

תרגיל DB

ליאור שפירא וברמל גרוס

חלק א: מצורפים קבצי SQL

חלק ב:

1. הרצנו את השאילתה: `SELECT actorId FROM playsin WHERE character = 'Sheriff';`

```
public=> select actorID from playsin where character like 'Sheriff';
public=> explain analyse select actorID from playsin where character like 'Sheriff';

QUERY PLAN
-----
Seq Scan on playsin  (cost=0.00..615.15 rows=50 width=4) (actual time=0.658..3.815 rows=50 loops=1)
  Filter: (("character")::text ~~ 'Sheriff'::text)
  Rows Removed by Filter: 32602
  Planning Time: 0.288 ms
  Execution Time: 3.861 ms
(5 rows)
```

ניתן לראות כי זמן חישוב השאילתה הינו $0.288 + 3.861 = 4.149$ שכן ללא אינדקס נעבוד על כל השורות של הטבלה ונבדוק האם ה- character זהה

לאחר מכן נוסיף את האינדקס `CREATE INDEX char_name ON playsin(character)`

```
ic=> explain analyse select actorID from playsin where character like 'Sheriff';

QUERY PLAN
-----
Bitmap Heap Scan on playsin  (cost=4.67..122.36 rows=50 width=4) (actual time=0.124..0.254 rows=50 loops=1)
  Filter: (("character")::text ~~ 'Sheriff'::text)
  Heap Blocks: exact=37
  -> Bitmap Index Scan on char_name  (cost=0.00..4.66 rows=50 width=0) (actual time=0.106..0.107 rows=50 loops=1)
       Index Cond: (("character")::text = 'Sheriff'::text)
  Planning Time: 0.440 ms
  Execution Time: 0.333 ms
(7 rows)
```

נעת ניתן לראות כי זמן חישוב השאילתה הוא $0.440 + 0.333 = 0.773$ שכן הוספת אינדקס מקצרת את זמן הריצה שכן נצטרך להגיע לעלים ולחפש בעלים המתאימים בלבד.

2. חישוב שאילתות:

הנחות:

- גודל בלוק הוא 1,000 בייטים.
- בטבלה Movies יש 10,000 שורות.
- כל שורה תופסת 150 בייטים.
- התכונה movieId תופסת 8 בייט.
- התכונה duration תופסת 8 בייט.
- התכונה genre תופסת 10 בייט.
- מצביע תופס 8 בייט.
- הערכים בduration בטבלה Movies מתפלגים אחיד בטווח [1,200]
- הערכים בgenre בטבלה מחולקים ל4 קטגוריות באופן אחיד.

ראשית נחשב את מספר הבלוקים שיש ב-movies: $\left\lceil \frac{10000}{6} \right\rceil = 1667 \rightarrow \left\lceil \frac{1000}{150} \right\rceil = 6$

א.

```
SELECT DISTINCT "exists"
FROM Movies
WHERE duration > 100
```

CREATE index on movies(duration)

(1) עלות חישוב השאילתה בהנחה שאין אינדקסים על הטבלה היא **1667 I/O** שכן נצטרך במקרה הגרוע ביותר לעבור על כל השורות בטבלה.

(2) כעת נחשב מהו דרגת הפיצול האופטימלית של האינדקס-

$$8 \cdot d + 8 \cdot (d - 1) \leq 1000 \rightarrow d = 63$$

(3) עלות חישוב השאילתה באמצעות האינדקס, בהנחה שדרגת הפיצול היא זו שחושבה בסעיף הקודם.

$$\log_{\frac{d}{2}} 10,000 \rightarrow \lceil \log_{32} 10,000 \rceil = 3$$
 גובה העץ: 3

חישוב עם אינדקס: אין צורך לחפש על העלים ולכן העלות על הטיול העלים היא 1 שכן רק נצטרך לבקר בעלה אחד כלומר סה"כ העלות תהיה **4 I/O** $3 + 1 = 4$

ב.

```
SELECT avg (duration)
```

```
FROM Movies
```

```
WHERE duration > 100
```

CREATE index on movies(duration)

(1) עלות חישוב השאילתה בהנחה שאין אינדקסים על הטבלה היא **1667 I/O** בבסעיף הקודם.

