

$$d_1 = (t_1 \ t_2 \ t_3 \ t_4 \ t_5 \ t_6)$$

$$d_2 = (t_3 \ t_4 \ t_5 \ t_1 \ t_8 \ t_2 \ t_2 \ t_2)$$

$$d_3 = (t_6 \ t_9 \ t_4 \ t_1 \ t_8 \ t_2 \ t_3 \ t_1 \ t_4)$$

$$q = (t_4 \ t_3)$$

$$u = 6$$

$$\lambda_T = 0.85$$

$$\lambda_0 = 0.1$$

$$\lambda_U = 0.05$$

$$w = 4$$

$$D = \{d_1, d_2, d_3\}$$

$|d|$  = length of a document

(a) what is the value of unigram feature function for "t<sub>4</sub>" in document d<sub>3</sub>?

$$f_T(q_i, d) = \log_2 P(q_i | d_d)$$

$$f_T(t_4, d_3) = \log_2 P(t_4 | d_3)$$

$$= \log_2 \left( \frac{c(t_4, d_3) + u P(t_4, 1)}{|d_3| + u} \right)$$

$$= \log_2 \left( \frac{2 + 6 \left( \frac{1+1+2}{6+8+9} \right)}{9+6} \right)$$

$$= \log_2 \left( \frac{2 + 6(4/23)}{15} \right)$$

$$= [-2.301]$$

$$(b) f_0(t_4, t_3, d_3) = \log_2 \left( \frac{c(t_4, t_3, d_3) + u P_0(t_4, t_3, 1)}{|d_3| + u} \right)$$

$$= \log_2 \left( \frac{1 + 6 \cdot \left( \frac{0+1+0}{6+8+9} \right)}{9+6} \right)$$

$$\frac{1+2+1}{6+8+9} = 0.1789$$

$$(c) f_T(t_4, t_3, d_3) = \log_2 \left( \frac{0+6 \left( \frac{1+2+1}{6+8+9} \right)}{9+6} \right) = -2.876$$

(c)  $f_T(t_4, t_3, d_3)$  = score( $d_3, q$ )

(d) Total retrieval score for  $d_3$  = score( $d_3, q$ )

$$= \lambda_T \sum_{i=1}^n f_T(q_i, d_3) + \lambda_D f_T(t_4, t_3, d_3) + \lambda_U(t_4, t_3, d_3)$$

$$= 0.85(-2.301 - 2.876) + 0.1(-5.845) + 0.05(-2.876)$$

$$= -5.129$$

$$(e) f_T(t_4, d_2) = \log_2 \left( \frac{1+6 \left( \frac{1+1+2}{6+8+9} \right)}{8+6} \right) = -2.7769$$

$$f_T(t_3, d_2) = \log_2 \left( \frac{2+6 \left( \frac{1+2+1}{6+8+9} \right)}{8+6} \right) = -2.2022$$

$$f_D(t_4, t_3, d_2) = \log_2 \left( \frac{1+6 \left( \frac{0+1+0}{6+8+9} \right)}{8+6} \right) = -3.4771$$

$$f_U(t_4, t_3, d_2) = \log_2 \left( \frac{2+6 \left( \frac{1+2+1}{6+8+9} \right)}{8+6} \right) = -2.2015$$

$$score(d_2, q) = 0.85(-2.7769 - 2.2022) + 0.1(-3.4771) + 0.05$$

$$= -4.689$$

$$(f) f_T(t_4, d_1) = \log_2 \left( \frac{1+6 \left( \frac{1+1+2}{6+8+9} \right)}{6+6} \right) = -2.5546$$

$$f_T(t_3, d_1) = \log_2 \left( \frac{1+6(0.1789)}{12} \right) = -2.5546$$

$$f_D(t_4, t_3, d_1) = \log_2 \left( \frac{0+6 \left( \frac{0+1+0}{6+8+9} \right)}{12} \right)$$

$$f_U(t_4, t_3, d_1) = \log_2 \left( \frac{1+6(0.174)}{12} \right) = -5.5395$$

$$score(d_1, q) = 0.85(-2.5546 - 2.5546) + 0.1(-5.5395) + 0.05$$

$$= -5.02367$$

## Normalization

Date \_\_\_\_\_  
Page \_\_\_\_\_

$$V' = \left( V - \frac{\min(A)}{\max(A)} - \min(A) \right) \left( \frac{\text{new\_max}(A) - \text{new\_min}(A)}{\max(A) - \min(A)} \right) + \text{new\_min}(A)$$

