשתתפות בהפרתה ולבסוף עושה איחוד עם כל הרשומות שחושבו עד כה עבור שאר התלויות, בסופו של דבר נרצה לסנן כפילויות ולהציג כל רשומה רק פעם אחת לכן נשתמש בdistinct.

שתי השאילתות הדינמיות מהוות את פלט הפונקציה.

## 2. **constraints\_check** -

הפונקציה מקבלת את dataframe של מסד הנתונים, קבוצת המחרוזות המהוות את התלויות ואת שתי השאילתות הדינאמיות שחושבו בפונקציה **build\_dynamic\_queries**. הפונקציה תריץ את השאילתות הללו באמצעות הפקודה sqldf מהספרייה pandasql ותחזיר את הפלט של שתי השאילתות במשתנים violatingPairs, violatingTuples יחד עם זמני הריצה של כל אחת מהשאילתות.

בעת, לאחר שהצגתי את שתי הפונקציות הללו אציג את המימושים של המדדים:

### 1. **first\_mesurer\_I\_D** :

הפונקציה תקבל את dataframe המכילה את כל הזוגות המשתתפות בהפרה כלשהי, במידה והdataframe ריק המשמעות היא שהמסד עקבי לכן יוחזר 0. אחרת, המסד איננו עקבי ויוחזר 1.

### 2. **Second\_mesurer\_I\_MI** :

הפונקציה תקבל את dataframe המכילה את כל הזוגות המשתתפות בהפרה כלשהי ותחזיר את מספר הרשומות בטבלה זו. (מספר הרשומות מהווה את מספר הזוגות)

### 3. `third_mesurer_I_P` :

הפונקציה תקבל את dataframe המכילה את כל הרשומות המשתתפות בהפרה כלשהי ותחזיר את מספר הרשומות בטבלה זו.

### 4. `fourth_mesurer_I_R` :

הפונקציה תקבל את dataframe (`uniquePairsDF`) המכילה את כל הזוגות המשתתפות בהפרה כלשהי ותשתמש בכלי לתכנון הלינארי `Gurobi optimizer` בכדי למצוא את המספר המינימאלי של הרשומות שיש למחוק מהמסד כדי שהמסד יהיה עקבי. הגדרת בעיית התכנון הלינארי מתבצעת באופן הבא :

*variants* : לכל רשומה  $i$  נגדיר משתנה בינארי (בעל ערך 0 או 1)  $x_i$  ייחודי .

*constraints* : לכל שתי רשומות  $i, j$  כך ש  $i, j \in \text{uniquePairsDF}$  (כלומר הרשומות משתתפות בהפרה כלשהי זו עם זו) נגדיר constraint מהצורה :

$$x_i + x_j \geq 1$$

*objective function* : פונקציית המטרה היא  $\sum_{i=1}^n x_i$  כאשר  $n$  הוא מספר הרשומות במסד הנתונים והמטרה היא להביא למינימום פונקציה זו. הפונקציה תקרא לפונקציית `optimize` של `Gurobi` ותחזיר כפלט את הערך שהתקבל יחד עם זמן הריצה של הפונקציה.

### 5. `fifth_mesurer_I_lin_R` :

הפונקציה הזו כמעט זהה לפונקציה המחשבת את המדד  $I_R$  מלבד הגדרת ה *variants* שיוגדרו להיות רציפים בין 0 ל-1 כלומר:  $0 \leq x_i \leq 1$  ולא בינאריים. באותו אופן כמו הפונקציה הקודמת, הפונקציה הזו תחזיר את הערך האופטימלי שחושב וזמן הריצה של הפונקציה.

### 6. `sixth_mesurer_I_MC` :

הפונקציה מקבלת את הכתובת המלאה לתיקיית הפרויקט ואת dataframe (`uniquePairsDF`) המכילה את כל הזוגות המשתתפות בהפרה כלשהי . הפונקציה תיצור את הגרף המשלים לגרף הקונפליקטים כלומר כל צומת מהווה רשומה כלשהי במסד הנתונים וישנה קשת בין שני צמתים רק במידה והם לא משתתפים יחד בהפרה . הגרף הנוצר יהיה בפורמט `(node degree edge)` :

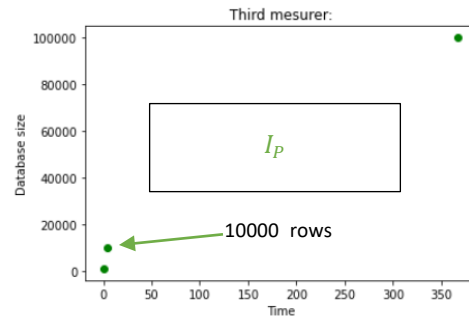
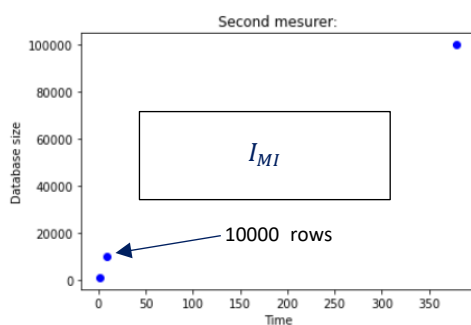
*number of nodes (i.e. number of tuples in the database)*  
*node degree (for each node in a separate row)*  
*node1 node2 (edge between the nodes)*

גרף זה יודפס אל תוך קובץ בעל סיומת `nde` ולאחר יצירתו הפונקציה תריץ את האלגוריתם `parallel_enum` ותחלץ את הפלט שהתקבל ותחזירו מהפונקציה יחד עם זמן הריצה הכולל.

