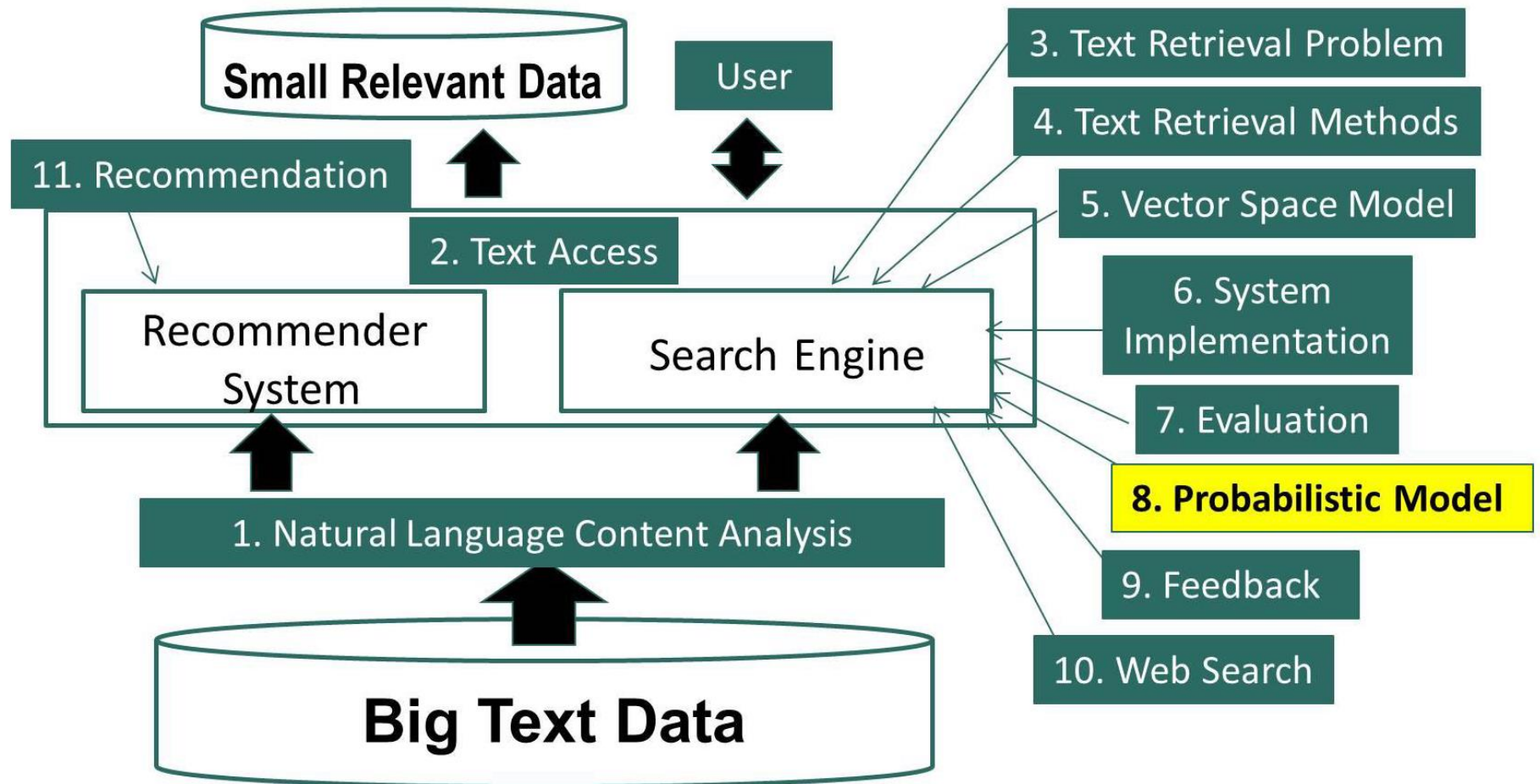


מודל אחזור הסתברותי



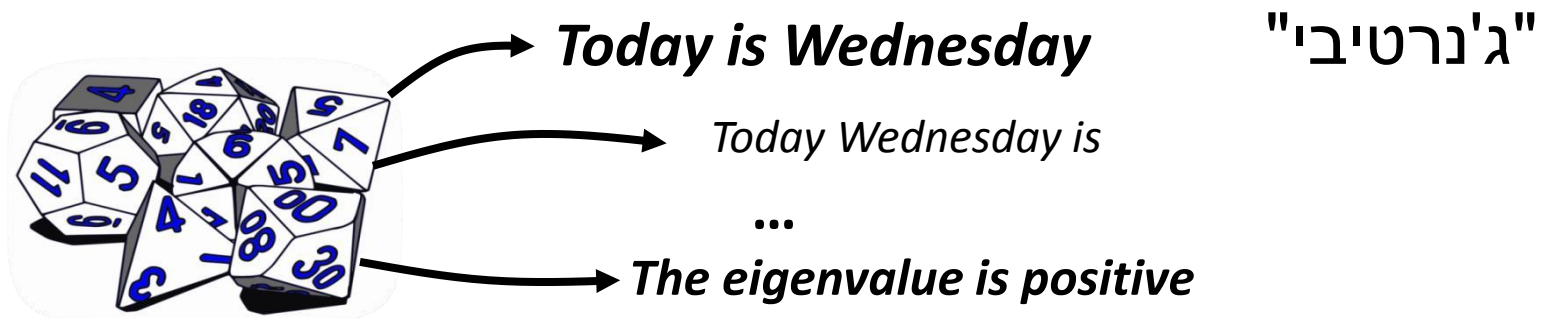
מהו מודל שפתי סטטיסטי (Language Model LM)?

- התפלגות הסתברות על פני רצפים של מילים
 - $p(\text{"Today is Wednesday"}) \approx 0.001$
 - $p(\text{"Today Wednesday is"}) \approx 0.0000000000000001$ *הסתברות נמוכה מאוד*
 - $p(\text{"The eigenvalue is positive"}) \approx 0.00001$

- ההתפלגות תלוייה בהקשר

– בהחלט ייתכן אם מדובר על דומיין של מתמטיקה
"The eigenvalue is positive" יהיה בעל הסתברות גבוהה
יותר מאשר "Today is Wednesday"

- ניתן גם להתייחס לזה כמנגנון הסתברותי המחולל (generate) טקסט, מסיבה זו הוא נקראה גם מודל



היעילות של ה-LM - Language Model

- מכמת את אי הוודאות בשפה הטבעית
- מאפשר לנו לענות על כל מיני שאלות, כגון:
 - בהינתן שאנו רואים את "John" ו-"feels" מה הסבירות שנראה "happy" לעומת "habit" כמילה הבאה?
(זיהוי דיבור, המילים נשמעות דומה)
 - בהינתן שאנו רואים במאמר חדשתי את המילה "בייסבול" שלוש פעמים ואת המילה "משחקים" פעם אחת, מה הסבירות שמדובר ב-"ספורט"?
(אחזור מידע, קיטלוג מסמכים)
 - בהינתן שמשתמש מעוניין בחדשות ספורט, באיזו סבירות המשתמש ישתמש במילה "בייסבול" בשאלתה?
(אחזור מידע)

מודל שפתי פשוט: יוניגרם מודל

היגיוני

- צור טקסט על ידי התייחסות לכל מילה באופן

עצמאי (דבר שאינו נכון בעולם האמיתי, יש קשר בין המילים במשפט)