1. For example:

$$V = \left( \frac{a_1}{\min}, \frac{a_2}{\min}, \frac{a_3}{\min}, \frac{a_4}{\max} \right) = (1000, 2000, 3000, 9000)$$

$$\bullet V = 1000$$

$$\min(A) = 1000$$

$$\max(A) = 9000$$

$$\text{new\_min}(A) = 0$$

$$\text{new\_max}(A) = 1$$

$$\therefore a_1 = \frac{1000 - 1000}{9000 - 1000} \cdot (1 - 0) + 0 = 0$$

$$\therefore a_2 = \frac{2000 - 1000}{9000 - 1000} \cdot (1 - 0) + 0 = 0.125$$

$$\therefore a_3 = \frac{3000 - 1000}{9000 - 1000} \cdot (1 - 0) + 0 = 0.25$$

$$\therefore a_4 = \frac{9000 - 1000}{9000 - 1000} \cdot (1 - 0) + 0 = 1$$

$$V' = (0, 0.125, 0.25, 1)$$

2.  $V = (2, 12, 7)$

$$V' = \left( \frac{2-2}{12-2} \cdot (1-0) + 0, \frac{12-2}{12-2}, \frac{7-2}{12-2} \right)$$

$$= (0, 1, 0.5)$$

K-MEAN clustering.

1.

	x	y
1	185	72
2	170	56
3	168	60
4	179	68
5	182	72
6	188	77

K=2

∴ Initial centroid

	x	y
c1	185	72
c2	170	56
<del>c3</del>	168	

∴ Euclidean distance  $[ (x, y), (a, b) ] :$ 

$$= \sqrt{(x-a)^2 + (y-b)^2}$$

$$ED_{c1-c1} = \sqrt{(185-185)^2 + (72-72)^2} = 0$$

$$ED_{c1-c2} = \sqrt{(170-185)^2 + (56-72)^2} \\ = 21.93$$

$$ED_{c2-c2} = 0$$

centroid

clusters	x	y	Assignment
c1	0	21.93	1
c2	21.93	0	2
<del>c3</del>			

Let's take next dataset  $(168, 60) = N$ 

$$ED_{c1-N} = \sqrt{(168-185)^2 + (60-72)^2} = 20.808$$

$$ED_{c2-N} = \sqrt{(168-170)^2 + (60-56)^2} = 4.472$$

N	x	y	Assi
	20.808	4.472	2 ( $\because y < x$ )

recalculate c2.

$$(170+168)/2$$

$$169$$

$$(60+56)/2$$

$$58$$

$$c2 = (169, 58)$$

Now  $N = (179, 68)$

$$ED_{C1-N} = \sqrt{(179-185)^2 + (68-72)^2} = 7.21103$$

$$ED_{C2-N} = \sqrt{(179-169)^2 + (68-56)^2} = 14.14214$$

$\therefore X < Y$  recalculate C1.

$$\therefore C_1 = ((185+179)/2, (72+68)/2) \\ = (182, 70)$$

Now  $N = (182, 72)$

$$ED_{C1-N} = \sqrt{(182-185)^2 + (72-70)^2} = 2$$

$$ED_{C2-N} = \sqrt{(182-169)^2 + (72-56)^2} = 19.10$$

$$\text{new } C_1 = ((182+182)/2, (70+72)/2) = (182, 71)$$

Now  $N = (188, 77)$

$$ED_{C1-N} = \sqrt{(188-182)^2 + (77-71)^2} = 8.4852$$

$$ED_{C2-N} = \sqrt{(188-169)^2 + (77-56)^2} = 26.87$$

$$\text{new } C_1 = ((182+188)/2, (71+77)/2) = (185, 74)$$

Final assignment:

Dataset	X	Y	Assignment
1	185	72	1
2	170	56	2
3	168	60	2
4	179	68	1
5	182	72	1
6	188	77	1

	Dataset	w	x	y	z
P1		3	1	0	4
P2		2	0	1.5	1
P3		1	3	4.5	5
P4		2.5	2.5	5	3.5
P5		6	3	2	0

	centroids	w	x	y	z
C1		3	5	4	4.5
C2		1.5	2.5	1	2
C3		2	3.5	1	0

For  $P_1 = (3, 1, 0, 4)$

$$ED_{P1-C1} = \sqrt{(3-3)^2 + (1-5)^2 + (0-4)^2 + (4-4.5)^2} = 32.25$$

$$ED_{P1-C2} = \sqrt{(3-1.5)^2 + (1-2.5)^2 + (0-1)^2 + (4-2)^2} = 9.5$$

$$ED_{P1-C3} = \sqrt{(3-2)^2 + (1-3.5)^2 + (0-1)^2 + (4-0)^2} = 24.25$$

$\therefore$  assignment = C2 ( $\because C_2 < C_1, C_2 < C_3$ )

$$\text{new } C_2 = \left( \frac{1.5+3}{2}, \frac{2.5+1}{2}, \frac{0+1}{2}, \frac{2+4}{2} \right) = (2.25, 1.75, 0.5, 3)$$

