

Analyzes

פרטי האינדקס:

כאשר מריצים את התוכנית מתווצר ארבעה קבצים בינאריים בפורמט הבא:

- קובץ index.txt המכיל בכל שורה ארבעה נתונים (מספר המזהה של כל review הוא מספר השורה בקובץ):
 - productId: האורך שלו קבוע (10 תווים) והוא נשמר בתחילת הקובץ.
 - Score: בגודל קבוע (תו בודד) הוא מספר של Byte1.
 - Helpfulness: נשמר בצורה של מחרוזת בתוך המחרוזת יש "\" כדי שנוכל להפריד בין המונה והמכנה וברגע שקוראים אחד מהם אנו ממירים אותו ל-integer.
- קובץ TotalFreq.txt שמכיל בשורה הראשונה מספר המילים בכל ה-Reviews ושאר השורות אורך של כל Review (מספר המזהה של כל review הוא מספר השורה בקובץ).
- קובץ words.txt מכיל את כל המילים שנמצאים בכל ה-reviews ללא חזרה כך שכל מילה נמצאת בשורה בודדת.
- קובץ data.txt מכיל את הנתונים של המילים כך שכל שתי שורות קשורות למילה בקובץ ה-words לפי הסדר שלהם כך שהשורה הראשונה שומרת מערך באיזה reviews מכילים את המילה והשורה השניה מתארת מערך של מספר הופעתה ב-reviews.

```
line1: <id1><score1><helpfulness1>
line2: <id2><score2><helpfulness2>
line3: <id3><score3><helpfulness3>
line4: <id4><score4><helpfulness4>
```

index.txt

```
line1: the words number in all reviews
line2: the words number in product1
line3: the words number in product2
line4: the words number in product3
```

TotalFreq.txt

```
line1: word1
line2: word2
line3: word3
line4: word4
```

words.txt

```
line1: list of the reviews that have the word1
line2: list of the numbers that the word1 appears in the reviews above
```

data.txt

קריאת הנתונים מהזיכרון/הדיסק

כאשר מריצים את התוכנית שום דבר לא נשמר בזיכרון, וברגע שצריכים לקרוא נתונים מסוימים ניגשים לקובץ המתאים בדיסק, קוראים את מה שצריכים ואז סוגרים את הקובץ שמסיימים ממנו.

הגודל הצפוי של האינדקס

אפשר לחשב גודל האינדקס דרך הנוסחאות הבאות:

- קובץ ה-index.txt: שומרים productid, score, helpfulness בגדלים קבועים שהם 5-3, 1, 10 בהתאם לכן הנוסחא היא 16 כפול מספר ה-review.
- קובץ ה-TotalFreq.txt בכל שורה שומרים מספר מסוג integer שגודלו 3 byte במומצע אז הנוסחא היא 3 כפול מספר ה-reviews.
- קובץ ה-words: ממוצע האותיות במילה הוא 7 אותיות, גודל השורה היא 2 bytes לכן הנוסחא: $(\text{מספר ה-reviews} * 2 \text{ bytes}) + (7 * \text{מספר המילים ללא חזרה})$
- קובץ ה-data.txt: עבור כל מילה שומרים שני מערכים כך שכל מערך מכיל מספרים של ה-reviews שהופיעה בהם המילה, הנוסחא היא: $2 * \text{מספר המילים} * (2 \text{ bytes} + 3 \text{ bytes} * \text{מספר ההופעה למילה})$