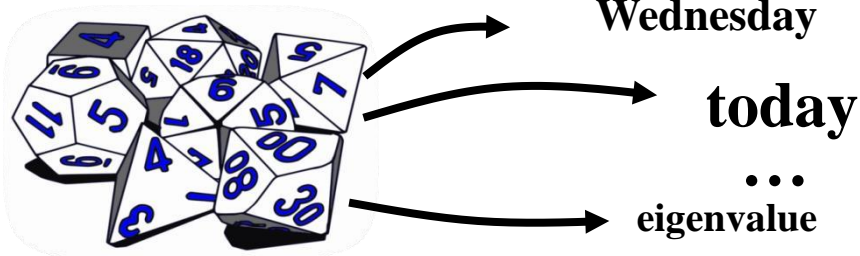
$$\rightarrow p(w_1 w_2 \dots w_n) = p(w_1)p(w_2)\dots p(w_n)$$

Parameters: $\{p(w_i)\}$ $p(w_1) + \dots + p(w_N) = 1$ (N is voc. size)

- טקסט הוא דגימה של מילים מהתפלגות מסויימת

הזרקה

כמה מילים המופיעות בפרק



$p(\text{"today is Wed"})$

$= p(\text{"today"})p(\text{"is"})p(\text{"Wed"})$

$= 0.0002 \times 0.001 \times 0.000015$

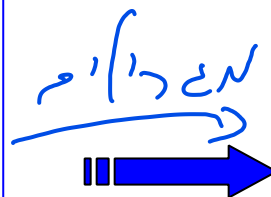
10 עשרות POW (הקטנה יקרה)

יצירת טקסט ע"י מודל יוניגרם

Unigram LM $p(w|\theta)$ **Sampling** → Document =?

Topic 1:
Text mining

text 0.2
mining 0.1
association 0.01
clustering 0.02
...
food 0.00001

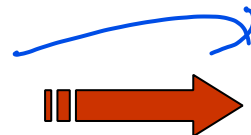


Text mining
article

הוא

Topic 2:
Health

...
food 0.25
nutrition 0.1
healthy 0.05
diet 0.02
...



Food nutrition
article

הערכת יוניגרם מודל

Unigram LM $p(w|\theta)$ **Estimation** ← Text Mining Paper d

Total #words=100

10/100 → text ?
 5/100 → mining ?
 3/100 → association ?
 3/100 → database ?
 ...
 1/100 → query ?

הנחיות למאמר -
 על מילון המילים
 וההסתברות



text 10
 mining 5
 association 3
 database 3
 algorithm 2
 ...
 query 1
 efficient 1

Maximum Likelihood (ML) Estimator:

$$p(w|\theta) = p(w|d) = \frac{c(w, d)}{|d|}$$

הסתברות מילון $c(w, d)$
 מספר מילים $|d|$
 מספר מילים d מההסתברות $p(w|d)$

האם זה הכי טוב?

כן, אבל ספציפית לנתונים אלו
 משקל המילים...

אורך מילים ב-Stops

LMs לייצוג נושא

המיון של אסוף מסמכים

המיון של המילים

General Background
English Text

B



the	0.03
a	0.02
is	0.015
we	0.01
...	
food	0.003
computer	0.00001
text	0.000006
...	

Computer Science
Papers

C



the	0.032
a	0.019
is	0.014
we	0.011
...	
computer	0.004
software	0.0001
text	0.00006
...	

Text mining
article

d



the	0.031
...	
text	0.04
mining	0.035
association	0.03
clustering	0.005
computer	0.0009
...	
food	0.000001

Background LM: $p(w|B)$

Collection LM: $p(w|C)$

Document LM: $p(w|d)$

בהתבסס על נתונים שונים נקבל LM שונה,
LM תלוי בהתפלגות של המילים בדאטה-סט

LMs לניתוח הקשרים שונים

אלו מילים קשורות סמנטית למילה "computer"?

Topic LM: $p(w | \text{"computer"})$

מסמכים המכילים את
"computer" המילה



the	0.032
a	0.019
is	0.014
we	0.008
computer	0.004
software	0.0001

איזה LM יכול לתת לנו את
המידע על מילים שכיחות ביותר
שאנו נרצה "למחוק" אותן?

כדי לדעת את המילים...

LMs לניתוח הקשרים שונים

נניח: אלו מילים קשורות סמנטית למילה "computer"?

Topic LM: $p(w | \text{"computer"})$

the	0.032
a	0.019
is	0.014
we	0.008
computer	0.004
software	0.0001

Normalized Topic LM:

$$\frac{p(w | \text{"computer"})}{p(w | \mathbf{B})}$$

Computer	400
software	150
program	104
...	
text	3.0
...	
the	1.1
a	0.99
is	0.9
we	0.8

Background LM: $p(w | \mathbf{B})$

the	0.03
a	0.02
is	0.015
we	0.01
computer	0.00001
text	0.000006
...	

מסמכים המכילים את
"computer"

