**11-442 / 11-642 / 11-742:**
**Search Engines**

**Best-Match Retrieval:**
**HW2 Implementation**

Jamie Callan

Carnegie Mellon University

callan@cs.cmu.edu

---

# Outline

**HW2 implementation**
- Indri default beliefs
- Window operator

## Indri Implementation

**Indri operators calculate scores <u>only for documents that contain a query term</u>**

- Use inverted or score lists
  - Similar to HW1
- Use document length, ctf, and corpus length for smoothing
  - Lookup from the index – see the HW2 web page

**But, one aspect is a little tricky to get right…**

$$\text{AND, COMBINE}: \quad p_{and}(q\,|\,d) \;=\; \prod_{q_i \in q} p(q_i\,|\,d)^{\frac{1}{|q|}}$$

$$\text{WAND, WEIGHT}: \quad p_{wand}(q\,|\,d) \;=\; \prod_{q_i \in q} p(q_i\,|\,d)^{\frac{w_i}{w}}, \quad w = \sum w_i$$

$$\text{OR}: \quad p_{or}(q\,|\,d) \;=\; 1 - \prod_{q_i \in q}\left(1 - p(q_i\,|\,d)\right)$$

$$\text{WSUM}: \quad p_{wsum}(q\,|\,d) \;=\; \sum_{q_i \in q}\frac{w_i}{w}p(q_i\,|\,d), \quad w = \sum w_i$$

$$\text{NOT}: \quad p_{not}(q\,|\,d) \;=\; 1 - p(q\,|\,d)$$

---

## Indri Review

**Your HW2 system will use two-stage smoothing to compute term weights**

$$p(q_i|d) = (1 - \lambda)\,\frac{tf_{q_i,d} + \mu\, p_{MLE}(q_i|C)}{length_d + \mu} + \lambda\, p_{MLE}(q_i|C)$$

$$p_{MLE}(q_i|C) = \frac{ctf_{q_i}}{length_C}$$

This is the #SCORE operator for the Indri retrieval model

Page 2

**Indri Implementation**

**Query:** #OR (a #AND (b c) )          **Document:** a

**Query terms b and c do not appear in this document**
**… what is the score of the #AND operator?**

- $tf_b = 0$          $tf_c = 0$
- Do the usual Indri score calculation

$$p(q_i|d) = (1 - \lambda)\frac{tf_{q_i,d} + \mu\, p_{MLE}(q_i|C)}{length_d + \mu} + \lambda\, p_{MLE}(q_i|C)$$

  – MLE p(t | d) scores are 0 for b and c, so only <u>smoothing</u> scores for b and c

**This is simple conceptually, but <u>how is it implemented</u>?**

- You don't want to calculate #AND scores for <u>every</u> document (<u>way</u> too slow)
  … just the documents that have at least one query term

© 2021 Jamie Callan

---

**Indri Implementation**

**Query:** #OR (a #AND (b c) )          **Document:** a

**Add a <u>new method</u> to <u>all</u> QrySop operators**
  double getDefaultScore (RetrievalModel r, long docid)

**When <u>any</u> QrySop operator calculates scores for the <u>Indri retrieval model</u>**
  if $q_i$ has a match for document d
  then
    call $q_i$.getScore
  else
    call $q_i$.getDefaultScore

© 2021 Jamie Callan

# Indri Implementation

**Query:** #OR (a #AND (b c) )        **Document:** a

**QrySopScore.getDefaultScore (RetrievalModel r, long docid)**

• The standard Indri <u>SCORE</u> calculation done with tf=0

   If r == RetrievalModel.Indri

$$p_{scoreDefault}(t \mid \text{docid}) = (1 - \lambda) \frac{0 + \mu \, p_{MLE}(t \mid C)}{length(\text{docid}) + \mu} + \lambda \, p_{MLE}(t \mid C)$$

**<span style="color:red">This is the main difference.</span>**
**<span style="color:red">Do the usual calculation, but with tf=0.</span>**

---

# Indri Implementation

**Query:** #AND (a #NEAR/3 (b c) )

**QrySopScore.getDefaultScore (RetrievalModel r, long docid)**

**What happens if #NEAR/3 (b c) doesn't occur in the collection?**

• Its ctf == 0

   … so $p_{MLE}(t|C) == 0$      (i.e., no smoothing weight)

   … so #AND returns 0 for all documents     $p_{and}(q \mid d) = \prod_{q_i \in q} p(q_i \mid d)^{\frac{1}{|q|}}$

