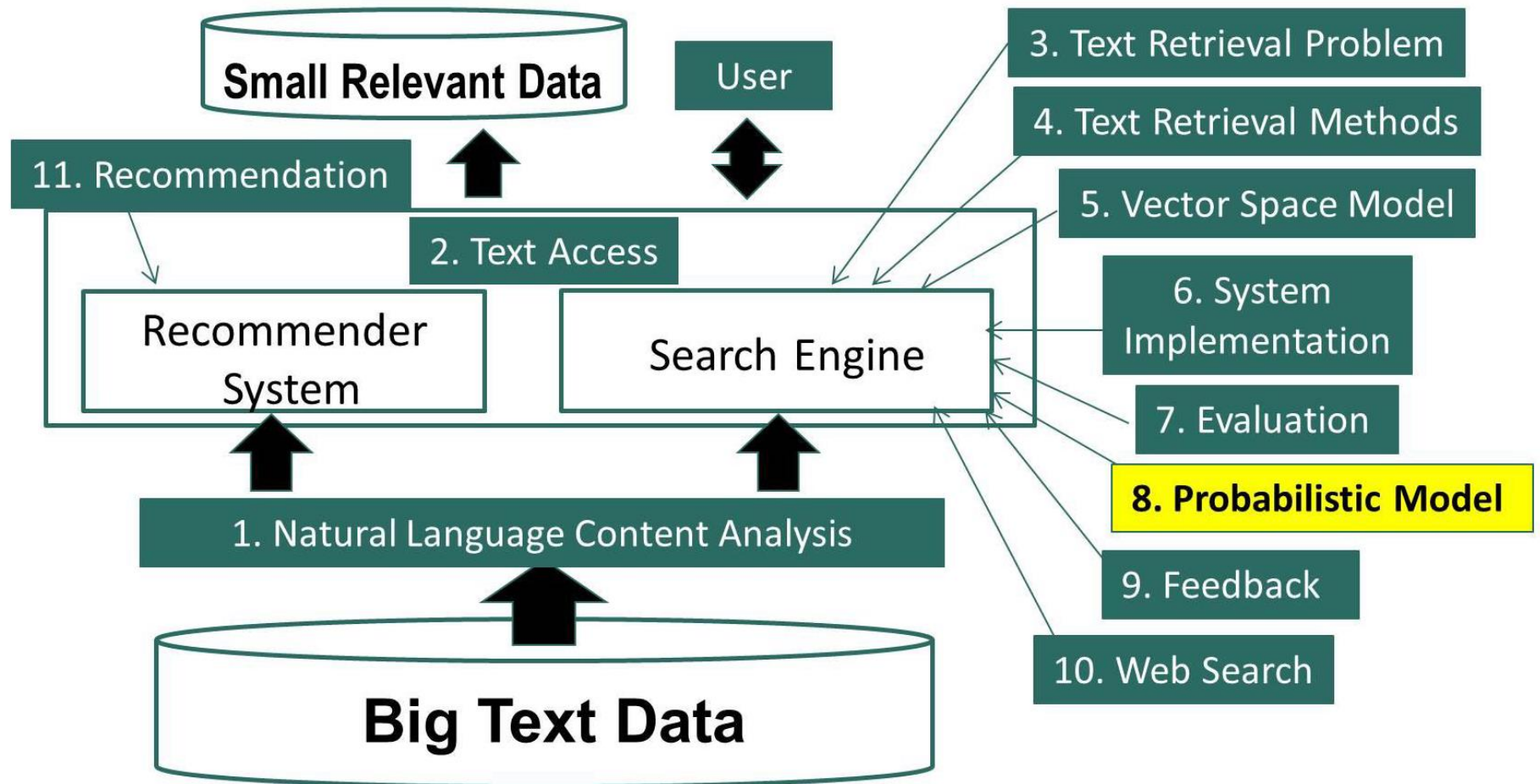


מודל אחזור הסתברותי



# שיטות אחזור מידע שונות

כיצד כן/לאן 1/0

- מודל הסתברותי:  $f(d,q) = p(R=1 | d,q)$ ,  $R \in \{0,1\}$

– מעריכים את ההסתברות של הרלוונטיות של המסמך לשאילתה.