תחילה, הרצתי את הקוד על מסדי נתונים המכילים כמיליון רשומות. אך נתקלתי בבעיה כאשר שמתי לב שהקוד איננו מסתיים. לאחר בדיקה לעומק של השאילתות שמוצאות את הרשומות המשתתפות בהפרה ייעלתי את ריצתן ע"י החלפת השאילתות משאילתות שמשתמשות ב *or* לשאילתות המשתמשות ב *union* ועושות למעשה איחוד בין טבלאות רבות. לאחר בדיקה ברחבי האינטרנט<sup>4</sup>, התוודעתי לכך ש *union* משפר את זמן הריצה של השאילתה ולעומת זאת שאילתה המכילה תנאים רבים עם *or* רצה יותר לאט. בהמשך, ניסיתי גם להוסיף לשאילתות אינדקסים כדי לאפשר זמן ריצה מהיר יותר (אינדקסים עבור עמודות במסד הנתונים מאפשרים מיון של הערכים של העמודות הללו וכך תתאפשר גישה מהירה יותר בעת חיפוש). הוספת האינדקסים אכן שיפרה במעט את זמן הריצה כפי שהוא מוצג ב *cost* (העלות המשוערכת של זמן הריצה של השאילתה).

Data Output	Explain	Messages	Notifications
<p>QUERY PLAN</p> <p>text</p>			
1	Hash Join (cost=99831.05..308066618.99 rows=16498067180 width=150)		
2	Hash Cond: ((t1.areacode)::text = (t2.areacode)::text)		
3	Join Filter: (((t1.state)::text <> (t2.state)::text) AND (t1.ctid <> t2.ctid))		
4	-> Seq Scan on inputdb t1 (cost=0.00..47487.02 rows=2000002 width=81)		
5	-> Hash (cost=47487.02..47487.02 rows=2000002 width=81)		
6	-> Seq Scan on inputdb t2 (cost=0.00..47487.02 rows=2000002 width=81)		

ניתן לראות בגרפים אלו את ההבדלים בזמני הריצה של הקוד על מסד הנתונים Tax : ניתן לראות כאן קפיצה משמעותית של זמן הריצה על 2 מדדים אלו בין מסד נתונים המכיל רק 10,000 שורות שרץ פחות מ-10 שניות לבין מסד נתונים המכיל 100,000 רשומות שרץ כ-400 שניות (כ-6 דקות)



לאחר הוספת המדדים הנוספים לקוד זמן הריצה הפך להיות ארוך מדי ולכן לא היה מנוס מלהקטין את מספר השורות במסד הנתונים ל-1000 שורות.

<sup>4</sup> [Using Union Instead of OR](#) , Posted on October 7, 2012 by Derek Dieter

בחלק זה אציג את תוצאות המדדים של הרצת הפרויקט על כל מסדי הנתונים הנכללים בתיקיית

הפרויקט המכילים 1000 שורות :

*Adult, Airport, Flight, Hospital, Stock, Voters, Tax, Food*

מספר ההפרות שהוכנסו הוא 50

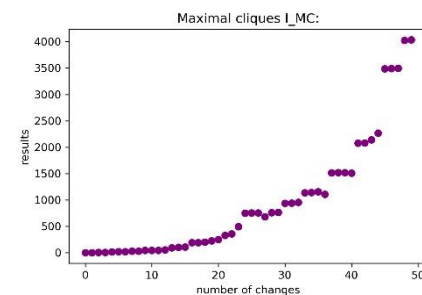
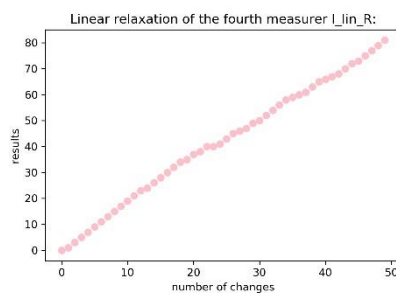
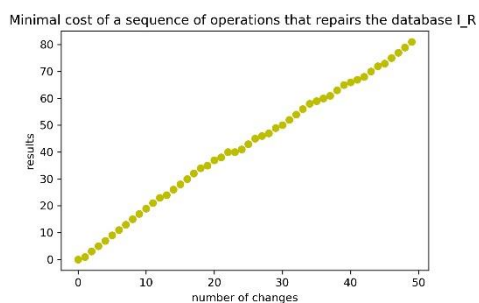
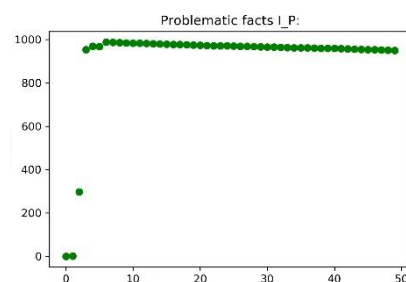
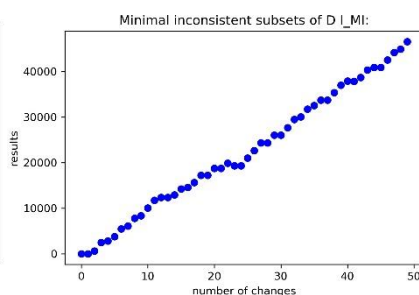
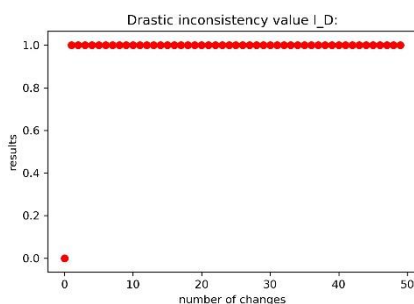
**מקרא:**  $I_D$  Drastic inconsistency

