

שאלות Data Integrity

1.

כדי לתת המלצות סרטים משתמשים ברעיון שסרטים שיש בהם את אותם התפקידים הם המלצות טובות.
הכלל ממומש בשאילתה הבאה:

```
select
f.movie_id as fid
s.movie_id as sid ,
count(distinct f.role) as roles ,
from
imdb_ajs.roles as f
join
imdb_ajs.roles as s
on
f.role = s.role
group by
fid
sid ,
having
count(distinct f.role) >= 3
```

הוטל עליכם לשפר את ביצועי השאילתה. עליכם ליצור ואריציה עם precision עדיף ווריאציה עם recall עדיף.
האילוץ הוא שהשאילתה נשארה כמו שהיא ומותר לכם לשנות רק את טבלת roles.

כיוונים אפשריים הם התייחסות ל:

1. התפקיד הנפוץ ביותר "" (שמשמש כשהתפקיד לא ידוע).
2. התפקיד Himself
3. התפקיד extra (ניצב)
4. תפקידי דמויי ניצב (נהג, פקיד קבלה)
5. תפקידים בעלי מופעים רבים בסרט (חיילים בסרט מלחמה)
6. תפקיד ראשי (והזיהוי שלו)
7. תפקידים ראשיים חוזרים (ג'יימס בונד, הארי פוטר)
8. תפקיד סופרמן/קלארק קנט
9. שמות דומים לאותו התפקיד (רוצח\מתנקש)
10. תפקידים דומים סמנטית (ביולוג/ארכיאולוג על תקן מדען מטורף)

1. הציעו 3 שיטות לכל ואריאציה.
2. ממשו ב sql או הסבירו מדוע sql אינה כלי מתאים למימוש הרעיון שלכם.

2.

בבנית ה gold standrd התבקשת לתת 200 זוגות של סרטים שמחציתם המלצות טובות ומחציתם המלצות בינוניות (בעלות קשר סביר אבל לא מספיק חזק להיות המלצה).

לקורס התגנב איש פלאי שמטרתו להזין המלצות לא מתאימות. המניעים להמלצות הלא טובות עשויות להיות חוסר הבנה, חוסר זמן, רצון לחמוד לצון, זדון וכדומה. חשוב מאד לאתר המלצות אלו כי המודלים שלכם ימדדו עליהן.

1. ציינו מוטיבציה וסוג המלצות לא מתאימות.
 2. הסבירו כיצד יראו הנתונים.
 3. הסבירו איך הנתונים ישפיעו על מדידת מערכת ההמלצות.
 4. הציעו דרך לזהות את הנתונים. ממשו ב sql או הסבירו מדוע sql אינה כלי מתאים למימוש הרעיון שלכם.
 5. הציעו דרך לתקן או להסיר את הנתונים. ממשו ב sql מדוע sql אינה כלי מתאים למימוש הרעיון שלכם.
- ענו על הסעיפים למעלה עבור שלושה מקרים.

3.

שמות השחקנים היא איזור מועד לשגיאות.

1. נניח שסרקתם את ויקיפדיה ובניתם רשימת שמות שחקנים. אלו בעיות הרשימה יכולה לעזור לזהות וכיצד? מה ההשלכות של בעיות כאלו על מודלי המלצה?
2. נניח שסרקתם את ויקיפדיה ובניתם רשימת שמות פוליטיקאים. אלו בעיות הרשימה יכולה לעזור לזהות וכיצד? מה ההשלכות של בעיות כאלו על מודלי המלצה?
3. רונלד רייגן היה שחקן ופוליטיקאי. נניח שהוא ואחרים נמצאים בשתי הרשימות. מה תהיה ההשפעה?
4. חישובו על שלושה מקרים שיכולים להוביל לשמות לא תקינים.
 - a. ציינו את הבעיה
 - b. הסבירו כיצד יראו הנתונים.
 - c. הסבירו איך הנתונים ישפיעו על מדידת מערכת ההמלצות.
 - d. הציעו דרך לזהות את הנתונים. ממשו ב sql או הסבירו מדוע sql אינה כלי מתאים למימוש הרעיון שלכם.
 - e. הציעו דרך לתקן או להסיר את הנתונים. ממשו ב sql מדוע sql אינה כלי מתאים למימוש הרעיון שלכם.
5. אנחנו מקבלים כמובן מאליו שהשמות הם באנגלית. האם זו בעיה בנתונים? כיצד הייתם תומכים בריבוי שפות?

המושג null records מציין את ה true negatives - רשומות שהן שליליות ומסווגות כשליליות. בהמלצות שלנו הרוב המוחלט של הזוגות הן רשומות שליליות ויסווגו כך משום שאין כל אינדיקציה לקשר בניהן. רשומות כאלו עלולות לפגוע בהערכת הביצועים, בוודאי כשהוזנו במקום רשומות מועילות יותר (המלצות בינוניות).

1. מה תהיה ההשפעה של null records על המדדים accuracy, precision, recall
2. ניתן לאתר null records בעזרת מודלים בעלי recall מאד גבוה. רשומה של המלצה שלילית שלא מסווגת כחיובית על ידי אף אחד מהמודל היא בהסתברות גבוהה null record
3. בנו מודל בעל recall גבוה על בסיס
 - a. במאי משותף
 - b. שחקן משותף
 - c. ז'אנר משותף
 - d. תפקיד משותף
 - e. שנה קרובה
 - f. על בסיס collaborative filtering "אנשים שאהבו את א', אהבו את ב"
 - g. בונס: למה מודל של: "אנשים שדירגו את א', דירגו גם את ב" יכול להועיל גם? ממשו
4. אחדו את כל המודלים לטבלה אחת
5. מיצאו את ההמלצות שהן null records
6. בידקו האם יש ממליצים הנוטים במיוחד להמליץ על null records