– BM25 הוא מודל הסתברותי קלאסי

– מודל שפתי, ספציפית, סבירות שאילתות (Query Likelihood) במודל זה אנו נתמקד

- ב- QL אנו נעריך את ההסתברות

$$p(R=1 | d,q) \approx p(q | d, R=1)$$

אם משתמש מעוניין במסמך  $d$  מה הסיכוי שהמשתמש יזין את השאילתה  $q$  (על מנת לאחזר את המסמך  $d$ )?

כשאנחנו כותבים שאלה אלו  
יש להם "רעיון" מסמך קטן  
זה כותבים את המילה שזוהי  
אנחנו יכולים  
מילים  
במילים  
המילים/מילים

הסיכוי  
המילה המילה  
אם  
המילה

# מודל אחזור הסתברותי (הרעיון המרכזי)

לא אף פעם  
חשבים על זה.

המרחב של המידע הנ"ל

| Query<br>(q) | Doc<br>(d) | Relevance<br>(R) |
|--------------|------------|------------------|
| q1           | d1         | 1                |
| q1           | d2         | 1                |
| q1           | d3         | 0                |
| q1           | d4         | 0                |
| q1           | d5         | 1                |
| ...          | ...        | ...              |
| q1           | d1         | 0                |
| q1           | d2         | 1                |
| q1           | d3         | 0                |
| q2           | d3         | 1                |
| q3           | d1         | 1                |
| q4           | d2         | 1                |
| q4           | d3         | 0                |

$$f(d,q) = p(R=1 | d,q) = \frac{\text{count}(q, d, R = 1)}{\text{count}(q, d)}$$

$$P(R=1 | q1, d1) = ? \text{ 1/2}$$

$$P(R=1 | q1, d2) = ? \text{ 2/2}$$

$$P(R=1 | q1, d3) = ? \text{ 0/2}$$

איך נתמודד עם מסמכים/שאלות שלא

נראו/הוזנו? (בעייתי, לא מתמודד עם בעיה זו, עדיין  
מציג את הרעיון הבסיסי) נצטרך להעריך בדרך כלשהי

בזמן שיש לנו את כל המידע הנ"ל, אפשר  
לחשב את ההסתברות לרלוונטיות.

# מודל אחזור Query Likelihood

| Query (q) | Doc (d) | Relevance (R) |
|-----------|---------|---------------|
| q1        | d1      | 1             |
| q1        | d2      | <u>1</u>      |
| q1        | d3      | 0             |
| q1        | d4      | 0             |
| q1        | d5      | 1             |
| ...       | ...     | ...           |
| q1        | d1      | <b>0</b>      |
| q1        | d2      | <u>1</u>      |
| q1        | d3      | 0             |
| q2        | d3      | 1             |
| q3        | d1      | 1             |
| q4        | d2      | 1             |
| q4        | d3      | 0             |