$I_{MI}$  Minimal inconsistent subsets of D  $I_P$  Problematic facts

$I_R$  Minimal tuples deletions  $I_R^{lin}$  Linear relaxation of the previous measurer

$I_{MC}$  Maximal cliques

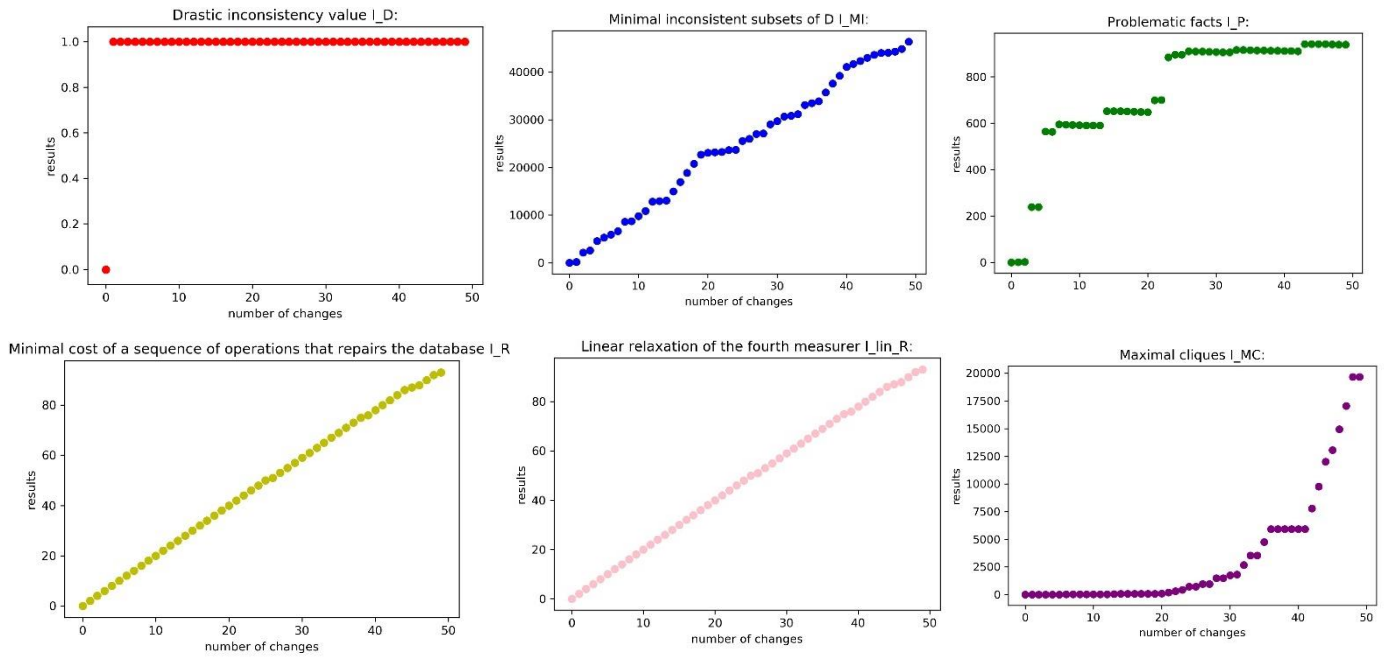
### Voters



זמני ריצה ממוצעים (בשניות) עבור כל אחד מהמדדים וסך זמן הריצה הכולל :

$I_{MI}$	$I_P$	$I_R$	$I_R^{lin}$	$I_{MC}$	Total run time
0.4404	0.6293	1.2769	1.0027	3.6970	300.9803