(2) כעת נחשב מהו דרגת הפיצול האופטימלית של האינדקס-

$$8 \cdot d + 8 \cdot (d - 1) \leq 1000 \rightarrow d = 63$$

(3) עלות חישוב השאילתה באמצעות האינדקס, בהנחה שדרגת הפיצול היא זו שחושבה בסעיף הקודם.

$$\log_{\frac{d}{2}} 10,000 \rightarrow \lceil \log_{32} 10,000 \rceil = 3$$
 גובה העץ: 3

$$\left\lceil \frac{200-100}{200} \right\rceil = \frac{1}{2} \rightarrow \left\lceil \frac{\frac{1}{2} \cdot 10000}{32-1} \right\rceil = 162$$
 חישוב עם אינדקס: 162

$$162 + 3 = 165$$
 אזי סה"כ העלות תהיה **165 I/O**

ג.

```
SELECT name
```

```
FROM Movies
```

```
WHERE movieid=200
```

CREATE index on movies(movieid)

(1) עלות חישוב השאילתה בהנחה שאין אינדקסים על הטבלה היא **1667 I/O**.

(2) כעת נחשב מהו דרגת הפיצול האופטימלית של האינדקס-

$$8 \cdot d + 8 \cdot (d - 1) \leq 1000 \rightarrow d = 63$$

(3) עלות חישוב השאילתה באמצעות האינדקס, בהנחה שדרגת הפיצול היא זו שחושבה בסעיף הקודם.

$$\log_{\frac{d}{2}} 10,000 \rightarrow \lceil \log_{32} 10,000 \rceil = 3$$
 גובה העץ:

מכיוון ש-*movieID* הוא מפתח הוא יופיע בעלה אחד בלבד וכן ניגש לטבלה פעם אחת אזי סה"כ העלות

$$3 + 1 + 1 = 5 \text{ I/O}$$

ד. נתונה השאילתה

```
SELECT avg(duration)
FROM Movies
WHERE genre = 'Drama'
create index on movies(genre)
```

(1) עלות חישוב השאילתה בהנחה שאין אינדקסים על הטבלה היא **1667 I/O**.

(2) בעת נחשב מהו דרגת הפיצול האופטימלית של האינדקס-

$$10 \cdot d + 8 \cdot (d - 1) \leq 1000 \rightarrow d = 56$$

(3) עלות חישוב השאילתה באמצעות האינדקס, בהנחה שדרגת הפיצול היא זו שחושבה בסעיף הקודם.

$$\log_{\frac{d}{2}} 10,000 \rightarrow \lceil \log_{28} 10,000 \rceil = 3$$
 גובה העץ:

הערכים מתפלגים אחיד על 4 קטגוריות ולכן יש סה"כ $\frac{10000}{4} = 2500$ שורות עם ג'אנר

דרמה. אזי עבור חישוב עם אינדקס נקבל כי יש $93 = \frac{2500}{\frac{d}{2}-1}$ בלוקים.

כעת הגישה לטבלה היא גישה למינימום מבין מספר השורות המתאימות למול מספר

$$\min(1667, 2500) = 1667$$
 ולכן

$$1667 + 93 + 3 = 1763 \text{ I/O}$$
 תהיה העלות

ה.

```
SELECT avg(duration)
FROM Movies
WHERE genre = 'Drama'
create index on movies(genre,duration)
```

(1) עלות חישוב השאילתה בהנחה שאין אינדקסים על הטבלה היא **1667 I/O**.

(2) בעת נחשב מהו דרגת הפיצול האופטימלית של האינדקס-

$$10 \cdot (d - 1) + 8 \cdot (d - 1) + 8d \leq 1000 \rightarrow \mathbf{d = 39}$$

(3) עלות חישוב השאילתה באמצעות האינדקס, בהנחה שדרגת הפיצול היא זו שחושבה בסעיף הקודם.

$$\log_{\frac{d}{2}} 10,000 \rightarrow \lceil \log_{20} 10,000 \rceil = 4$$

גובה העץ:

$$\text{יש סה"כ } \frac{10000}{4} = 2500 \text{ שורות המתאימות לג'אנר דרמה.}$$

$$\text{עבור חישוב עם אינדקס: } \left\lceil \frac{2500}{\frac{d}{2}-1} \right\rceil = 132, \text{ בנוסף הפעם יש אינדקס ולכן לא יעלה לנו}$$

$$132 + 3 = \mathbf{136 \text{ I/O}}$$

לגשת לטבלה. אזי סה"כ העלות תהיה