המשתמש מעוניין במסמך

$$f(d,q) = p(R=1 | d,q) \approx p(\mathbf{q} | \mathbf{d}, R=1)$$

מה ההסתברות שהמשתמש  
יזין את השאילתה q

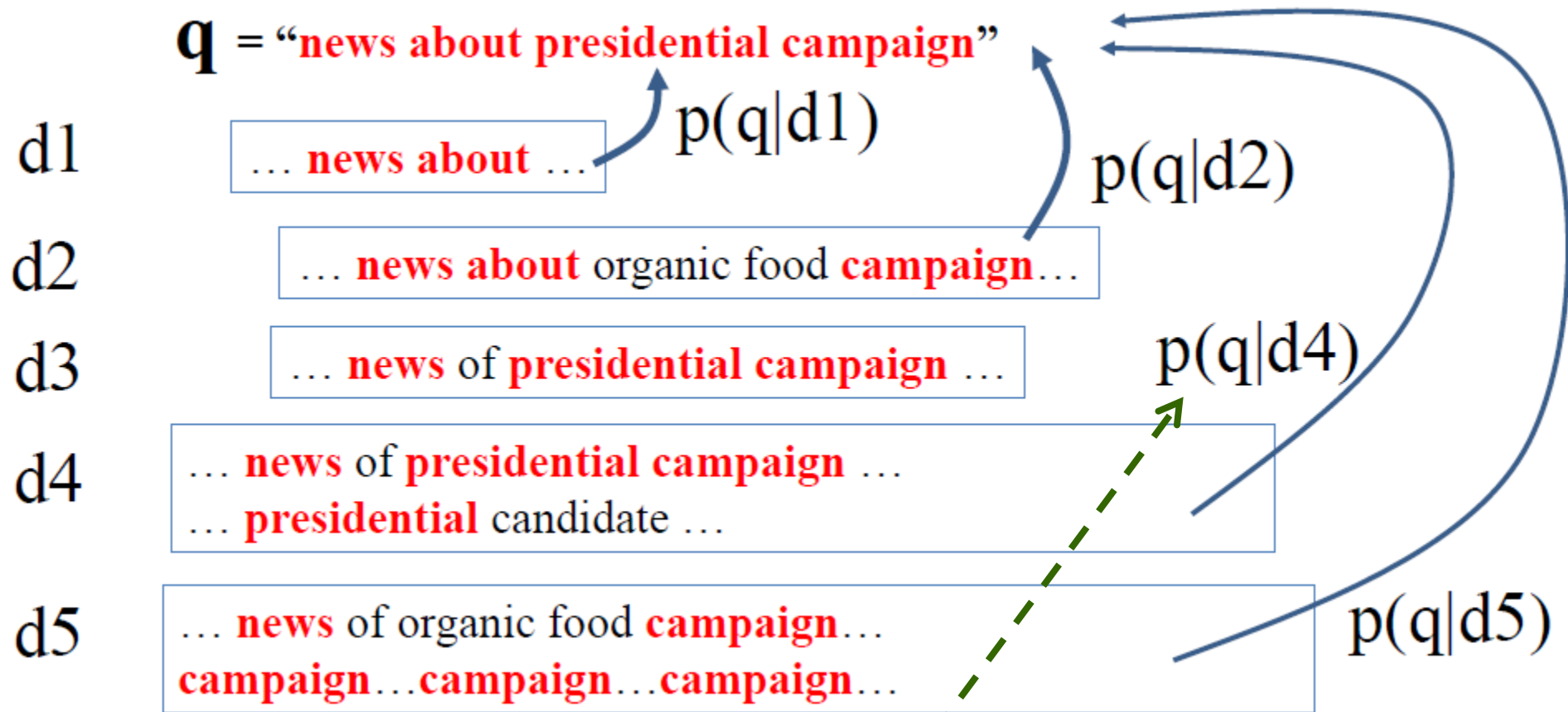
הנחה: משתמש בונה שאילתה על  
סמך "מסמך רלוונטי בעיני-רוחו"

מהו ככל-הנראה "מסמך רלוונטי בעיני-רוחו"