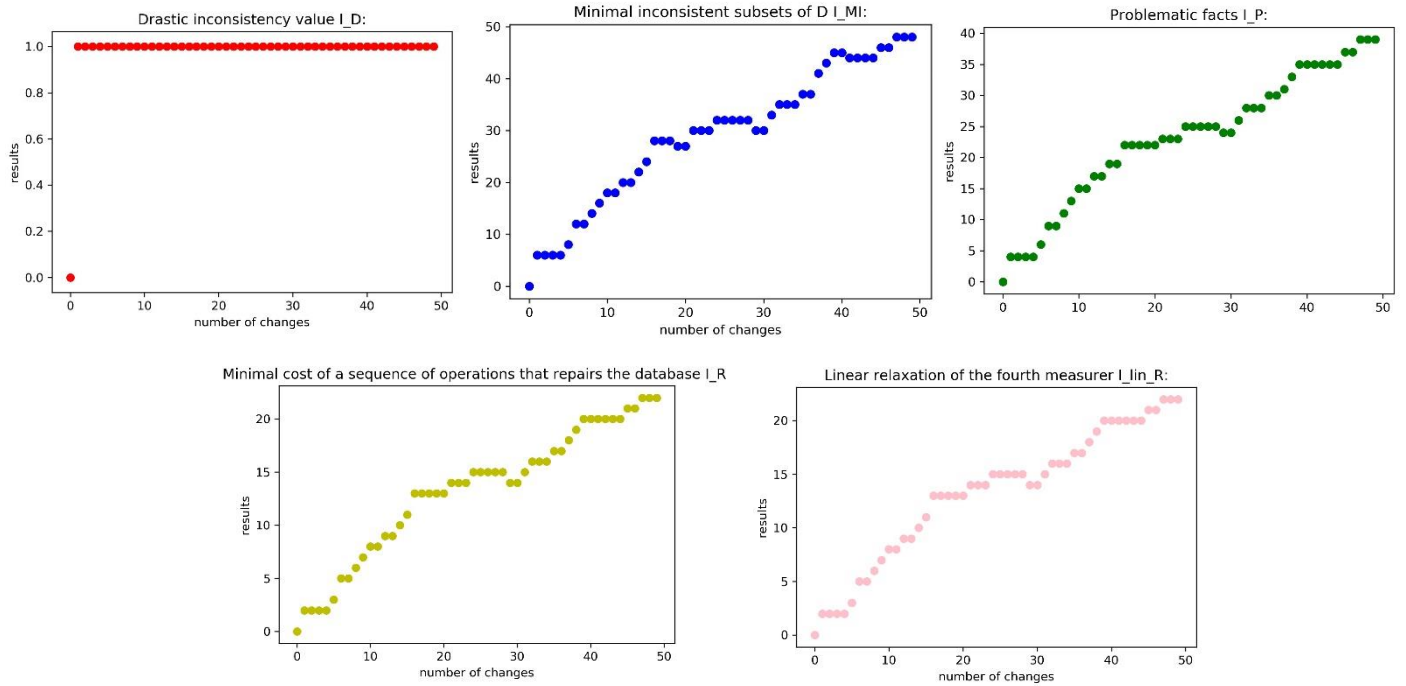
## Adult



זמני ריצה ממוצעים (בשניות) עבור כל אחד מהמדדים וסך זמן הריצה הכולל :

$I_{MI}$	$I_P$	$I_R$	$I_R^{lin}$	$I_{MC}$	<i>Total run time</i>
0.2377	0.2905	1.1616	0.8595	9.6623	589.4423

## Food



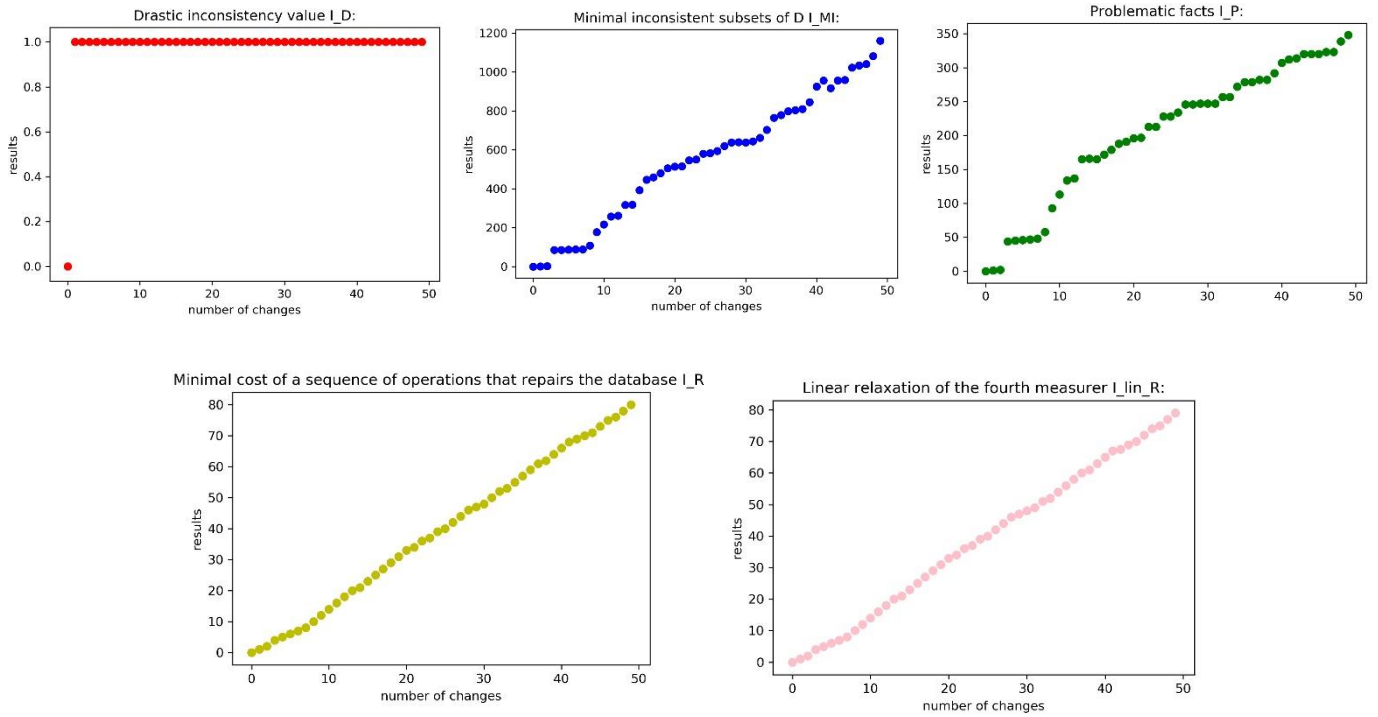
זמני ריצה ממוצעים (בשניות) עבור כל אחד מהמדדים וסך זמן הריצה הכולל :

