

Language Models

G&F

Sec. 2.3

a term-document matrix

	D ₁	D ₂	D ₃
a	1	1	1
arrived	0	1	1
damaged	1	0	0
delivery	0	1	0
five	1	0	0
gold	1	0	1
in	1	1	1
of	1	1	1
shipped	1	0	1
silver	0	2	0
truck	0	1	1
length	<u>7</u>	<u>8</u>	<u>7</u>

22 tokens

Given a query g we want to calculate for each document d_i the probability of d_i generating g

$$SC(g, d_i) = \prod_{t_j \in g} P(t_j | M_{d_i}) \prod_{t_j \notin g} (1 - P(t_j | M_{d_i}))$$

prob. of M_{d_i} generating a term in the query

prob. of M_{d_i} NOT generating a term NOT in the query

How to estimate $P(t_j | M_{D_i})$
 for a given term and document i ?

$$P(t_j | M_{D_i}) = p_{M_i}(t_j | M_{D_i}) \\ = \frac{tf(t_j, D_i)}{\text{length}(D_i)}$$

which is easy enough to compute,
 except that if a term t_j
 doesn't happen to occur in
 document d_i , $P(t_j | M_{D_i}) = 0$
 and that makes the whole
 product zero :- (

So what do we do when
 this happens?

Choices include:

$$1. \quad P(t_j | M_{D_i}) = \frac{\sum_{\text{all } d_i} tf_j}{\sum_{\text{all } d_i} \sum_{\text{all } t_j} tf_{ij}} = \frac{\text{occurrences of } t_j}{\text{occurrences of all terms}}$$

G&F
call this $\frac{c_{tj}}{cs}$

OR

2. $P(t_j | M_{d_i})$ ~~is~~ based on risk t_j as shown on p. 4 3/

For each term t_i in document d_j , calculate t_i

$P_{M_h}(t_i M_{d_i})$	length (d_j)			Average Prob.
	D_1	D_2	D_3	
a	1/7	1/8	1/7	$(\frac{2}{7} + \frac{1}{8})/3$
arrived	0	1/8	1/7	$(\frac{1}{8} + \frac{1}{7})/2$
damaged	1/7	0	0	1/7
delivery	0	1/8	0	1/8
fire	1/7	0	0	1/7
gold	1/7	0	1/7	1/7
in	1/7	1/8	1/7	$(\frac{2}{7} + \frac{1}{8})/3$
of	1/7	1/8	1/7	$(\frac{2}{7} + \frac{1}{8})/3$
shipment	1/7	0	1/7	1/7
silver	0	1/4	0	1/4
truck	0	1/8	1/7	$(\frac{1}{8} + \frac{1}{7})/2$

G&F Table 2.12 p. 49

Page (t) Average

Prob.

To calculate average probability,

$$\sum_{d_i: \text{ted}} P_{M_i}(t_i | M_{d_i})$$

$$= \frac{\text{sum of each row.}}{\text{\# of non-zero entries}}$$

↑
G&F Table 2.13 p. 50

but were not done!

Given these average probabilities, and the length of each document, we can compute

$$\overline{f_{t,d}} = P_{\text{avg}}(t) \times \text{length}(d_j)$$

$\overline{f_t}$	D_1	D_2	D_3	$P_{\text{avg}}(t)$
arrived		$8 \cdot \left(\frac{1}{8} + \frac{1}{7}\right) \cdot 2$		$\left(\frac{1}{8} + \frac{1}{7}\right) / 2$
shipped silver length	$\frac{1}{7}$	$\frac{2}{8}$	$\frac{1}{7}$	$\frac{1}{7}$ $\frac{1}{4}$

Given $\overline{f_t}$ for term t_i and document d_j ,

$$R_{t,d} = \left(\frac{1}{1 + \overline{f_t}} \right) \left(\frac{\overline{f_t}}{1 + \overline{f_t}} \right)^{\overline{f_{t,d}}}$$

$$R_{\text{arrived}, D_2} = \left(\frac{1}{1 + 1.071} \right) \left(\frac{1.071}{1 + 1.071} \right)^1$$

$$= \frac{1}{2.071} \cdot \frac{1.071}{2.071} = \frac{1.071}{4.289} = .249$$

$$R_{\text{silver}, D_2} = \left(\frac{1}{1 + 2} \right) \left(\frac{2}{1 + 2} \right)^2 = \left(\frac{1}{3} \cdot \left(\frac{2}{3} \right)^2 \right) = \frac{4}{27} = .148$$

Now we can (Finally!) calculate

$$P(t|M_d) = P_{ml}(t|M_d)^{(1-R_{td})} \times K$$

↑
 from p. 3
 ↓
 $P_{avg}(t)$

↑
 from ps. 4
 ←
 R_{td}

and this gives G&F Table 2.16

for $g =$ "gold silver truck" and document d_1 , for example,

$$P(g|M_{d_1}) = \prod_{t \in g} P(t|M_{d_1}) \times \prod_{\text{other terms}} (1 - P(t))$$

$\cdot 143 \times \cdot 091 \times \cdot 091$
 gold
 silver
 truck
 $(\cdot 00118) \times$

$\cdot \prod (1 - P(t, M_{d_1})) =$
 all other terms arrived damaged delivery
 $(1 - \cdot 141) (1 - \cdot 09) (1 - \cdot 143) (1 - \cdot 045)$
 size in of shipment
 $(1 - \cdot 143) (1 - \cdot 141) (1 - \cdot 141) (1 - \cdot 143)$

$= (\cdot 00118)$
 $(\cdot 86) (\cdot 91) (\cdot 86) (\cdot 95) (\cdot 86) (\cdot 86)$
 $(\cdot 86) (\cdot 86) (\cdot 91) (\cdot 91) =$
 $\cdot 0003$

which is close to Table 2.17