**q** = "news about presidential campaign"

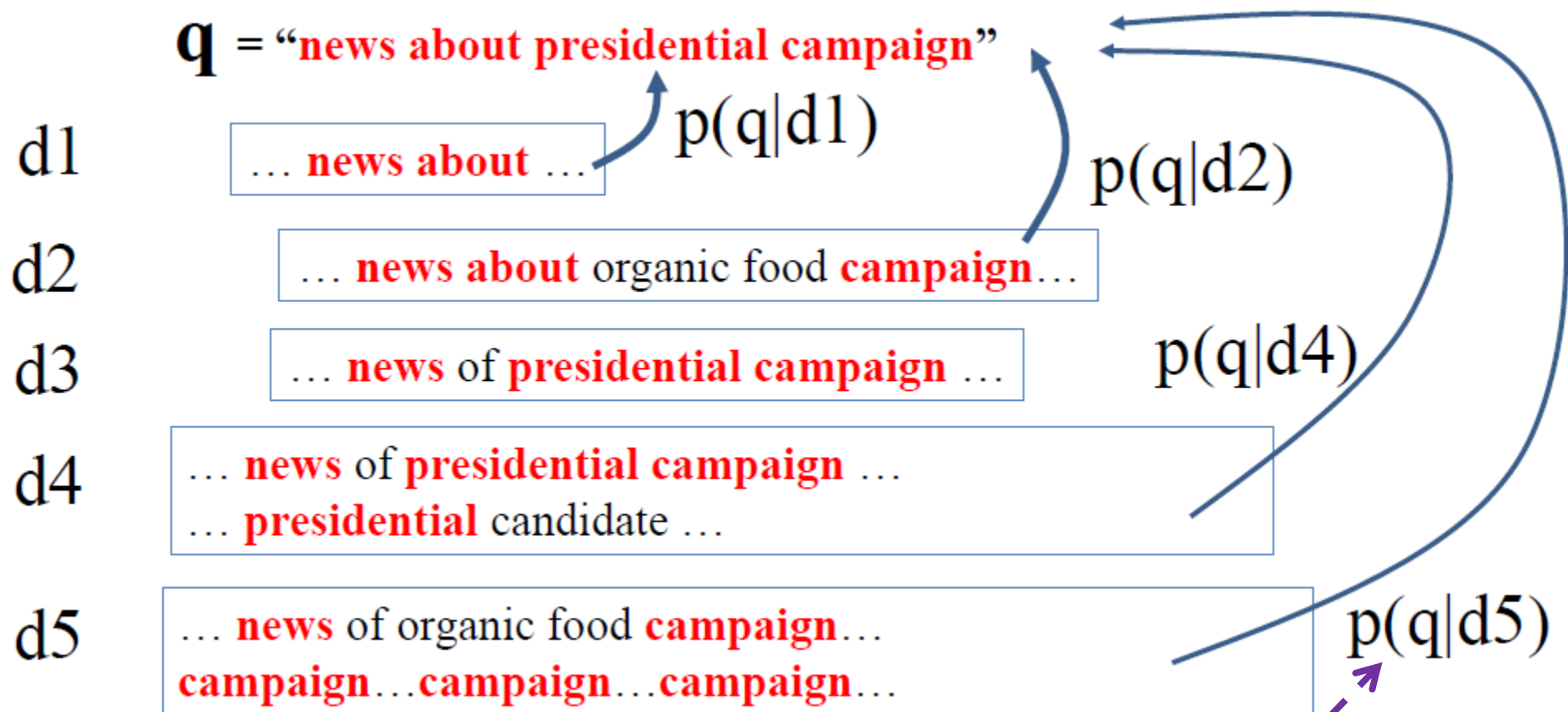
האם זה נכון?

# מהו ככל-הנראה "מסמך רלוונטי בעיני-רוחו"



נניח שהמשתמש מעוניין  
במסמך  $d4$  מה הסיכוי  
שהשאלתה תהיה  $q$ ?