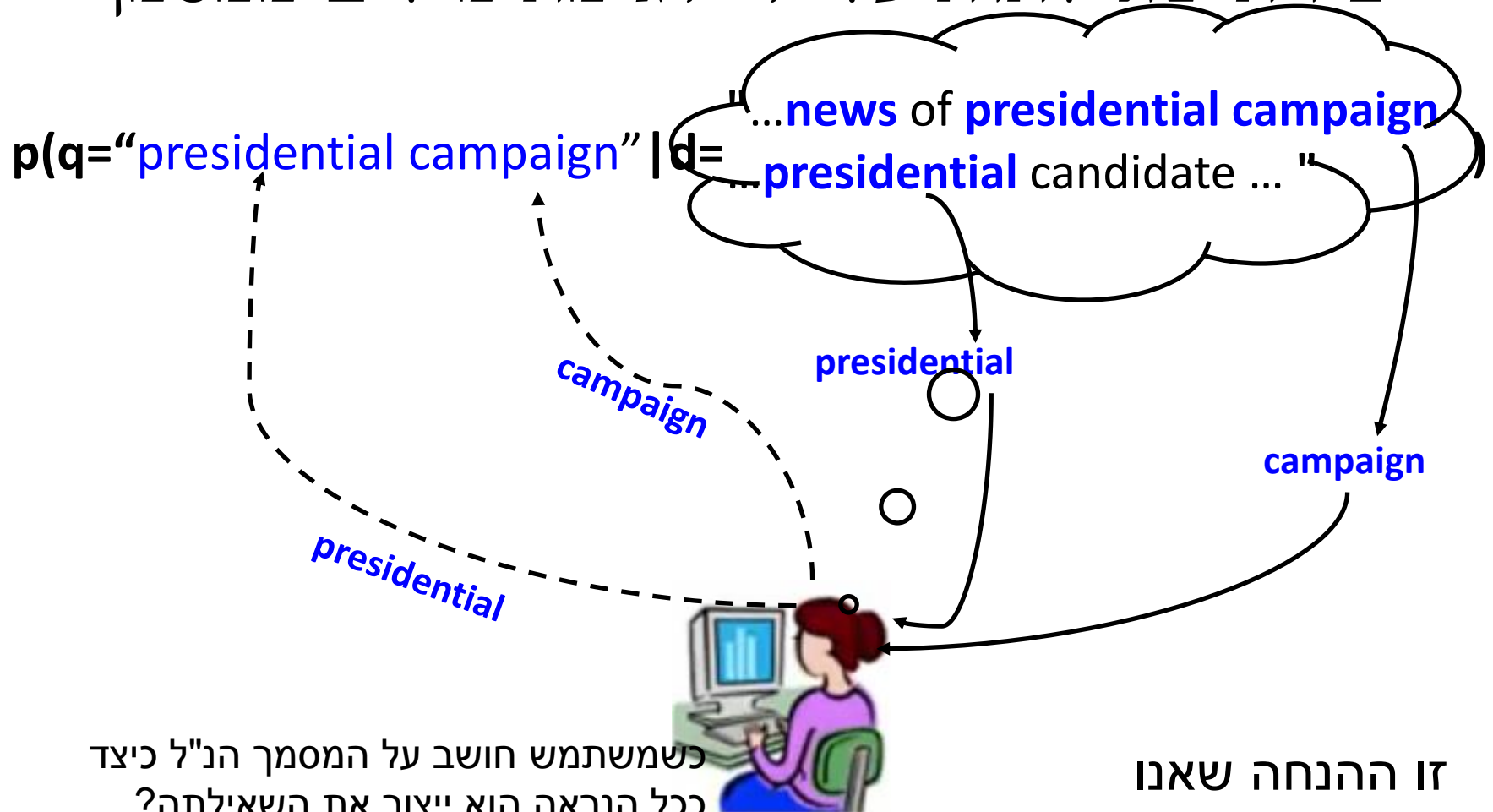
General Background
English Text

B

סיכום ביניים

- LM – מודל שפתי = התפלגות הסתברות על טקסט
- יוניגראם מודל = התפלגות המילים (בטקסט/מאגר נתון)
- שימושים ב-LM
 - מציאת נושא
 - מציאת קשרים בין מילים

יצירת שאילות על ידי דגימת מילים ממסמך



כשמשמש חושב על המסמך הנ"ל כיצד
ככל הנראה הוא ייצור את השאילתה?

זו ההנחה שאנו
מניחים על דרך
יצירת השאילתה

Unigram Query Likelihood

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

$$\begin{aligned} & p(\text{"presidential"} \mid d) * p(\text{"campaign"} \mid d) \\ &= \frac{c(\text{"presidential"}, d)}{|d|} * \frac{c(\text{"campaign"}, d)}{|d|} \end{aligned}$$

הנחה:

אין תלות בין המילים

האם יש הגיון ב- Query Likelihood

$$p(q = \text{"presidential campaign"}|d) = \frac{c(\text{"presidential"}, d)}{|d|} * \frac{c(\text{"campaign"}, d)}{|d|}$$

$$P(q|d4 = \dots \text{news of } \textbf{presidential campaign} \dots \textbf{presidential} \text{ candidate} \dots) = \frac{2}{|d4|} * \frac{1}{|d4|}$$

$$P(q|d3 = \dots \text{news of } \textbf{presidential campaign} \dots) = \frac{1}{|d3|} * \frac{1}{|d3|}$$

$$P(q|d2 = \dots \text{news about organic food} \textbf{campaign} \dots) = \frac{0}{|d2|} * \frac{1}{|d2|} = 0$$

כמו שחשבנו $d4 > d3 > d2$



האם יש הגיון ב- Query Likelihood

q = “presidential campaign **update**”
...אולי...

$$P(q | d4 = \dots \text{news of } \textbf{presidential campaign} \dots \textbf{presidential candidate} \dots) = \frac{2}{|d4|} * \frac{1}{|d4|} * \frac{0}{|d4|} = 0$$

$$P(q | d3 = \dots \text{news of } \textbf{presidential campaign} \dots) = \frac{1}{|d3|} * \frac{1}{|d3|} * \frac{0}{|d3|} = 0$$

$$P(q | d2 = \dots \text{news about organic food} \dots \textbf{campaign} \dots) = \frac{0}{|d2|} * \frac{1}{|d2|} * \frac{0}{|d2|} = 0$$

בגלל ההנחה שכל מילה בשאילתה חייבת להופיע במסמך

בגלל איזו הנחה עלתה הבעיה? כיצד נתקן זאת?

בשביל לתקן את זה אנחנו צריכים לקחת בחשבון שיהיו

מילים בשאילתה שאינן מופיעות במסמך.

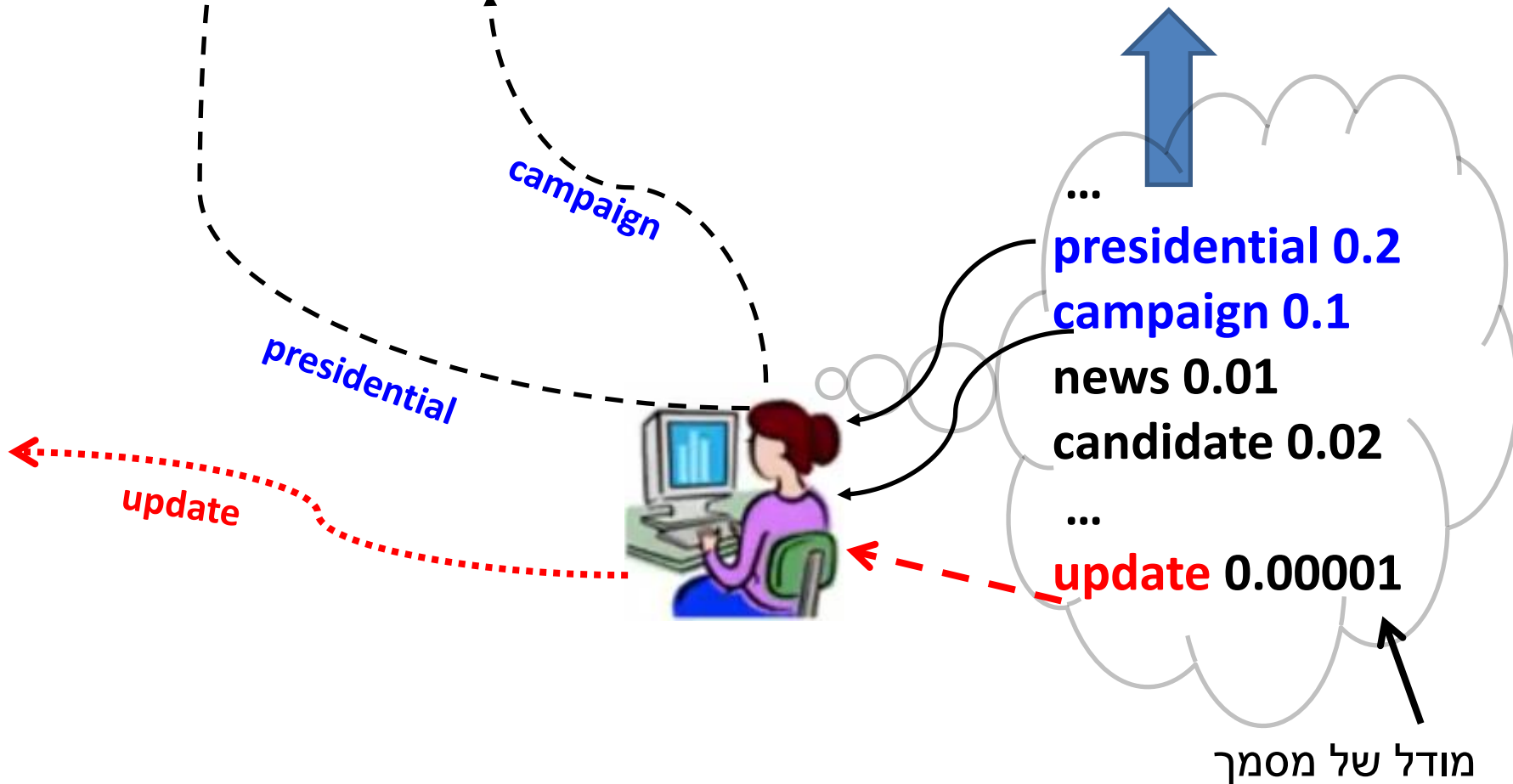
נכנסו מילים חדשות
משפטים מלאים

מודל משופר: המשתמש בוחר מילים מתוך

באקראיות גבוהה - מילים...