• This behavior is exact match, not best match

   – Indri is a best match model

**Indri Implementation**

**Query:** #AND (a #NEAR/3 (b c) )

**QrySopScore.getDefaultScore (RetrievalModel r, long docid)**

**What happens if #near/3 (b c) doesn't occur in the collection?**

**Solution:** <u>Extra smoothing</u> for terms that have ctf = 0

  If r == RetrievalModel.Indri

    If ctf $(t) = 0$

      calculate $p_{MLE}(t|C)$ using ctf $(t) = 0.5$

**Undocumented behavior (ctf < 1)**

$$p_{scoreDefault}(t \mid \text{docid}) = (1-\lambda)\frac{0 + \mu \, p_{MLE}(t \mid C)}{length(\text{docid}) + \mu} + \lambda \, p_{MLE}(t \mid C)$$

---

**Indri Implementation**

**Query:** #OR (a #AND (b c) )        **Document:** a

**QrySopAnd.getDefaultScore (RetrievalModel r, long docid)**

• The standard Indri <u>AND</u> calculation done on the <u>default score</u> of each argument

  If r == RetrievalModel.Indri

$$p_{andDefault}(q \mid d) = \prod_{q_i \in q} \left( p_{q_i default}(q_i \mid d) \right)^{\frac{1}{|q|}}$$

**This is the only difference.**
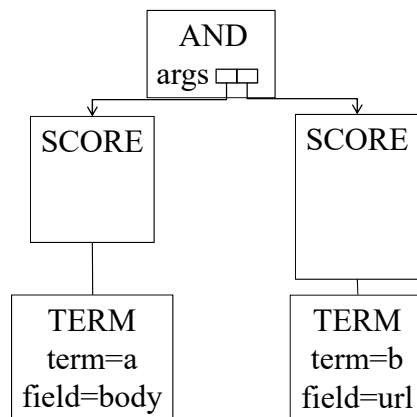**Call the i<sup>th</sup> query argument's getDefaultScore method.**

Page 5

# QryEval Example

---

# The initial query:
# #AND (a.body b.url)

```
                    +-------------+
                    |    AND      |
                    | args [  ][  ] |
                    +-------------+
               /                    \
    +-------------+            +-------------+
    |   SCORE     |            |   SCORE     |
    |             |            |             |
    +-------------+            +-------------+
         |                          |
    +-------------+            +-------------+
    |   TERM      |            |   TERM      |
    |  term=a     |            |  term=b     |
    |  field=body |            |  field=url  |
    +-------------+            +-------------+
```

## Query Initialization

```
        AND
        args ▭▭
       ↙        ↘
   SCORE        SCORE
     │            │
┌──────────────┐  ┌──────────┐
│ field=body   │  │ TERM     │
│ df=4,  ctf=8 │  │ term=b   │
│ ┌──────────┐ │  │ field=url│
│ │docid=3, tf=2│ └──────────┘
│ └──────────┘ │
│ docid=18, tf=4│
│  :   :   : : │
└──────────────┘
```
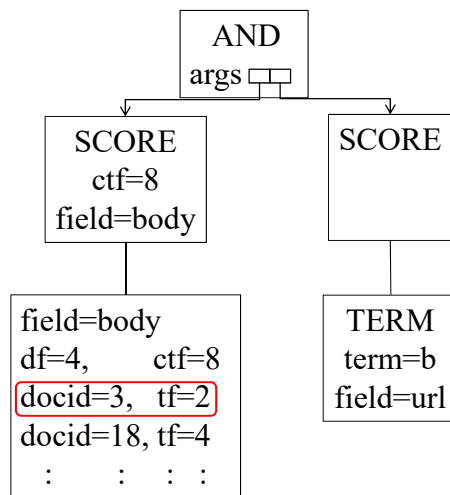
**AND operator initialization**
AND initializes its first arg.
SCORE initializes its arg.
The result is an inverted list.