For  $P_2 = (2, 0, 1.5, 1)$

$$ED_{P2-C1} = \sqrt{(2-3)^2 + (0-5)^2 + (1.5-4)^2 + (1-4.5)^2} = 6.670$$

$$ED_{P2-C2} = \sqrt{(2-1.5)^2 + (0-1.75)^2 + (1.5-0.5)^2 + (1-3)^2} = 2.850$$

$$ED_{P2-C3} = \sqrt{(2-2)^2 + (0-3.5)^2 + (1.5-1)^2 + (1-0)^2} = 3.674$$

$\therefore$  assignment = C2 ( $\because C_2 < C_1, C_2 < C_3$ )

$$\text{new } C_2 = \left( \frac{2.25+2}{2}, \frac{1.75+0}{2}, \frac{0.5+0.5}{2}, \frac{1+3}{2} \right)$$

$$= (2.125, 0.875, 1, 2)$$

For  $P_3 = (1, 3, 4.5, 5)$

$$ED_{P3-C1} = \sqrt{(1-3)^2 + (3-5)^2 + (4.5-4)^2 + (5-4.5)^2} = 2.915$$

$$ED_{P3-C2} = \sqrt{(1-1.5)^2 + (3-0.875)^2 + (4.5-1)^2 + (5-2)^2} = 5.20$$

$$ED_{P3-C3} = \sqrt{(1-2)^2 + (3-3.5)^2 + (4.5-1)^2 + (5-0)^2} = 6.20$$

assignment = C1 ( $\because C_1 < C_2, C_1 < C_3$ )

$$\text{new } C_1 = \left( \frac{3+1}{2}, \frac{5+3}{2}, \frac{4+4.5}{2}, \frac{4.5+5}{2} \right) = (2, 4, 4.25, 4.75)$$

$$\text{then } C_1 = \left( \frac{2.5+2}{2}, \frac{2.5+4}{2}, \frac{5+4.25}{2}, \frac{3.5+4.75}{2} \right) \\ = (2.25, 3.25, 4.625, 4.125)$$



$$\text{For } P_4 = (3.5, 2.5, 5, 3.5)$$

$$ED_{P_4-C_1} = \sqrt{(3.5-2)^2 + (2.5-4)^2 + (5-4.75)^2 + (3.5-4.75)^2} \\ = 2.150$$

$$ED_{P_4-C_2} = \sqrt{(3.5-2.125)^2 + (2.5-0.875)^2 + (5-1)^2 + (3.5-2)^2} \\ = 4.586$$

$$ED_{P_4-C_3} = \sqrt{(3.5-2)^2 + (2.5-3.5)^2 + (5-1)^2 + (3.5-0)^2} \\ = 5.431$$

assignment = C1 ( $\because C_1 < C_2, C_1 < C_3$ )

$$\text{For } P_5 = (6, 3, 2, 0)$$

$$ED_{P_5-C_1} = \sqrt{(6-2.25)^2 + (3-3.25)^2 + (2-4.25)^2 + (0-4.125)^2} \\ = 6.167$$

$$ED_{P_5-C_2} = \sqrt{(6-2.125)^2 + (3-0.875)^2 + (2-1)^2 + (0-2)^2} \\ = 4.952$$

$$ED_{P_5-C_3} = \sqrt{(6-2)^2 + (3-3.5)^2 + (2-1)^2 + (0-0)^2} \\ = 4.153$$

$\therefore$  assignment = C3 ( $\because C_3 < C_1, C_3 < C_2$ )

$$N = 1000$$

$$k_1 = 1.25$$

$$b = 0.8$$

$$\text{avgd}1 = 50$$

BH25

Date \_\_\_\_\_

Page \_\_\_\_\_

	+1	+2	+3	+4	+5	+6	length
d1	3	0	8	1	10	5	21
d2	4	3	3	2	1	7	20
C	100	50	80	93	100	95	1000

$$BH25(q, d) = \sum_{t \in q} \frac{f_{t,d} \cdot (1 + k_1)}{f_{t,d} + k_1(1 - b + b \cdot \left(\frac{1}{\text{avgd}1}\right))} \cdot \log_{10} \frac{N}{n_t}$$

q = "t2"

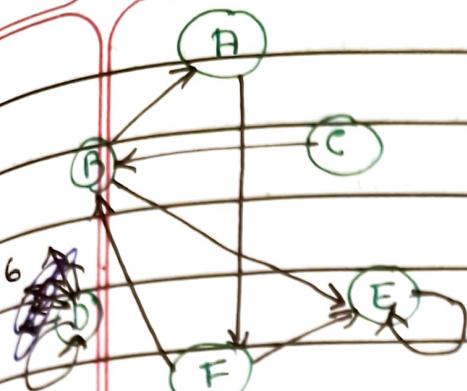
$$\textcircled{1} \quad BH25(q, d1) = 0$$