מודל של המסמך

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



חישוב השאילתה

Query q = "data mining alg`"

מסמך

d_1

Text mining
article

$P(w/d_1)$

text 0.2
mining 0.1
association 0.01
clustering 0.02
...
food 0.00001

$$\begin{aligned} P(\text{"data mining alg`"} / d_1) = & p(\text{"data"} / d_1) * \\ & p(\text{"mining"} / d_1) * \\ & p(\text{"alg`"} / d_1) \end{aligned}$$

d_2

Food nutrition
article

$P(w/d_2)$

...
food 0.25
nutrition 0.1
healthy 0.05
diet 0.02
...

$$\begin{aligned} P(\text{"data mining alg`"} / d_2) = & p(\text{"data"} / d_2) * \\ & p(\text{"mining"} / d_2) * \\ & p(\text{"alg`"} / d_2) \end{aligned}$$

סיכום ביניים: דירוג מבוסס QL

$$q = w_1 w_2 \dots w_n$$

$$p(q | d) = p(w_1 | d) * \dots * p(w_n | d)$$

$$f(q, d) = \log(p(q|d)) = \sum_{i=1}^n \log(p(w_i|d)) = \sum_{w \in V} c(w, q) \log(p(w|d))$$

Handwritten notes: $\log(a \cdot b \cdot c \dots)$ above the first sum; $\log(a) + \log(b) + \log(c) \dots$ above the second sum.

מדוע לעבור לחישוב עם \log ?

הסתברויות נמוכות מובילות למכפלה מאוד נמוכה ולחריגה (under flow) בתוצאות

Document language model

Handwritten notes: "לפי - שאלה" and "הסתברות" with arrows pointing to the $p(w|d)$ term in the equation above.

המעבר ממכפלה לחיבור של \log פשוטה

הסיגמא הראשונה מבוצעת על כל המילים בשאילתה

הסיגמא השנייה מבוצעת על כל המילים במילון אבל עם מכפלה של נוכחות המילה בשאילתה

Retrieval problem \rightarrow Estimation of $p(w_i | d)$

Handwritten notes: "בהתאמה" and "קלאס" with arrows pointing to the retrieval problem and the probability term respectively.

Different estimation methods \rightarrow different ranking functions

Handwritten notes: "הצרכה שלנו" and "לדאגה שלקול" with arrows pointing to the estimation methods and ranking functions respectively.

פונקציית דירוג מבוסס QL

$$q = w_1 w_2 \dots w_n \quad p(q | d) * p(w_1 | d) * \dots * p(w_n | d)$$

$$f(q, d) = \log(p(q|d)) = \sum_{i=1}^n \log(p(w_i|d)) = \sum_{w \in V} c(w, q) \log(p(w|d))$$

כיצד אנו נעריך את $p(w|d)$?

יש מילון של מופעים ב אקסטקט
ובן במאלה.

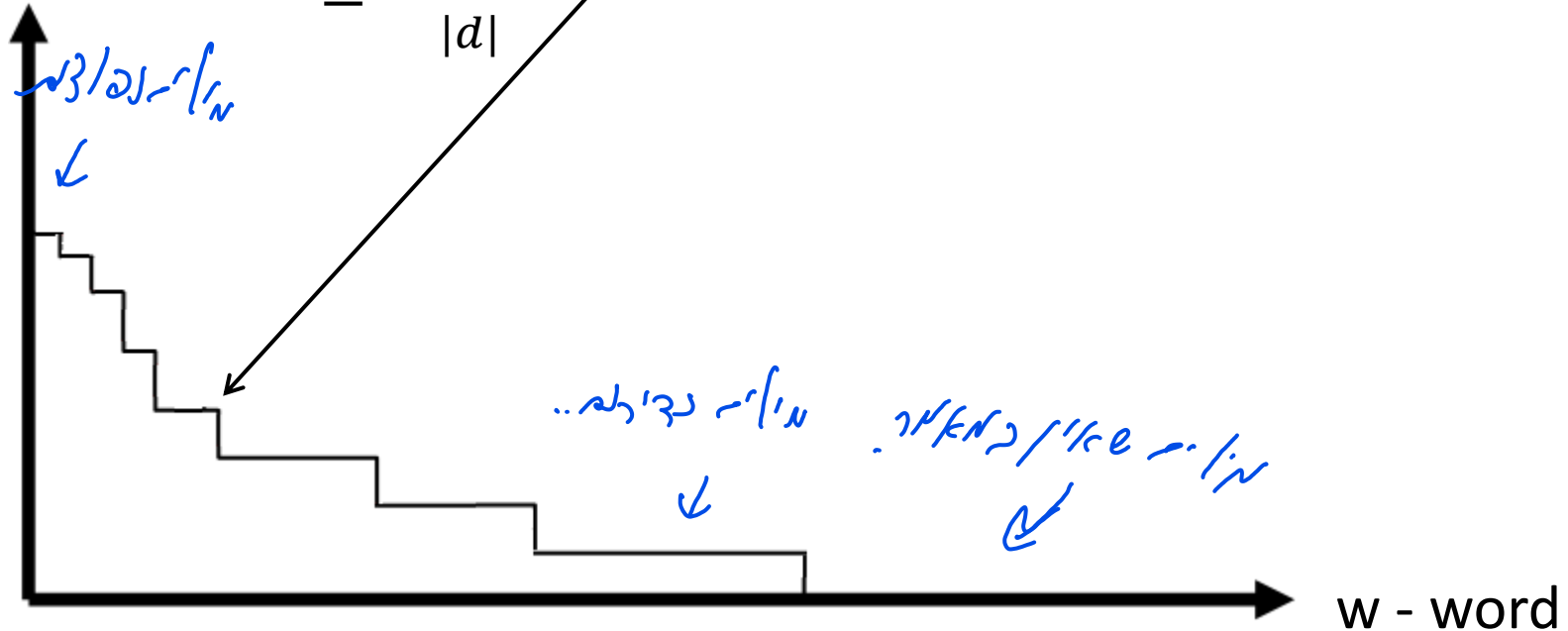
חישוב $p(w|d)$

Max. Likelihood Estimate

$$p_{ML}(w|d)$$

$$= \frac{c(w,d)}{|d|}$$

$p(w|d)$



קו ישר: מילים בעלי שכיחות זהה

(נאמר ה' - מבוטא כ'דאנקאצא...) נידב אמת ערב'י אמי'י, של, מופ'אל.