$I_{MI}$	$I_P$	$I_R$	$I_R^{lin}$	$I_{MC}$	<i>Total run time</i>
0.1544	0.1566	0.1564	0.1555	<i>time out</i> <sup>5</sup>	22.6719

<sup>5</sup> Time out means that the running time is more than 24 hours



## Tax

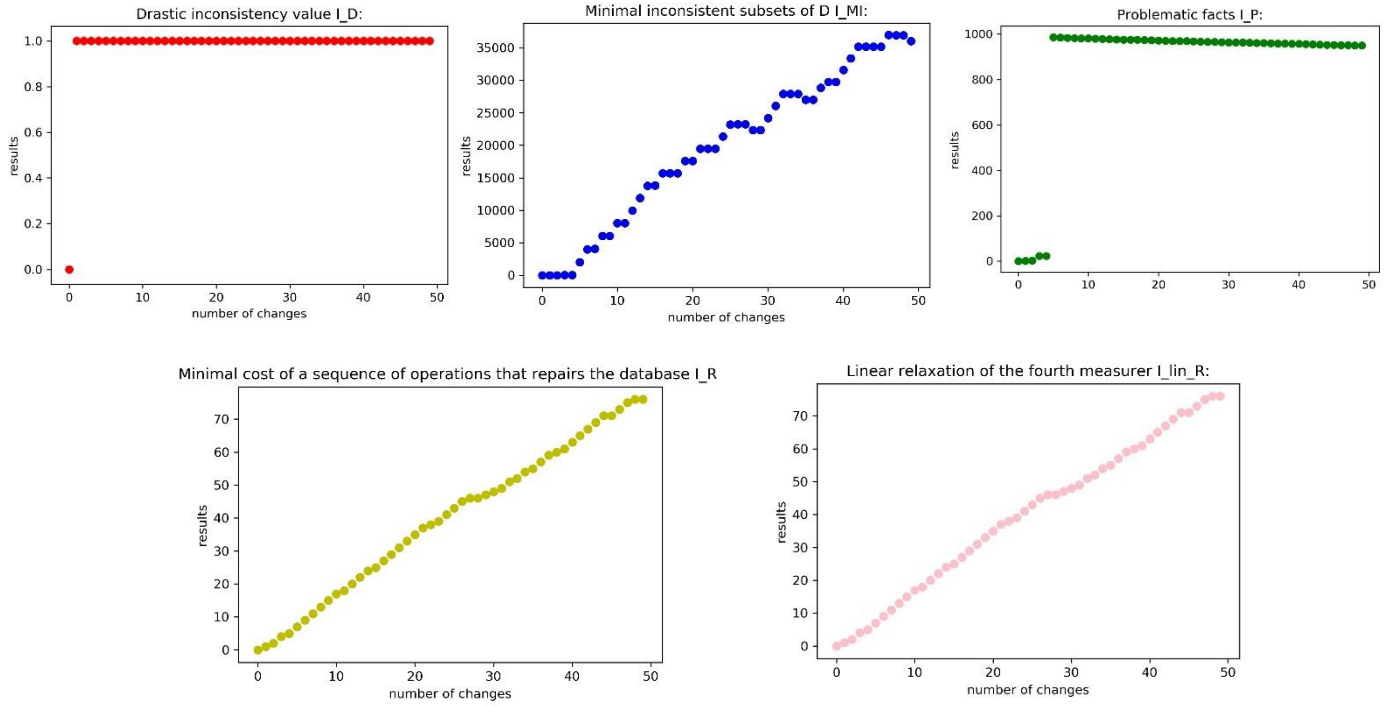


זמני ריצה ממוצעים (בשניות) עבור כל אחד מהמדדים וסך זמן הריצה הכולל :