$$\textcircled{2} \quad BH25(q, d2) = \frac{3 \cdot (1 + 1.25) \cdot \log_{10} \left( \frac{1000}{50} \right)}{3 + 1.25(1 - 0.8 + 0.8 \left( \frac{20}{50} \right))} = \frac{8.781952}{3.65} \\ = 2.40601$$

q = "t2 t2 +5"

$$\textcircled{3} \quad BH25(q, d1) = 0 + \frac{10 \cdot (1 + 1.25) \cdot \log_{10} \left( \frac{1000}{100} \right)}{10 + 1.25(1 - 0.8 + 0.8 \left( \frac{91}{50} \right))} = \frac{23.5}{10.67} \\ = 2.1087$$

$$\textcircled{4} \quad BH25(q, d2) = 2.40601 + 2.40601 + \frac{1 \cdot (1 + 1.25) \cdot \log_{10} \left( \frac{1000}{100} \right)}{1 + 1.25(1 - 0.8 + 0.8 \left( \frac{90}{50} \right))} \\ = 6.175$$



	in	out	
A	1	1	B/2
B	2	2	C + E/2
C	0	1	0
D	1	1	D/2
E	3	1	B/2 + F/2 + D/2
F	1	2	A/1

Q      1  
A  $0.167$        $(\frac{0.2}{6}) + 0.8(\frac{0.167}{2}) = 0.1$

B  $0.167$        $(\frac{0.2}{6}) + 0.8(0.167 + \frac{0.167}{2}) = 0.233$

C  $0.167$        $(\frac{0.2}{6}) + 0.8(0) = 0.033$

D  $0.167$        $(\frac{0.2}{6}) + 0.8(\frac{0.167}{2}) = 0.167$

E  $0.167$        $(\frac{0.2}{6}) + 0.8(\frac{0.167}{2} + \frac{0.167}{2} + 0.167) = 0.3$

F  $0.167$        $(\frac{0.2}{6}) + 0.8(0.167) = 0.167$

2  
A  $(\frac{0.2}{6}) + 0.8(\frac{0.233}{2}) = 0.127$

B  $(\frac{0.2}{6}) + 0.8(0.033 + \frac{0.167}{2}) = 0.127$

C  $(\frac{0.2}{6}) + 0.8(0) = 0.033$

D  $(\frac{0.2}{6}) + 0.8(0.167) = 0.167$

E  $(\frac{0.2}{6}) + 0.8(\frac{0.233}{2} + \frac{0.167}{2} + 0.167) = 0.327$

F  $(\frac{0.2}{6}) + 0.8(0.1) = 0.113$

? 0.433

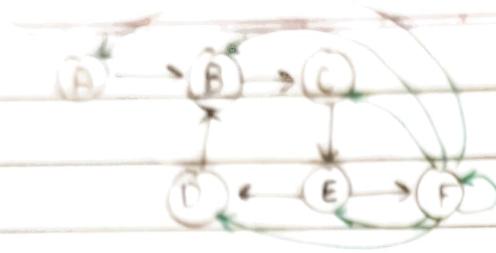
## PageRank

$$q = 0.4$$

$$T = 6$$

$$1 - q = 0.6$$

Bounce  
Page



F does not have outgoing link so sink sink need to be fixed.

<u>out</u>	$P(i_{in}(out))$	$i=0$	$i=1$
A	$F/6$	$\frac{1}{6} = 0.167$	$\left(\frac{0.4}{6}\right) + 0.6\left(\frac{0.167}{6}\right) = 0.083$
B	$A + \frac{F}{6} + D$	$0.167$	$\left(\frac{0.4}{6}\right) + 0.6\left(0.167 + \frac{0.167}{6} + 0.167\right) = 0.283$
C	$B + \frac{E}{6}$	$0.167$	$\left(\frac{0.4}{6}\right) + 0.6\left(0.167 + \frac{0.167}{6}\right) = 0.183$
D	$E/2 + \frac{F}{6}$	$0.167$	$\left(\frac{0.4}{6}\right) + 0.6\left(0.167\right) + \frac{0.167}{2} = 0.133$
E	$C + \frac{F}{6}$	$0.167$	$\left(\frac{0.4}{6}\right) + 0.6\left(0.167 + \frac{0.167}{6}\right) = 0.183$
F	$F/6 + E/2$	$0.167$	$\left(\frac{0.4}{6}\right) + 0.6\left(0.167\right) + \frac{0.167}{2} = 0.133$

$i=2$

- A  $\left(\frac{0.4}{6}\right) + 0.6\left(\frac{0.133}{6}\right) = 0.079$
- B  $\left(\frac{0.4}{6}\right) + 0.6\left(0.083 + \frac{0.133}{6} + 0.133\right) = 0.209$
- C  $\left(\frac{0.4}{6}\right) + 0.6\left(0.283 + \frac{0.133}{6}\right) = 0.249$
- D  $\left(\frac{0.4}{6}\right) + 0.6\left(0.183 + \frac{0.133}{6}\right) = 0.135$
- E  $\left(\frac{0.4}{6}\right) + 0.6\left(0.183 + \frac{0.133}{6}\right) = 0.189$
- F  $\left(\frac{0.4}{6}\right) + 0.6\left(0.133 + \frac{0.183}{6}\right) = 0.135$