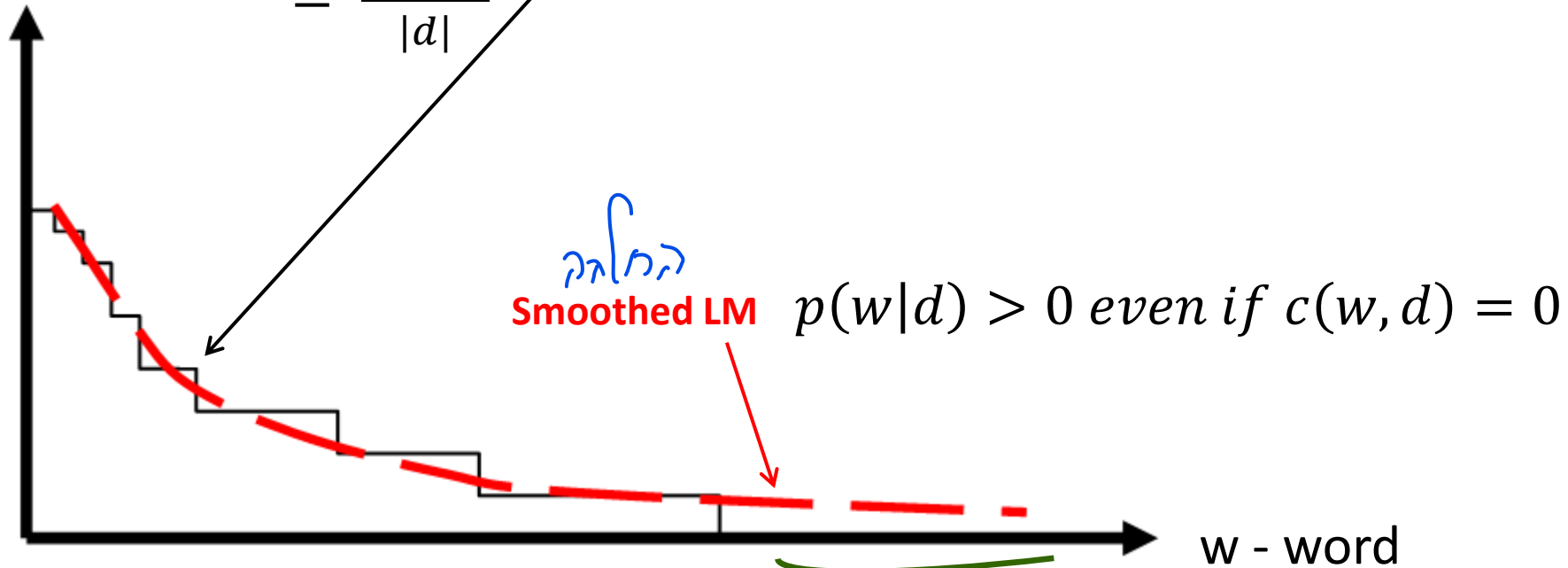
חישוב $p(w|d)$

Max. Likelihood Estimate

$$p_{ML}(w|d)$$

$$= \frac{c(w,d)}{|d|}$$

$p(w|d)$



אם המחבר היה כותב עוד מילים במסמך אז המילים המופיעות כאן היו בהסתברות חיובית, לכן לעשות smoothing זה הדבר הנכון

כיצד נעשה smoothing

- השאלה המרכזית: איזו הסתברות ניתן למילה שאינה במסמך? *הכיון - ההסתברות של המילה בטקסט..*
 - ההסתברות של מילה שאינה במסמך תהיה פרופורציונאלית להסתברות שלה הניתנת על ידי רפרנס LM (=מילני) *השפה*
- אפשרות אחת: רפרנס LM = אוסף LM

Discounted ML estimate

$$p(w|d) = \begin{cases} p_{seen}(w|d) & \text{if } w \text{ is seen in } d \\ \alpha_d p(w|C) & \text{otherwise} \end{cases}$$

מילני דלוקס *לא מילני*

Collection language model

כמה מילני דלוקס קווקטין

α – מקדם של הקולקציה

כתיבה מחודשת של ה-smoothing

✓ - מילון רשפה

$$f(q, d) = \log(p(q|d)) = \sum_{i=1}^n \log(p(w_i|d)) = \sum_{w \in V} c(w, q) \log(p(w|d)) =$$

סכום המילים הקיימות

$$\sum_{w \in V, c(w, d) > 0} c(w, q) \log p_{seen}(w|d) + \sum_{w \in V, c(w, d) = 0} c(w, q) \log \alpha_d p(w|d)$$

מילים בשאלתה המופיעות במסמך

מילים בשאלתה שאינן מופיעות במסמך

$$\sum_{w \in V} c(w, q) \log \alpha_d p(w|C) - \sum_{w \in V, c(w, d) > 0} c(w, q) \log \alpha_d p(w|d)$$

כל המילים בשאלתה

המילים בשאלתה שנמצאות במסמך

כתיבה מחודשת של ה-smoothing

$$\sum_{w \in V, c(w,d) > 0} c(w, q) \log p_{\text{seen}}(w|d) + \sum_{w \in V, c(w,d) = 0} c(w, q) \log \alpha_d p(w|C)$$

\uparrow מילים בשאלתה המופיעות במסמך \uparrow מילים בשאלתה שאינן מופיעות במסמך

$$\sum_{w \in V} c(w, q) \log \alpha_d p(w|C) - \sum_{w \in V, c(w,d) > 0} c(w, q) \log \alpha_d p(w|C)$$

\uparrow כל המילים בשאלתה \uparrow המילים בשאלתה שנמצאות במסמך



$$\sum_{w \in V, c(w,d) > 0} c(w, q) \log p_{\text{seen}}(w|d) - \sum_{w \in V, c(w,d) > 0} c(w, q) \log \alpha_d p(w|C)$$

$$= \sum_{w \in V, c(w,d) > 0} c(w, q) \log \frac{p_{\text{seen}}(w|d)}{\alpha_d p(w|C)}$$

כתיבה מחודשת של ה-smoothing

$$\sum_{w \in V, c(w,d) > 0} c(w,q) \log p_{seen}(w|d) - \sum_{w \in V, c(w,d) > 0} c(w,q) \log \alpha_d p(w|C)$$
$$= \sum_{w \in V, c(w,d) > 0} c(w,q) \log \frac{p_{seen}(w|d)}{\alpha_d p(w|C)}$$



$$\sum_{w \in V} c(w,q) \log \alpha_d p(w|C) = \overset{\text{קרי} \cdot}{|q| \log \alpha_d} + \sum_{w \in V} c(w,q) \log p(w|C)$$

כתיבה מחודשת של ה-smoothing

$$\begin{aligned}
 & \sum_{w \in V, c(w,d) > 0} c(w,q) \log p_{\text{seen}}(w|d) + \sum_{w \in V, c(w,d) = 0} c(w,q) \log \alpha_d p(w|C) \\
 &= \sum_{w \in V} c(w,q) \log \alpha_d p(w|C) + \sum_{w \in V, c(w,d) > 0} c(w,q) \log p_{\text{seen}}(w|d) \\
 &\quad - \sum_{w \in V, c(w,d) > 0} c(w,q) \log \alpha_d p(w|C) =
 \end{aligned}$$

$$\sum_{w \in V, c(w,d) > 0} c(w,q) \log p_{\text{seen}}(w|d) + \sum_{w \in V, c(w,d) > 0} c(w,q) \log \frac{p_{\text{seen}}(w|d)}{\alpha_d p(w|C)}$$

$$\boxed{= \sum_{w \in V, c(w,d) > 0} c(w,q) \log \frac{p_{\text{seen}}(w|d)}{\alpha_d p(w|C)} + |q| \log \alpha_d + \sum_{w \in V} c(w,q) \log p(w|C)}$$

כך נכתב המרחיב את הסכום על כל המילים (כולן) ולא רק על המילים שראו.