$I_{MI}$	$I_P$	$I_R$	$I_R^{lin}$	$I_{MC}$	<b>Total run time</b>
0.0704	0.0738	0.0853	0.0841	<i>time out</i> <sup>6</sup>	15.4795

<sup>6</sup> Time out means that the running time is more than 24 hours

## Airport

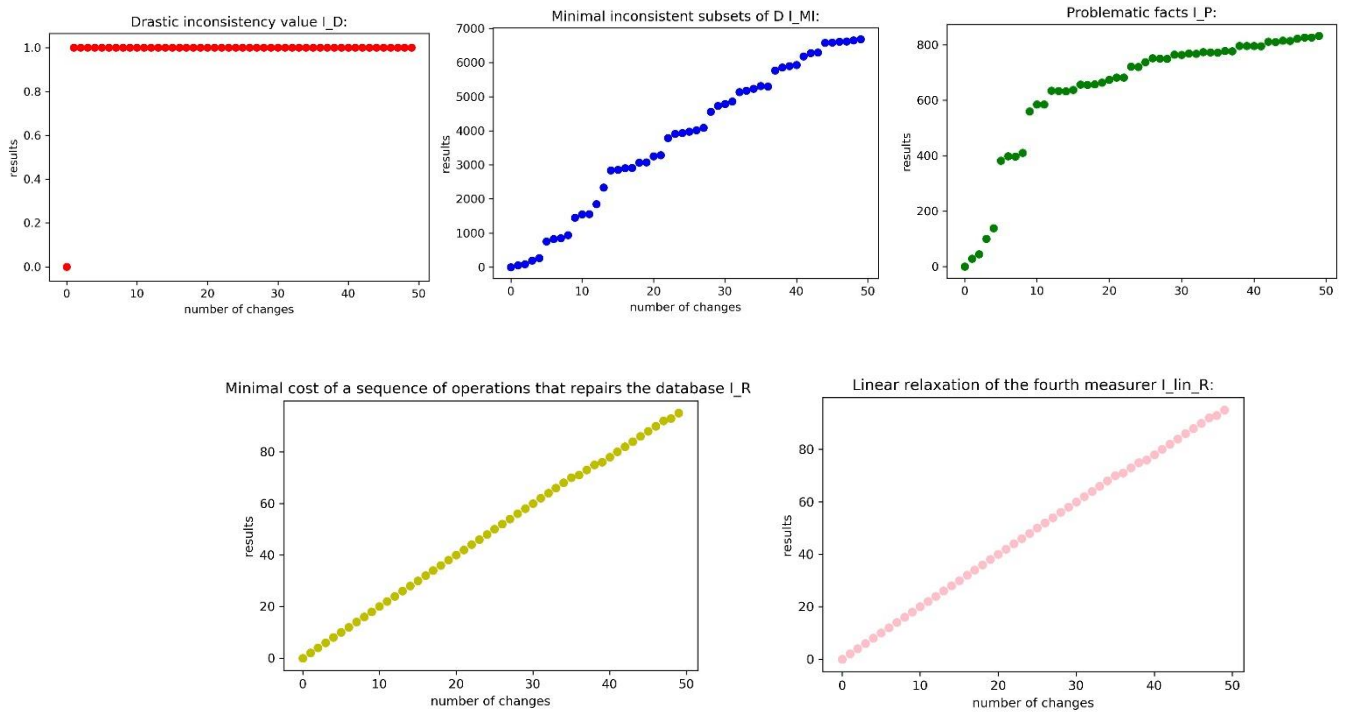


זמני ריצה ממוצעים (בשניות) עבור כל אחד מהמדדים וסך זמן הריצה הכולל :

$I_{MI}$	$I_P$	$I_R$	$I_R^{lin}$	$I_{MC}$	<i>Total run time</i>
0.2170	0.2795	0.9614	0.6876	<i>time out</i> <sup>7</sup>	94.8695