Instrument True Class

		A	B	
1	N	Y	N	-
2	N	N	Y	-
3	Y	Y	Y	-
4	N	Y	N	+ TP = 3
5	Y	N	Y	+ FN = 1
6	Y	Y	Y	-
7	N	Y	N	-
8	Y	Y	N	+ 1 3

$$\underline{\underline{A}} \quad P = \frac{TP}{TP+FP} = \frac{3}{3+3} = \frac{3}{6} = 0.5$$

$$R = \frac{TP}{FN+TP} = \frac{3}{1+3} = \frac{0.75}{\text{same}}$$

$$F1 = \frac{2PR}{P+R} = \frac{2 \cdot (0.5)(0.75)}{0.5+0.75} = 0.6$$

$$\underline{\underline{B}} \quad P = \frac{3}{1+3} = \boxed{0.75}$$

$$R = \frac{3}{4} = 0.75 \quad \text{same}$$

$$F1 = \frac{2(0.75)(0.75)}{0.75+0.75} = \boxed{0.75} \quad 0.75$$

actual predicted

		Predicted	
		+	-
		actual	
	N	Y	
	Y	Y	
	N	Y	+ 3 2
	N	Y	- 3 2
	N	N	
	Y	Y	
	Y	Y	
	N	N	
	Y	N	

# Classification Evaluation

Date \_\_\_\_\_  
Page \_\_\_\_\_

Predicted.

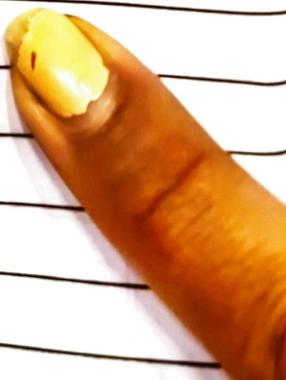
	N	Y	-	+
-	Y	Y	TN = 3	FP = 3
+	Y	Y	FN = 1	TP = 3

Actual

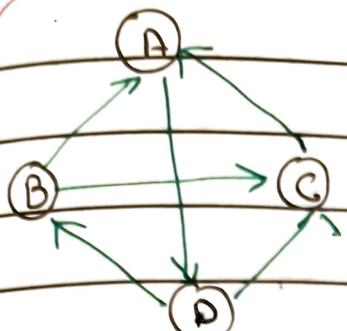
$$P = \frac{TP}{TP + FP} = \frac{3}{3+3} = 0.5$$

$$R = \frac{TP}{FN + TP} = \frac{3}{1+3} = 0.75$$

$$F1 = \frac{2 \cdot PR}{P + R} = \frac{2(0.5)(0.75)}{0.5 + 0.75} = 0.6.$$



# PageRank



$$\alpha = 0.2$$

$$T = 4$$

$$1 - \alpha = 0.8$$

P(i<sub>m</sub>)

$$A \quad B/2 + C/1$$

0

1

2

$$B \quad D/2$$

$$0.25$$

$$0.15$$

$$0.31$$

$$C \quad B/2 + D/2$$

$$0.25$$

$$0.25$$

$$0.21$$

$$D \quad A$$

$$0.25$$

$$0.25$$

$$0.33$$

$$(0.2) + 0.8(0.25 + 0.25) = 0.35$$

$$(0.2) + 0.8(0.15 + 0.25) = 0.15$$

$$(0.2) + 0.8(0.25 + 0.25) = 0.35$$

$$(0.2) + 0.8(0.35) = 0.25$$

# Language Modeling Retrieval

## Disjunctive smoothing



	$d_1$	$d_2$	$d_3$	$d_4$	$ C $	$P(t, C)$
$t_1$	1	1	2	1	5	$\frac{5}{22} = 0.227$
$t_2$	0	2	0	1	3	$\frac{3}{22} = 0.136$
$t_3$	2	0	1	0	3	$\frac{3}{22} = 0.136$
$t_4$	4	0	1	2	7	$\frac{7}{22} = 0.318$
$t_5$	1	2	1	0	4	$\frac{4}{22} = 0.182$
	8	5	5	4	22	

① Empirical- $P(t_5 | d_1) = \frac{C_{t_5, d_1}}{|d_1|} = \frac{1}{8} = 0.125$

② background- $P(t_4 | G) = \frac{\sum_{d'} C_{t_4, d'}}{\sum_{d'} |d'|} = \frac{7}{22} = 0.318$