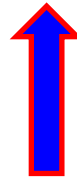
כתיבה מחודשת של ה-smoothing

הנוסחה הראשונה



$$\sum_{w \in V, c(w,d) > 0} c(w, q) \log(p_{seen}(w|d)) + \sum_{w \in V, c(w,d) = 0} c(w, q) \log(\alpha_d(p(w|C)))$$

$$= \sum_{w \in V, c(w,d) > 0} c(w, q) \log \frac{p_{seen}(w|d)}{\alpha_d p(w|C)} + |q| \log \alpha_d + \sum_{w \in V} c(w, q) \log(p(w|C))$$



הנוסחה הסופית

תזכורת: נוסחת ה-IDF-TF

מספר מופעים של מילה/מונח/מאפיין
ספציפי בשאילתה

Total # of docs in collection

$$f(q, d) = \sum_{i=1}^N x_i y_i = \sum_{w \in q \cap d} c(w, q) c(w, d) \log \frac{M + 1}{df(w)}$$

All matched query words in d

Doc Frequency

$$\text{IDF}(w) = \log \frac{M + 1}{df(w)}$$

נוסחה מעט שונה עם אותם מרכיבים

$$\sum_{w \in V, c(w,d) > 0} c(w, q) \log p_{seen}(w|d) + \sum_{w \in V, c(w,d) = 0} c(w, q) \log \alpha_d p(w|C)$$

$$= \sum_{w \in V, c(w,d) > 0} c(w, q) \log \frac{p_{seen}(w|d)}{\alpha_d p(w|C)} + |q| \log \alpha_d + \sum_{w \in V} c(w, q) \log p(w|C)$$

$$\sum_{w_i \in q \cap d} \left[\log \frac{p_{seen}(w_i|d)}{\alpha_d p(w_i|C)} \right] + n * \log \alpha_d + \sum_{i=1}^n \log p(w_i|C)$$

היתרונות של הכתיבה המחודשת

Smoothing עם $p(w|C)$ מכיל תכונות של TF-IDF ממושקל עם נרמול בגודל מסמך

מספר המסמכים הכולל באוסף

$$IDF(w) = \log \frac{M + 1}{df(w)}$$

בכמה מסמכים נמצאת המילה

$TF(w) = c(w, d)$

מספר מופעי המילה בשאלתה

מספר מופעי המילה במסמך

$$FT - IDF = \sum_{w_i \in q \cap d} c(w, q) c(w, d) \log \frac{M + 1}{df(w)}$$

כגודל המופעים במסמך כך גודל המשקל

$$\sum_{w_i \in q \cap d} c(w, q) \left[\log \frac{p_{seen}(w_i|d)}{\alpha_d p(w_i|C)} \right] + n * \log \alpha_d + \sum_{i=1}^n \log p(w_i|C)$$

נרמול באורך המשפט, בהמשך יהיה ברור יותר

פופולריות של מונח באוסף

היתרונות של הכתיבה המחודשת

Smoothing עם $p(w|C)$ מכיל תכונות של TF-IDF ממושקל עם נרמול בגודל מסמך

מספר המסמכים הכולל באוסף

$$IDF(w) = \log \frac{M + 1}{df(w)}$$

בכמה מסמכים נמצאת המילה

$$TF(w) = c(w, d)$$

מספר מופעי המילה בשאלתה

מספר מופעי המילה במסמך

$$FT - IDF = \sum_{w_i \in q \cap d} c(w, q) c(w, d) \log \frac{M + 1}{df(w)}$$

מספר קבוע עבור כל המסמכים,

אינו משפיע על הדירוג
(אין טימן להחזיק)

$$\sum_{w_i \in q \cap d} c(w, q) \left[\log \frac{p_{seen}(w_i|d)}{\alpha_d p(w_i|C)} \right] + n * \log \alpha_d + \sum_{i=1}^n \log p(w_i|C)$$

העניין הוא
שלא צריך דירוג
מסומן
אין הצטרפות
של המידע הזה...

סיכום ביניים על Smoothing

- Smoothing על $p(w|d)$ הכרחי בשביל QL

- בכלליות: Smoothing עם $p(w|C)$

- הנחה: מילה שאינה מופיעה ב- d ספציפי היא
פרופורציונאלית ל- $p(w|C)$

- מוביל אותנו לנוסחת דירוג כללית בעלת תכונות של TF-IDF
ממושקל ועם נרמול בגודל המסמך

- הציון/דירוג מבוסס בעיקר על סכום משקולות המונחים
המופיעים בשאילתה

- עדיין נותרה השאלה כיצד בדיוק יתבצע ה-
Smoothing

Query Likelihood + Smoothing with $p(w|C)$

$$\log p(q|d) = \sum_{w_i \in q \cap d} c(w, q) \left[\log \frac{p_{\text{seen}}(w_i|d)}{\alpha_d p(w_i|C)} \right] + n * \log \alpha_d + \sum_{w \in V} c(w, q) \log p(w|C)$$

$$f(q, d) = \sum_{w_i \in q \cap d} c(w, q) \left[\log \frac{p_{\text{seen}}(w_i|d)}{\alpha_d p(w_i|C)} \right] + n * \log \alpha_d$$

איך נחשב בדיוק סיכוי של מילה?

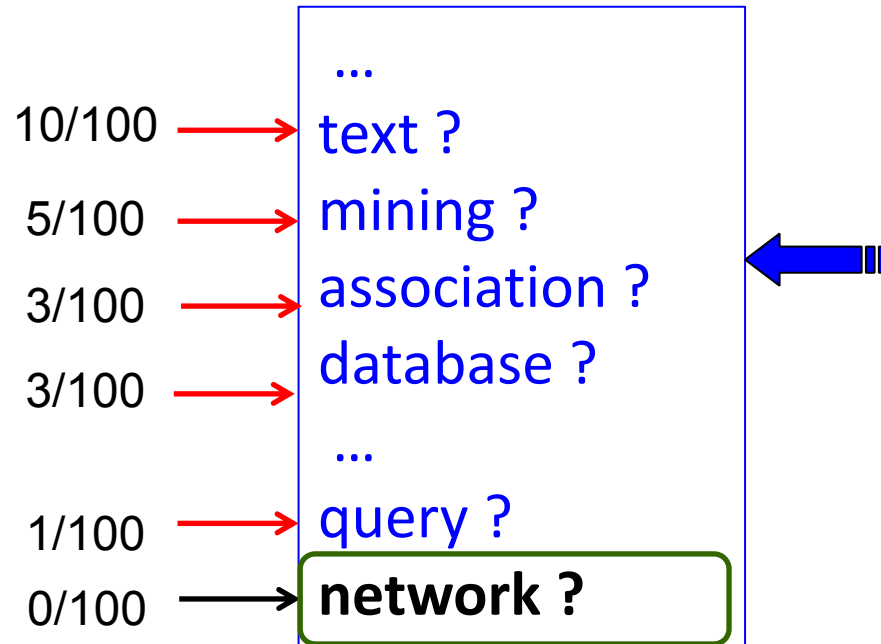
$P_{\text{seen}}(w_i|d)?$
 $\alpha_d?$

How to smooth $p(w|d)$?

איך נקבע את אלפא?

אינטרפולציה לינארית (Jelinek-Mercer) Smoothing

Unigram LM $p(w|\theta)$



Document d
Total #words=100

text 10
mining 5
association 3
database 3
algorithm 2
...
query 1
efficient 1

Collection LM
 $P(w|C)$

the 0.1
a 0.08
..
computer 0.02
database 0.01
.....
text 0.001
network 0.001
mining 0.0009
...

$$p(w|d) = \underbrace{(1 - \lambda)}_{\text{weight}} \frac{c(w, d)}{|d|} + \lambda \underbrace{p(w, C)}_{\text{collection LM}} \quad \lambda \in [0, 1]$$

משקל
של
הקולקציה

אינטרפולציה לינארית (Jelinek-Mercer) Smoothing

Collection LM
 $P(w|C)$

the	0.1
a	0.08
..	
computer	0.02
database	0.01
.....	
text	0.001
network	0.001
mining	0.0009
...	