<sup>7</sup> Time out means that the running time is more than 24 hours

## Flight

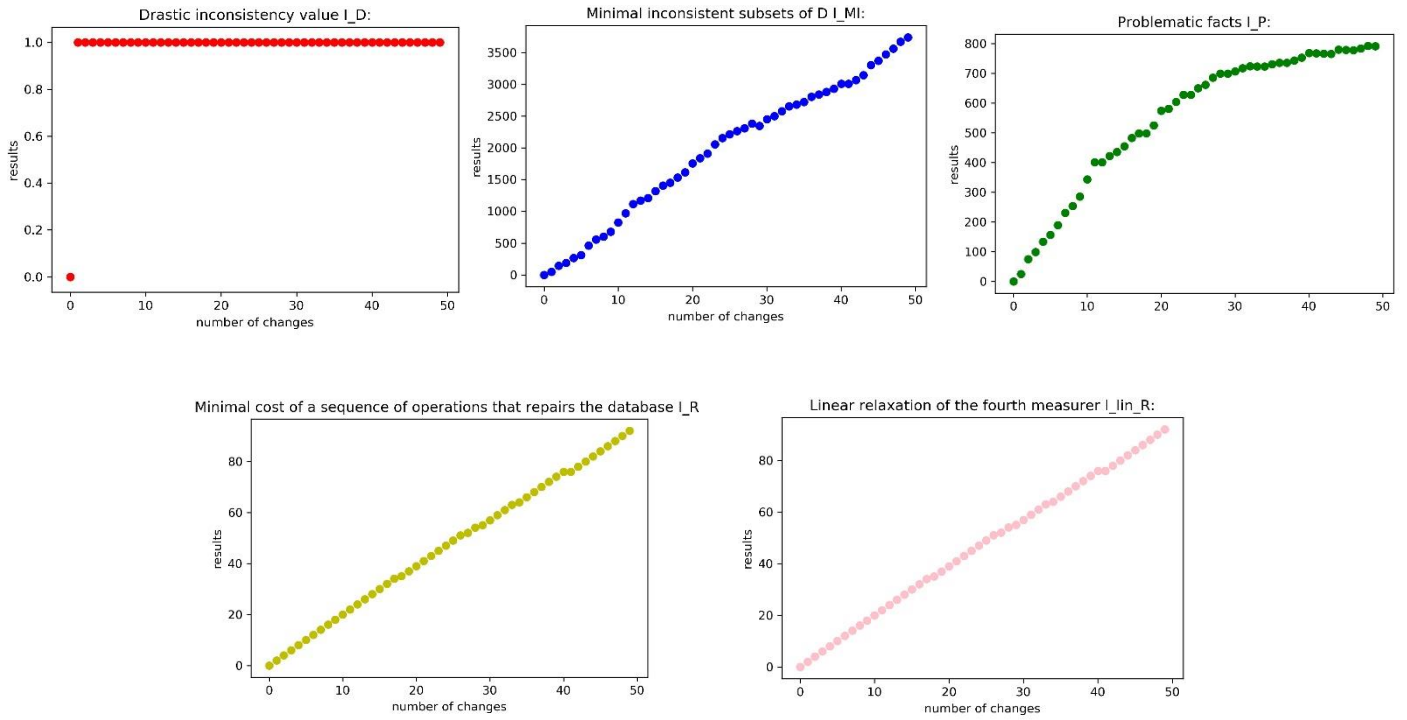


זמני ריצה ממוצעים (בשניות) עבור כל אחד מהמדדים וסך זמן הריצה הכולל :

$I_{MI}$	$I_P$	$I_R$	$I_R^{lin}$	$I_{MC}$	Total run time
0.1844	0.2094	0.3030	0.2738	time out <sup>8</sup>	37.5367

<sup>8</sup> Time out means that the running time is more than 24 hours

## Hospital

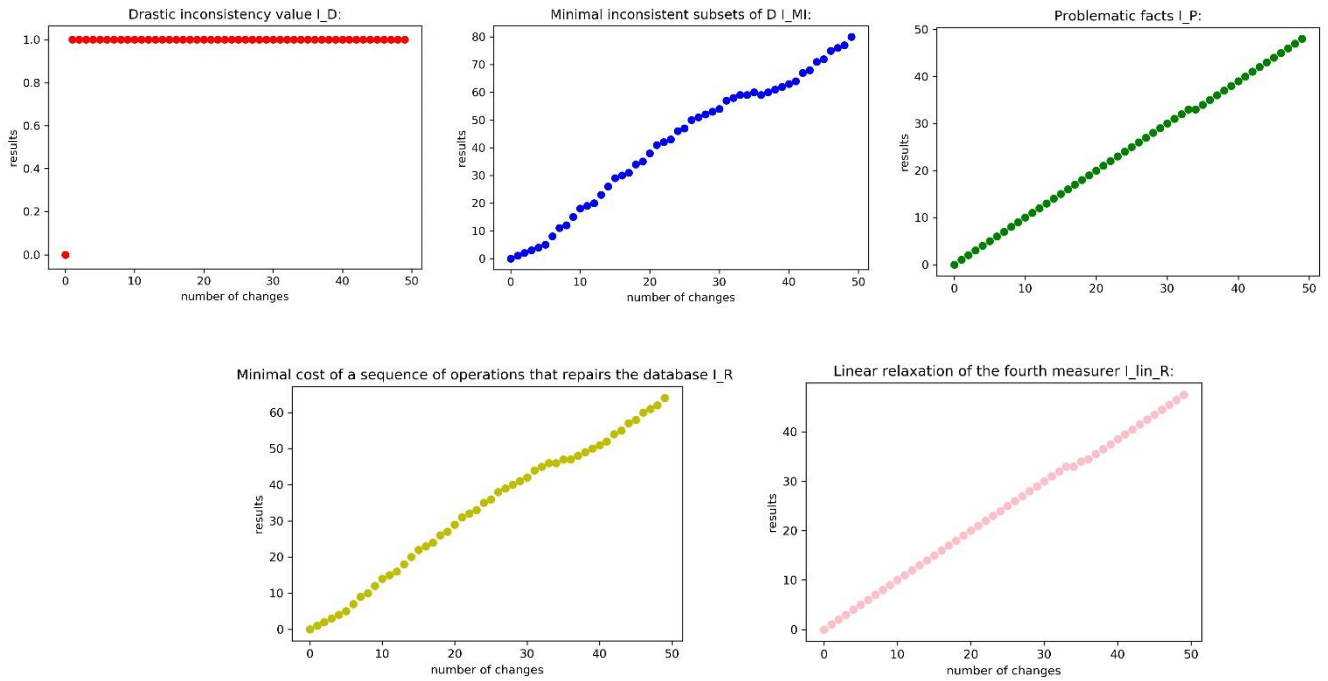


זמני ריצה ממוצעים (בשניות) עבור כל אחד מהמדדים וסך זמן הריצה הכולל :