③ smoothed- $P(t_2 | d_3) = \frac{C_{t_2, d_3} + \lambda P(t_2 | C)}{|d_3| + \lambda} = \frac{0 + 6(\frac{3}{22})}{5 + 6} = 0.074$

$$P(t_1 | d_2) = \frac{1 + 6(\frac{5}{22})}{5 + 6} = 0.214 \quad P(t_4 | d_2) = \frac{0 + 6(\frac{3}{22})}{5 + 6} = 0.173$$

④  $P(t_2 | d_1) = \frac{2 + 6(\frac{3}{22})}{5 + 6} = 0.356 \quad P(t_5 | d_2) = \frac{8 + 6(\frac{4}{22})}{5 + 6} = 0.380$

$$P(t_3 | d_2) = \frac{0 + 6(\frac{5}{22})}{5 + 6} = [0.074] \text{ lowest}$$

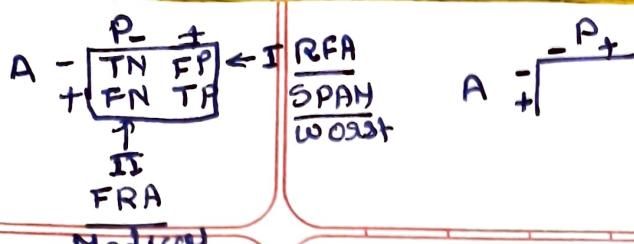
⑤  $Q = "t_1 t_3 t_5"$

$$\begin{aligned} \text{Score}(d_1, Q) &= \log\left(\frac{1 + 6(\frac{5}{22})}{5 + 6}\right) \cdot 1 + \log\left(\frac{2 + 6(\frac{3}{22})}{5 + 6}\right) \cdot 1 + \log\left(\frac{1 + 6(\frac{4}{22})}{5 + 6}\right) \cdot 1 \\ &= -2.2944 \end{aligned}$$

$$\begin{aligned} \text{Score}(d_2, Q) &= \log\left(\frac{1 + 6(\frac{3}{22})}{5 + 6}\right) \cdot 1 + \log\left(\frac{0 + 6(\frac{3}{22})}{5 + 6}\right) \cdot 1 + \log\left(\frac{2 + 6(\frac{4}{22})}{5 + 6}\right) \cdot 1 \\ &= -2.3476 \end{aligned}$$

$$\begin{aligned} \text{Score}(d_3, Q) &= \log\left(\frac{2 + 6(\frac{5}{22})}{5 + 6}\right) \cdot 1 + \log\left(\frac{1 + 6(\frac{3}{22})}{5 + 6}\right) \cdot 1 + \log\left(\frac{1 + 6(\frac{4}{22})}{5 + 6}\right) \cdot 1 \\ &= -2.0173 \leftarrow \text{highest} \end{aligned}$$

$$\begin{aligned} \text{Score}(d_4, Q) &= \log\left(\frac{1 + 6(\frac{5}{22})}{4 + 6}\right) \cdot 1 + \log\left(\frac{0 + 6(\frac{3}{22})}{4 + 6}\right) \cdot 1 + \log\left(\frac{0 + 6(\frac{4}{22})}{4 + 6}\right) \cdot 1 \\ &= -2.6757 \end{aligned}$$



Date \_\_\_\_\_  
Page \_\_\_\_\_

Medical  
Worst.

$$ACC = \frac{TP + TN}{TP + TN + FN + FP} = \text{F of CC items}$$

$$ER = \frac{FN + FP}{TP + FN + TP + FP} = \text{F of ICC items.}$$

$$P = \frac{TP}{TP + FP} = \text{F of CC items} \rightarrow \text{pos from total pos predicted.}$$

$$R = \frac{TP}{TP + FN} = \text{F of CC items} \rightarrow \text{pos from actual pos.}$$

sensitivity  
True Pos Rate

$$\text{Specificity} = \frac{TN}{FP + TN}$$

$$F1 = \frac{2 \cdot PR}{PR} = \text{mean of P \& R.}$$

$$FPR = E \cdot I = \frac{FP}{FP + TN} = \text{F of WC items - pos out of actual - neg}$$

$$FNR = F-II = \frac{FN}{FN + TP} = \text{F of WC items - neg out of actual pos}$$

$$P(t|C) = \frac{n(t, C)}{\sum_{t'} n(t', C) + |C|}$$

NB

Date \_\_\_\_\_  
Page \_\_\_\_\_

	$t_1$	$t_2$	$t_3$	$t_4$	$t_5$	$t_6$	$t_7$	$t_8$	class
$d_1$	1	9	7	5	2	1	3	0	C1
$d_2$	0	2	1	3	0	1	0	2	C2
$d_3$	1	1	0	0	4	5	2	1	C1
$d_4$	2	1	4	0	5	2	7	0	C3
$d_5$	0	1	5	0	1	4	0	2	C2
$d_6$	6	3	0	3	2	1	3	0	C3
$d_7$	2	1	3	1	8	1	0	9	C1
$d_8$	1	1	0	1	0	9	0	0	C3

$$\frac{n(t, C)}{\text{Total } P(C)}$$

$n(t, C)$	$\{ C_1 \}$	4	11	10	6	14	7	5	10	67	$31/8 = 0.375$
	$\{ C_2 \}$	0	3	6	3	1	5	0	4	22	$2/8 = 0.25$
	$\{ C_3 \}$	9	5	4	4	7	12	10	0	51	$31/8 = 0.375$