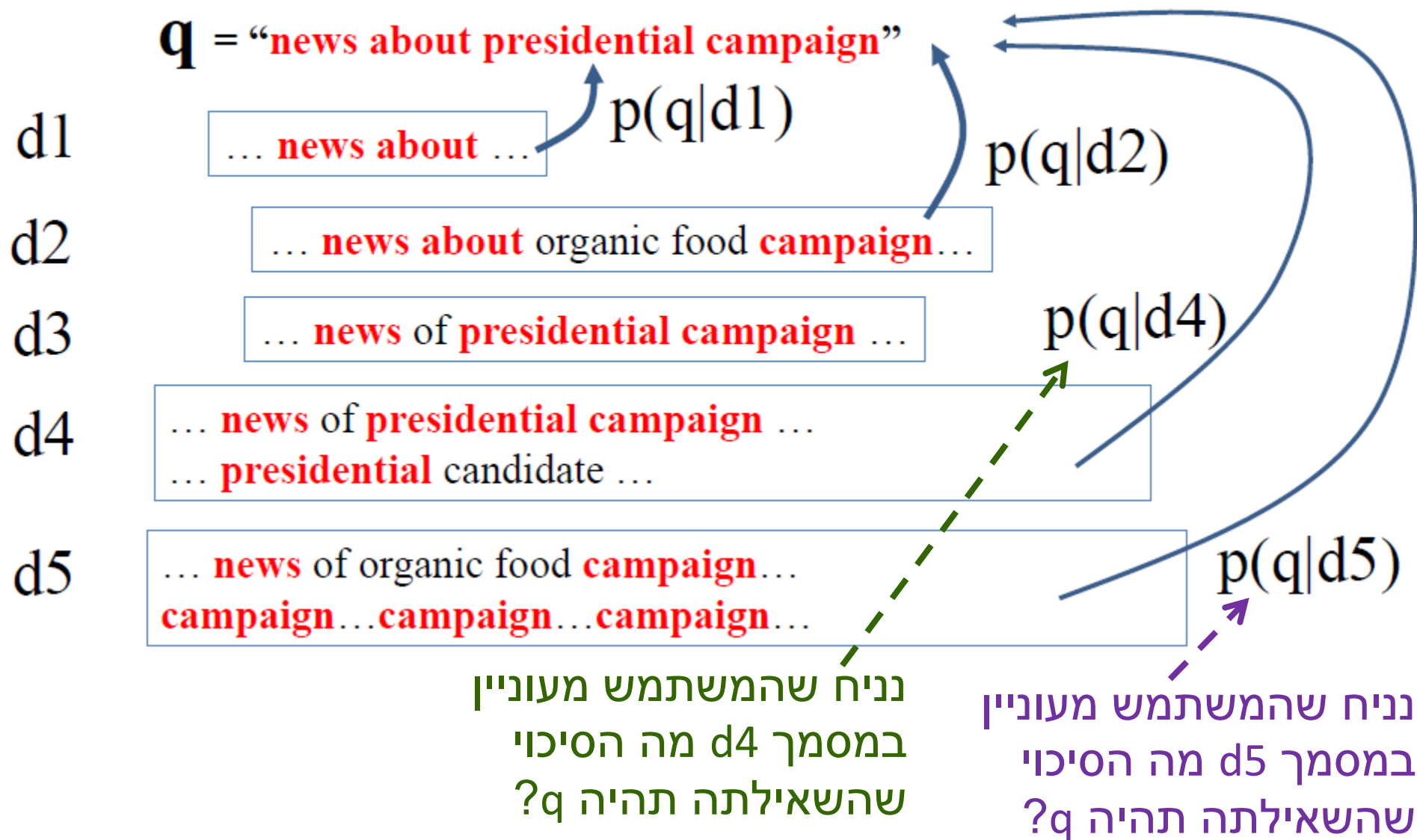
# מהו ככל-הנראה "מסמך רלוונטי בעיני-רוחו"



נניח שהשתמש מעוניין  
במסמך d5 מה הסיכוי  
שהשאלתה תהיה q?