כך מקבלים הסתברות
שאיננה שלילית גם
למונחים שאינם
במסמך

אינטרפולציה לינארית

← דוגמה

$$(1 - \lambda) + \lambda = 1$$

$$p(w|d) = (1 - \lambda) \frac{c(w, d)}{|d|} + \lambda p(w, C) \quad \lambda \in [0, 1]$$

Max` likelihood

Smoothing
parameter

Collection LM

אינטרפולציה לינארית (Jelinek-Mercer)

Smoothing

Unigram LM $p(w|\theta)?$

	...
10/100	→ text ?
5/100	→ mining ?
3/100	→ association ?
3/100	→ database ?
	...
1/100	→ query ?
0/100	→ network ?

Document d
Total #words=100

text 10
mining 5
association 3
database 3
algorithm 2
...
query 1
efficient 1

Collection LM
 $P(w|C)$

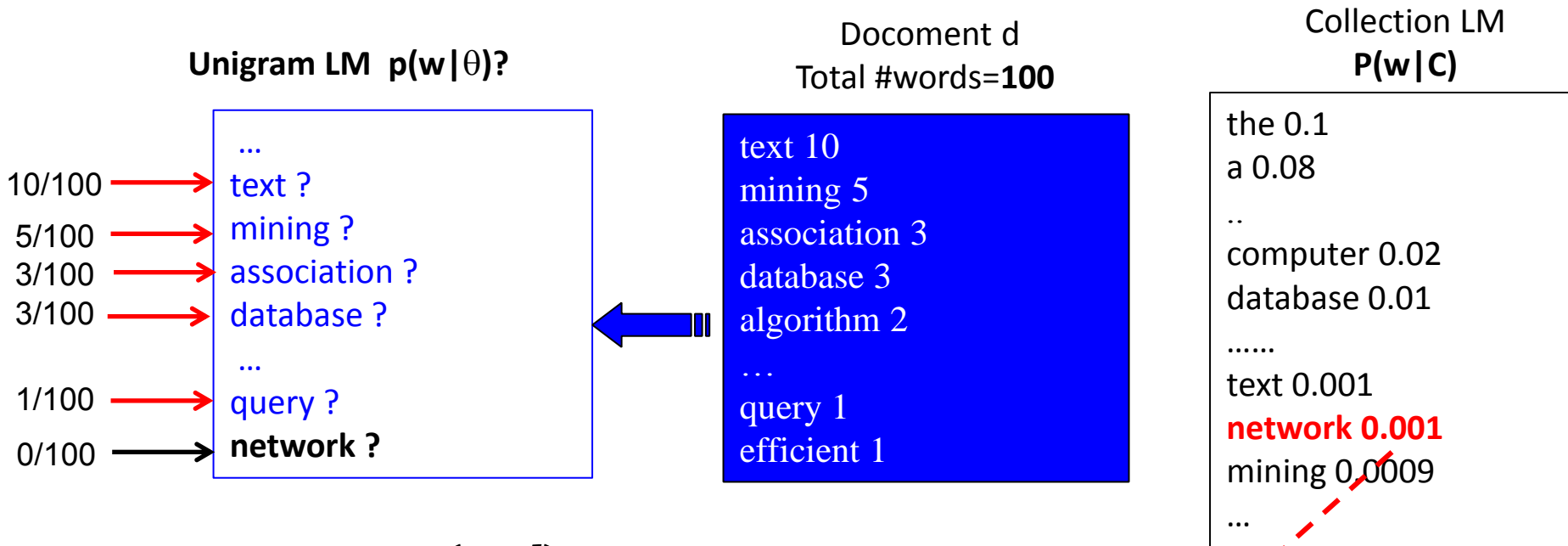
the 0.1
a 0.08
..
computer 0.02
database 0.01
.....
text 0.001
network 0.001
mining 0.0009
...

$$p(w|d) = (1 - \lambda) \frac{c(w, d)}{|d|} + \lambda p(w, C) \quad \lambda \in [0, 1]$$

$$p(\text{"text"}|d) = (1 - \lambda) \frac{10}{100} + \lambda * 0.001$$

אינטרפולציה לינארית (Jelinek-Mercer)

Smoothing



$$p(w|d) = (1 - \lambda) \frac{c(w, d)}{|d|} + \lambda p(w, C) \quad \lambda \in [0, 1]$$

zero

$$p(\text{"network"}|d) = \lambda * \mathbf{0.001}$$

על מנת להימנע מבעיה של חלוקת 0 ב-0

$$p(w|d) = (1 - \lambda) \frac{c(w, d)}{|d|} + \lambda p(w, C) \quad \lambda \in [0, 1]$$

Max' likelihood

Constant
Smoothing
parameter

Collection LM

$$p(w|d) = \frac{c(w, d) + \mu p(w|C)}{|d| + \mu} = \frac{|d|}{|d| + \mu} \frac{c(w, d)}{|d|} + \frac{\mu}{|d| + \mu} p(w|C) \quad \mu \in [0, \infty)$$

גודל מסמך
ממוצע

dynamic

Smoothing
parameter

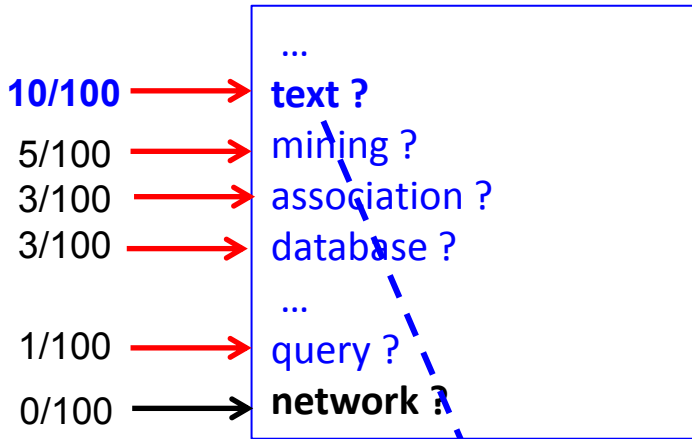
ככל שהמסמך גדול יותר
כך המקדם קטן יותר

$$\frac{|d|}{|d| + \mu} + \frac{\mu}{|d| + \mu} = 1$$

ומילגמל/קפ"ס ולכן ניצב קבועה מעלה כמ"ק.

Dirichlet Prior (Bayesian) Smoothing

Unigram LM $p(w|\theta)?$



Document d
Total #words=100

text 10
mining 5
association 3
database 3
algorithm 2
...
query 1
efficient 1

Collection LM
 $P(w|C)$

the 0.1
a 0.08
..
computer 0.02
database 0.01
.....
text 0.001
network 0.001
mining 0.0009
...

$$p(w|d) = \frac{c(w, d) + \mu p(w|C)}{|d| + \mu} = \frac{|d|}{|d| + \mu} \frac{c(w, d)}{|d|} + \frac{\mu}{|d| + \mu} p(w|C)$$

$\mu \in [0, \infty)$

$$p(\text{"text"}|d) = \frac{10 + \mu * 0.001}{100 + \mu}$$

Dirichlet Prior (Bayesian) Smoothing

Unigram LM $p(w|\theta)?$

...	
10/100	→ text ?
5/100	→ mining ?
3/100	→ association ?
3/100	→ database ?
...	
1/100	→ query ?
0/100	→ network ?