$p(t, C)$	$\{ C_1 \}$	$\frac{4+1}{67+3} = 0.071$	$0.171$	$0.157$	$0.1$	$0.213$	$0.114$	$0.086$	$0.157$
	$\{ C_2 \}$	$\frac{1}{22} = 0.045$	$0.16$	$0.28$	$0.16$	$0.08$	$0.24$	$0.04$	$0.2$
	$\{ C_3 \}$	$\frac{9+1}{51+3} = 0.186$	$0.11$	$0.093$	$0.093$	$0.148$	$0.24$	$0.203$	$0.018$

$$P(C_2 | "t_2 + 3 + t_6")$$

$$= P(C_2) \cdot P("t_2" | C_2) \cdot P("t_6" | C_2) \cdot P("t_3" | C_2)$$

$$= 0.25 \times 0.16 \times 0.28 \times 0.24$$

$$= 0.0087$$

$$P(C_1 | "t_2 + 7 + t_8") \text{ is the highest? } C_1$$

$$P(C_1 | "t_2 + 7 + t_8") = 0.375 \times 0.171 \times 0.086 \times 0.157 = 0.00011$$

$$P(C_2 | "t_2 + 7 + t_8") = 0.25 \times 0.16 \times 0.04 \times 0.2 = 0.00032$$

$$P(C_3 | "t_2 + 7 + t_8") = 0.375 \times 0.11 \times 0.203 \times 0.018 = 0.000015$$

$$P(t|c) = \frac{n(t,c)}{\sum_i n(t^i, c) + 1(c)} \quad \underline{NB}$$

Date \_\_\_\_\_  
Page \_\_\_\_\_

	$t_1$	$t_2$	$t_3$	$t_4$	$t_5$	$t_6$	$t_7$	Class
$d_1$	2	0	1	2	0	2	4	$C_1$
$d_2$	0	0	0	0	3	2	2	$C_3$
$d_3$	3	4	0	2	0	0	2	$C_2$
$d_4$	4	0	3	1	1	1	0	$C_3$
$d_5$	1	0	0	3	1	2	0	$C_2$
$d_6$	0	1	1	0	3	4	1	$C_1$

	$n(t^i, c)$	$P(c)$
$c_1$	8	$\frac{8}{21} = 0.38$
$c_2$	4	$\frac{4}{21} = 0.19$
$c_3$	4	$\frac{4}{21} = 0.19$
	10	56

Smoothed $P(c t)$		$c_1$	$c_2$	$c_3$
$c_1$	$\frac{3}{21} = 0.143$	0.433	0.175	0.125
$c_2$	$\frac{5}{21} = 0.238$	0.238	0	$0.386$
$c_3$	$\frac{5}{20} = 0.25$	0	0.2	0.1

$$\textcircled{1} \quad P(c_2) = 0.333$$

$$\textcircled{2} \quad P("t_4" | c_2) = 0.285$$

$$\textcircled{3} \quad P(c_1 | "t_1") = P(c_1) \cdot P("t_1" | c_1) = 0.33 \times 0.143 = 0.041$$

$$\begin{aligned} \textcircled{4} \quad P(c_3 | "t_1 t_4 t_5") &= P(c_3) \cdot P("t_1" | c_3) \cdot P("t_4" | c_3) \cdot P("t_5" | c_3) \\ &= 0.33 \times 0.25 \times 0.1 \times 0.25 \\ &= 0.0020625 \end{aligned}$$

$$\textcircled{5} \quad P(c_1 | "t_4 t_4 t_5") = 0.33 \times 0.175 \times 0.175 \times 0.167 = 0.000861$$

$$P(c_2 | "t_4 t_4 t_5") = 0.33 \times 0.285 \times 0.285 \times 0.095 = 0.002546$$

$$P(c_3 | "t_4 t_4 t_5") = 0.33 \times 0.1 \times 0.1 \times 0.25 = 0.000825$$

# Weighted TF-IDF

Date \_\_\_\_\_  
Page \_\_\_\_\_

N=4

	$t_1$	$t_2$	$t_3$	$t_4$	$t_5$	$t_d$
$d_1$	0	5	6	10	1	12
$d_2$	1	6	3	13	5	18
$d_3$	7	0	3	18	0	18
$d_4$	4	4	9	0	2	19
$N_t$	12	15	21	11	8	21
$n_t$	3	3	4	2	3	3

$$w_{t,d} = \frac{f_{t,d}}{\sum_{t \neq d} f_{t,d}} \cdot \log_2 \frac{N+1}{n_t}$$