# מהו ככל-הנראה "מסמך רלוונטי בעיני-רוחו"



# סיכום ביניים

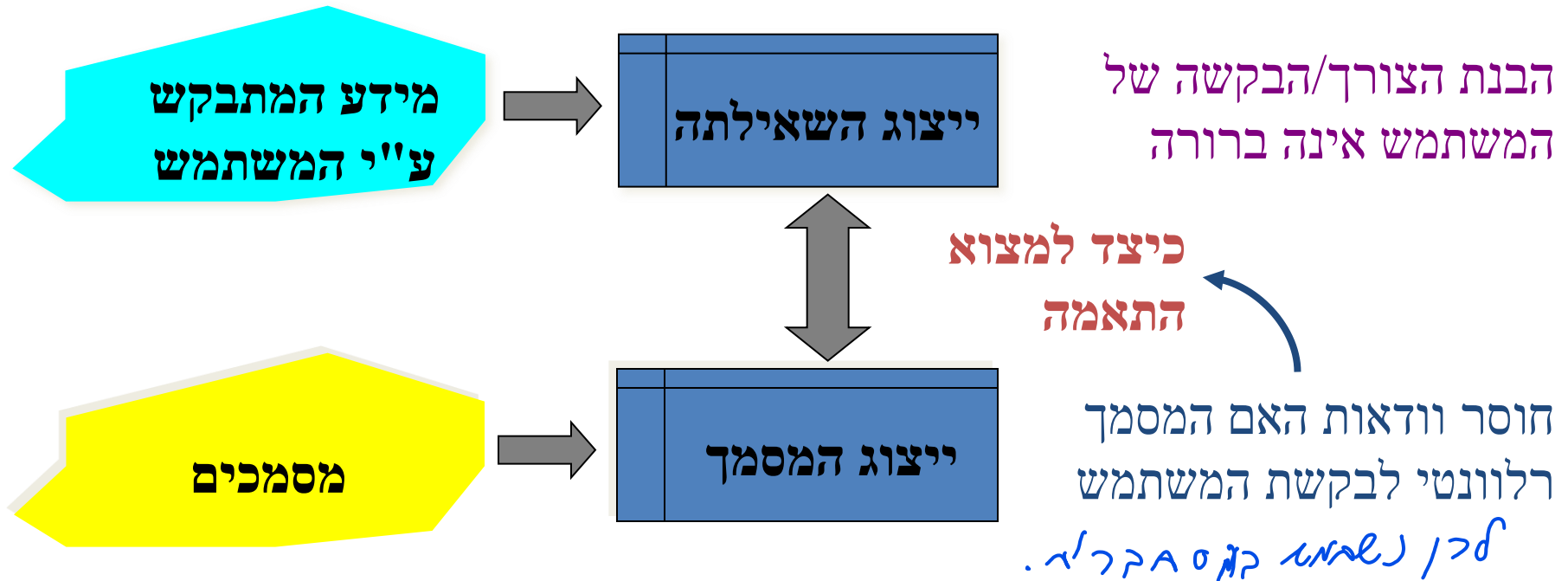
- הרעיון הכללי של המודל ההסתברותי, אנו מניחים שיש ערך אקראי  $R$  "וסביבו" בונים את פונקציית הדירוג.
- בגלל בעיות מסויימות אנו משתמשים בפונקציית קירוב, ההסתברות של  $q$ , ספציפית, בהינתן המסמך (כך לכל אחד מהמסמכים הרלוונטיים)
- $\text{Relevance}(q,d) = p(R=1 | q,d) \rightarrow p(q | d,R=1)$
- פונקציית הדירוג, **Query likelihood**:  $f(q,d)=p(q | d)$
- אבל, כיצד נחשב  $p(q | d)$ ? כיצד מחשבים הסתברות של טקסט בכלל?, לצורך זה יש את המודל השפתי, Language Model

$p(q=\text{"presidential campaign"} | d=\text{"...news of presidential campaign ... presidential candidate ..."})$

מה הסיכוי שהוא יזין את השאילתה הזו

ספציפית, אם המשתמש מעוניין במסמך זה

# מדוע להשתמש במודל הסתברותי בשביל IR



במערכות IR מסורתיות, מנסים להתאים בין כל מסמך לשאילתה כאשר מרחב מונחי האינדקס במרחב אינו ברור לחלוטין.

ההסתברויות בבסיסה "נכנסת לעבודה" במצבי אי וודאות.  
האם אנו יכולים להשתמש בהסתברויות בכדי לכמת את אי וודאות החיפוש שלנו?

# נושאים הסתברותיים ב- IR

- מודל הסתברותי קלאסי של אחזור מידע
  - עקרון הדירוג ההסתברותי
  - Binary independence model ( $\approx$  Naïve Bayes)
  - (Okapi) BM25
- רשתות בייסיאניות לאחזור טקסט
  - בהנתן אירוע שהתרחש, ניבוי הסבירות של כל אחת מכמה סיבות אפשריות גרם לאירוע.
- מודל שפתי ל- IR
- שיטות הסתברותיות הן אחד הנושאים הוותיקים ביותר אך גם אחד הנושאים החמים כרגע ב- IR
  - מסורתית: רעיונות פשוטים/מדויקים, אך בעבר לא הציגו את הביצועים הטובים ביותר
  - היום המצב שונה

# בעיית דירוג המסמכים

- יש בידינו מאגר של מסמכים
- המשתמש מזין שאילתה
- צריך להחזיר רשימה של מסמכים
- **שיטת דירוג מסמכים חיונית למערכות IR מודרניות**
  - באיזה סדר נציג למשתמש את התוצאות?
  - אנו מעוניינים שהתוצאות הטובות ביותר יהיו ראשונות
- **אידיאלי, דרג מסמכים בהתאם לרלוונטיות של מסמך בהתאם לשאילתה**
  - $P(R=1 | \text{document}_i, \text{query})$