$I_{MI}$	$I_P$	$I_R$	$I_R^{lin}$	$I_{MC}$	Total run time
0.0916	0.1051	0.1541	0.1385	time out <sup>9</sup>	22.2117

<sup>9</sup> Time out means that the running time is more than 24 hours

## Stock



זמני ריצה ממוצעים (בשניות) עבור כל אחד מהמדדים וסך זמן הריצה הכולל :

$I_{MI}$	$I_P$	$I_R$	$I_R^{lin}$	$I_{MC}$	Total run time
0.0328	0.0313	0.0344	0.0342	time out <sup>10</sup>	9.8281

<sup>10</sup> Time out means that the running time is more than 24 hours

מתוך ההתנהגות של הגרפים ראינו שלכל מדד ישנה התנהגות שונה על פני מסדים שונים. הממד הטרוויאלי ביותר  $I_D$  התנהג כמצופה והגיע לערך 1 לאחר הכנסת ההפרה הראשונה ונשאר יציב לאורך כל הריצה. הסקנו מכך שממד זה לא מספק מספיק מידע על מידת חוסר העקביות במסד הנתונים.

המדד  $I_{MI}$  נטה להיות יותר לינארי ביחס למדדים האחרים. הממד  $I_P$  לעומתו, בדרך כלל הגיע לרוויה (1000 רשומות המשתתפות בהפרה) מהר מאוד. ניתן לראות שהמדדים  $I_R, I_R^{lin}$  לאורך כל ההרצות שמרו על תבנית כמעט לינארית ויציבה. ולבסוף, הממד  $I_{MC}$  לא היה יציב במיוחד ולא הראה תבנית אחידה על פני המסדים שעליהם הוא רץ.

אך התנהגות הגרפים איננה הדבר היחיד שחקרנו, זמני הריצה של המדדים השונים הם חשובים לא פחות. זמני הריצה חשובים מאוד כדי להכריע איזה מדד מתאים לשימוש של המשתמש. למשל, אם המשתמש זקוק לנתונים אודות עקביות המסד באופן כמעט מיידי חלק מהמדדים לא יספקו לו תוצאה בפרק זמן סביר.

מדד ה-  $I_{MC}$  כפי שניתן לראות בתוצאות, היווה את צוואר הבקבוק של זמן הריצה. על חלק ממסדי הנתונים, מדד זה רץ במשך שעות ואף יותר מיממה.

אחריו הממד  $I_R$  שאמנם לא היה ארוך משמעותית מבחינת זמן הריצה שלו אך ניתן לראות בכל הניסויים כי  $I_R^{lin}$  היה בעל זמן ריצה קטן יותר בהשוואה אליו, למרות ששניהם היו בעלי תבנית לינארית. נוכל להסיק כי ה-  $relaxation$  של הממד  $I_R$  (כלומר העובדה ש"הקלנו" על הטווח האפשרי עבור המשתנים בבעיית התכנון הלינארי בממד  $I_R^{lin}$ ) היה מוצלח והביא למדד שרץ בסיבוכיות פולינומיאלית.

המסקנה שניתן להגיע אליה מתוך תוצאות הניסויים היא שהמדד  $I_R^{lin}$  הוא מדד יציב, שואף להיות כמעט תמיד לינארי וניתן לחישוב בזמן פולינומיאלי ועל כן הוא נחשב מדד איכותי ומועיל למדידת חוסר עקביות במסדי נתונים.

[inconsistency-measurer](#)

*The project's Google Drive*

[Principles of Progress Indicators for Database Repairing](#)

*Ester Livshits, Ihab F. Ilyas, Benny Kimelfeld, Sudeepa Roy, 13 Apr 2019.*

[Gurobi optimizer site](#)

*Provided me with installation info and code examples for the LP tool.*

[ipywidgets documentation](#)

*Provided me with usage examples of widgets in jupyter lab.*

[Python official site](#)

*Provided me with usage of all libraries presented in the report.*

[Pandas official site](#)

*Provided me with usage and code examples.*

[Matplotlib official site](#)

*Provided me with usage and code examples.*

[Parallel enum algorithm GitHub](#)

*The algorithm calculating maximal cliques in a graph, is used for calculating the measurer  $I_{MC}$ .*

[Using Union Instead of OR](#)

*Article regarding time optimization of SQL queries using Union instead of OR*

[Stackoverflow](#)

*Provided me with solutions and information which was not found in the previous sources.*