$w_{d_1}$	$\frac{0}{12} \cdot \log_2 \left( \frac{4+1}{3} \right)$	$\frac{5}{12} \cdot \log_2 \left( \frac{5}{3} \right)$	$\frac{6}{12} \cdot \log_2 \left( \frac{5}{4} \right)$	0	$\frac{1}{12} \cdot \log_2 \left( \frac{5}{3} \right)$
$w_{d_2}$	$\frac{1}{18} \cdot \log_2 \left( \frac{5}{3} \right)$	$\frac{6}{18} \cdot \log_2 \left( \frac{5}{3} \right)$	$\frac{3}{18} \cdot \log_2 \left( \frac{5}{4} \right)$	$\frac{3}{18} \cdot \log_2 \left( \frac{5}{3} \right)$	$\frac{5}{18} \cdot \log_2 \left( \frac{5}{3} \right)$
$w_{d_3}$	$\frac{7}{18} \cdot \log_2 \left( \frac{5}{3} \right)$	0	$\frac{3}{18} \cdot \log_2 \left( \frac{5}{4} \right)$	$\frac{8}{18} \cdot \log_2 \left( \frac{5}{3} \right)$	0
$w_{d_4}$	$\frac{4}{19} \cdot \log_2 \left( \frac{5}{3} \right)$	$\frac{4}{19} \cdot \log_2 \left( \frac{5}{3} \right)$	$\frac{9}{19} \cdot \log_2 \left( \frac{5}{4} \right)$	0	$\frac{2}{19} \cdot \log_2 \left( \frac{5}{3} \right)$

$w_{d_1}$	0	0.307	0.161	0	0.061
$w_{d_2}$	0.041	0.246	0.054	0.220	0.205
$w_{d_3}$	0.287	0	0.054	0.587	0
$w_{d_4}$	0.155	0.155	0.182	0	0.078

$$\text{Cosine Similarity} = \frac{\sum_t w_{t,d} \cdot w_{t,d'}}{\sqrt{\sum_t (w_{t,d})^2} \cdot \sqrt{\sum_t (w_{t,d'})^2}} = \frac{d \cdot d'}{\|d\| \cdot \|d'\|}$$

$$\text{Sim}(d_2, d_4) = \frac{[0.041, 0.246, 0.054, 0.220, 0.205] \cdot [0.155, 0.155, 0.155, 0, 0.078]}{\sqrt{0.041^2 + 0.246^2 + 0.054^2 + 0.220^2 + 0.205^2} \cdot \sqrt{0.155^2 + 0.155^2 + 0.155^2 + 0 + 0.078^2}} = 0.672$$

$$\text{Sim}(w_{d_2}, w_{d_4}) = \frac{[0.041, 0.246, 0.054, 0.220, 0.205] \cdot [0.155, 0.155, 0.155, 0, 0.078]}{\sqrt{0.041^2 + 0.246^2 + 0.054^2 + 0.220^2 + 0.205^2} \cdot \sqrt{0.155^2 + 0.155^2 + 0.155^2 + 0 + 0.078^2}} = 0.627$$

# BM 25

Date \_\_\_\_\_

Page \_\_\_\_\_

	$t_1$	$t_2$	$t_3$	$t_4$	$t_5$	$t_6$	length
$d_1$	3	0	15	2	10	5	25
$d_2$	4	3	5	2	1	5	20
C	100	50	80	93	100	25	1000

$$\text{avgdl} = 50$$

$$K_1 = 1.2$$

$$b = 0.75$$

$$BM25_{(q, d)} = \sum_{t \in q} \frac{f_{t,d} \cdot (1 + K_1)}{f_{t,d} + K_1(1 - b + b \cdot \frac{\text{avgdl}}{f_{t,d}})} \cdot \log \frac{N}{n_t}$$

$q$  = "terme"

$$BM25_{(q, d_1)} = 0$$

$$BM25_{(q, d_2)} = \frac{3 \cdot (1 + 1.2)}{3 + 1.2(1 - 0.75 + 0.75(\frac{20}{50}))} \cdot \log_{10} \left( \frac{1000}{50} \right)$$

$$= 2.346$$

$q$  = "terme terme terms"

$$BM25_{(q, d_1)} = 0 + 0 + \frac{10 \cdot (1 + 1.2) \cdot \log_{10} \left( \frac{1000}{100} \right)}{10 + 1.2(1 - 0.75 + 0.75(\frac{25}{50}))} = 2.046$$

$$BM25_{(q, d_2)} = 2.346 + 2.346 + \frac{1 \cdot (1 + 1.2) \cdot \log_{10} \left( \frac{1000}{100} \right)}{1 + 1.2(1 - 0.75 + 0.75(\frac{20}{50}))} = 6.01750$$