Document d
Total #words=100

text 10
mining 5
association 3
database 3
algorithm 2
...
query 1
efficient 1

Collection LM
 $P(w|C)$

the 0.1
a 0.08
..
computer 0.02
database 0.01
.....
text 0.001
network 0.001
mining 0.0009
...

$$p(w|d) = \frac{c(w, d) + \mu p(w|C)}{|d| + \mu} = \frac{|d|}{|d| + \mu} \frac{c(w, d)}{|d|} + \frac{\mu}{|d| + \mu} p(w|C)$$

$\mu \in [0, \infty)$

zero (pointing to the $|d|$ term in the numerator of the second fraction)

$$p(\text{"network"}|d) = \frac{\mu * \mathbf{0.001}}{100 + \mu}$$

פונקציית דירוג עבור (Dirichlet Prior) Smoothing

$$f(q, d) = \sum_{w_i \in q \cap d} c(w, q) \log \left[\frac{p_{\text{seen}}(w_i | d)}{\alpha_d p(w_i | C)} \right] + n * \log \alpha_d$$

Maximum Likelihood (ML) :

$$p(w|d) = \frac{c(w, d)}{|d|}$$

$$p(w|d) = (1 - \lambda) \frac{c(w, d)}{|d|} + \lambda p(w, C) \quad \lambda \in [0, 1]$$

$$\frac{p_{\text{seen}}(w|d)}{\alpha_d p(w|C)} = \frac{(1 - \lambda) p_{ML}(w|d) + \lambda p(w|C)}{\lambda p(w|C)} = 1 + \frac{1 + \lambda}{\lambda} \frac{c(w, d)}{|d| p(w|C)}$$

$$f_{JM}(q, d) = \sum_{w_i \in q \cap d} c(w, q) \log \left[1 + \frac{1 + \lambda}{\lambda} \frac{c(w, d)}{|d| p(w|C)} \right] + n * \log \alpha_d$$

פונקציית דירוג עכור Smoothing

$$f(q, d) = \sum_{w_i \in q \cap d} c(w, q) \log \left[\frac{p_{\text{seen}}(w_i | d)}{\alpha_d p(w_i | C)} \right] + n * \log \alpha_d$$

$$\frac{p_{\text{seen}}(w | d)}{\alpha_d p(w | C)} = \frac{(1 - \lambda) p_{ML}(w | d) + \lambda p(w | C)}{\lambda p(w | C)} = 1 + \frac{1 + \lambda}{\lambda} \frac{c(w, d)}{|d| p(w | C)}$$

$$f_{JM}(q, d) = \sum_{w_i \in q \cap d} c(w, q) \log \left[1 + \frac{1 + \lambda}{\lambda} \frac{c(w, d)}{|d| p(w | C)} \right] + n * \log \alpha_d$$

$\lambda = \alpha_d \rightarrow n * \log \alpha_d$ is constant ← נוסחה קבועה

$$f_{JM}(q, d) = \left[\sum_{w_i \in q \cap d} c(w, q) \log \left[1 + \frac{1 + \lambda}{\lambda} \frac{c(w, d)}{|d| p(w | C)} \right] \right] + \cancel{n * \log \alpha_d}$$

פונקציית דירוג עבור Smoothing

$$f_{JM}(q, d) = \sum_{w_i \in q \cap d} c(w, q) \log \left[1 + \frac{1 + \lambda}{\lambda} \frac{c(w, d)}{|d| p(w|C)} \right]$$

TF = $c(w, d)$

IDF = $p(w|C)$

$|d|$ נרמול באורך המשפט

הנוסחה הזו "תופשת" את כל הרעיונות הנ"ל

בעוד ב VSM – השתמשנו בהיוריסטיקות בשביל שלושת הרעיונות במודל

ההסתברותי זה נובע ישירות (לאחר הנחות בסיסיות)

אין קשר בין המילים (לדוגמה)
קריאה המיועדת למטרה אחרת

$|d| p(w|C)$ - משמעות של הסיכוי לקבל את המילים במסמך מתוך כל הקולקציה של המסמכים

מספר מופעי המילה בפועל מנורמל במספר המופעים של המילה שהיינו מצפים לקבל מכל הקולקציה

$$\frac{c(w, d)}{|d| p(w|C)}$$

פונקציית דירוג עבור (Dirichlet Prior) Smoothing

$$f(q, d) = \sum_{w_i \in q \cap d} c(w, q) \log \left[\frac{p_{\text{seen}}(w_i | d)}{\alpha_d p(w_i | C)} \right] + n * \log \alpha_d$$

$$\underline{p(w|d)} = \frac{c(w,d) + \mu p(w|C)}{|d| + \mu} = \frac{|d|}{|d| + \mu} \frac{c(w,d)}{|d|} + \frac{\mu}{|d| + \mu} p(w|C) \quad \mu \in [0, \infty)$$

$$\frac{p_{\text{seen}}(w|d)}{\alpha_d p(w|C)} = \frac{\frac{c(w,d) + \mu p(w|C)}{|d| + \mu}}{\frac{\mu p(w|C)}{|d| + \mu}} = 1 + \frac{c(w,d)}{\mu p(w|C)} \quad \alpha_d = \frac{\mu}{|d| + \mu}$$

$$f_{\text{DIR}}(q, d) = \left[\sum_{w_i \in q \cap d} c(w, q) \log \left[1 + \frac{c(w, d)}{\mu p(w|C)} \right] \right] + n * \log \frac{\mu}{|d| + \mu}$$

$$p(w|d) = \frac{c(w,d) + \mu p(w|C)}{|d| + \mu} = \frac{|d|}{|d| + \mu} \frac{c(w,d)}{|d|} + \frac{\mu}{|d| + \mu} p(w|C) \quad \mu \in [0, \infty)$$

$$\frac{p_{seen}(w|d)}{\alpha_d p(w|C)} = \frac{\frac{c(w,d) + \mu p(w|C)}{|d| + \mu}}{\frac{\mu p(w|C)}{|d| + \mu}} = 1 + \frac{c(w,d)}{\mu p(w|C)} \quad \alpha_d = \frac{\mu}{|d| + \mu}$$

אנחנו משווים את מספר המופעים האמתי במשפט עם דגימה של μ מילים מהמאגר

$$\alpha_d = \frac{\mu}{|d| + \mu}$$

היכן הנרמול באורך המשפט?

$$f_{DIR}(q, d) = \left[\sum_{w_i \in q \cap d} c(w, q) \log \left[1 + \frac{c(w, d)}{\mu p(w|C)} \right] \right] + n * \log \frac{\mu}{|d| + \mu}$$

סיכום הסמoothing

- למדנו שתי שיטות של smoothing
 - Jelinek-Mercer: אינטרפולציה לינארית, פרמטר קבוע
 - Dirichlet Prior: אינטרפולציה אדפטיבית
- שתיהן מובילות לפונקציות IR עם הנחות יסוד מוגדרות (פחות היוריסטיות)
 - מכילות תבונות של FT-IDF עם נרמול באורך המסמך
 - בעלות פרמטר smoothing אחד בלבד

סיכום

- פונקציות דירוג יעילות המתקבלות באמצעות מודל הסתברותי "נקי" תחת ההנחות הבאות:

1. $\text{Relevance}(q,d) = p(R=1 | q,d) \approx p(q | d, R=1) \approx p(q | d)$

2. מילות השאילתה אינן תלויות בניהן (BoW) → תחילה שאני

3. Smoothing with $p(w | C)$ → הלאה נייב

4. שתי שיטות של Smoothing

- פחות היוריסטיות יחסית ל-VSM
- ישנן לא מעט הרחבות לפונקציות הדירוג