13

© 2021 Jamie Callan

## Query Initialization

```
        AND
        args ▭▭
       ↙        ↘
   SCORE        SCORE
   ctf=8          │
  field=body      │
     │            │
┌──────────────┐  ┌──────────┐
│ field=body   │  │ TERM     │
│ df=4,  ctf=8 │  │ term=b   │
│ ┌──────────┐ │  │ field=url│
│ │docid=3, tf=2│ └──────────┘
│ └──────────┘ │
│ docid=18, tf=4│
│  :   :   : : │
└──────────────┘
```
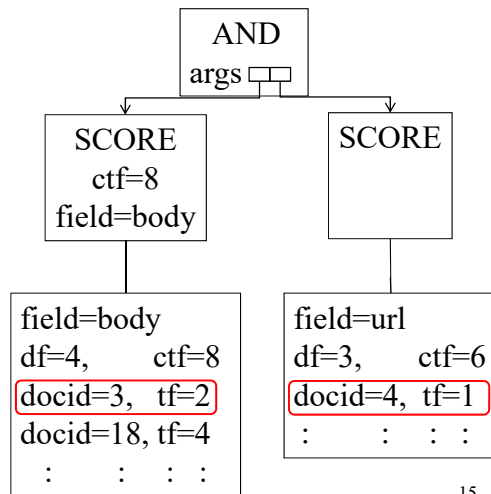
**AND operator initialization**
AND initializes its first arg.
SCORE initializes its arg.
The result is an inverted list.

The SCORE operator
caches information that it
will need later.
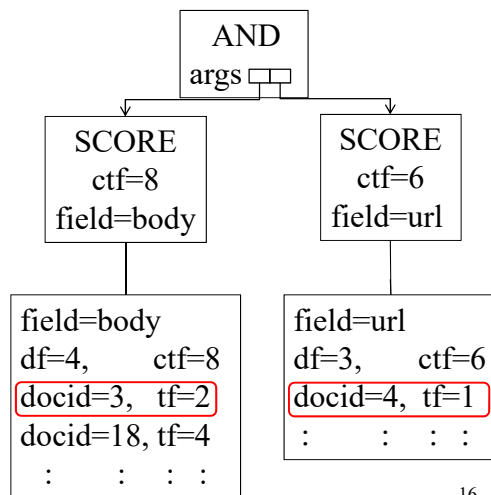
14

© 2021 Jamie Callan

Page 7

## Query Initialization

```
         AND
         args ⎡⎤⎡⎤
        ↙        ↓
   SCORE        SCORE
   ctf=8
   field=body
      |            |
┌────────────┐  ┌────────────┐
│ field=body │  │ field=url  │
│ df=4,  ctf=8│ │ df=3,   ctf=6│
│ docid=3, tf=2│ │ docid=4, tf=1│
│ docid=18, tf=4│ │ :    :  : : │
│ :    :  : :  │ └────────────┘
└────────────┘
```

**AND operator initialization**
AND initializes its second arg.
SCORE initializes its arg.
The result is an inverted list.

15

---

## Query Initialization

```
         AND
         args ⎡⎤⎡⎤
        ↙        ↓
   SCORE        SCORE
   ctf=8         ctf=6
   field=body    field=url
      |            |
┌────────────┐  ┌────────────┐
│ field=body │  │ field=url  │
│ df=4,  ctf=8│ │ df=3,   ctf=6│
│ docid=3, tf=2│ │ docid=4, tf=1│
│ docid=18, tf=4│ │ :    :  : : │
│ :    :  : :  │ └────────────┘
└────────────┘
```
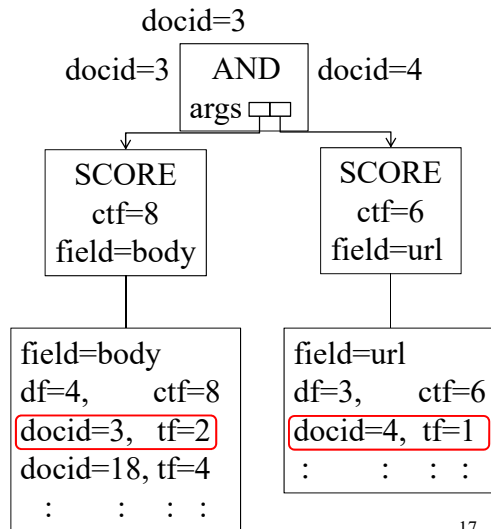
**AND operator initialization**
AND initializes its second arg.
SCORE initializes its arg.
The result is an inverted list.

The SCORE operator
caches information that it
will need later.

16

Page 8

## Call to docIteratorHasMatch & getScore
### (First Time)

docid=3

docid=3 | AND | docid=4
args

SCORE
ctf=8
field=body

SCORE
ctf=6
field=url

field=body
df=4,       ctf=8
docid=3,   tf=2
docid=18, tf=4
:      :    :   :

field=url
df=3,       ctf=6
docid=4,  tf=1
:      :    :   :

**AND Operator Evaluation**
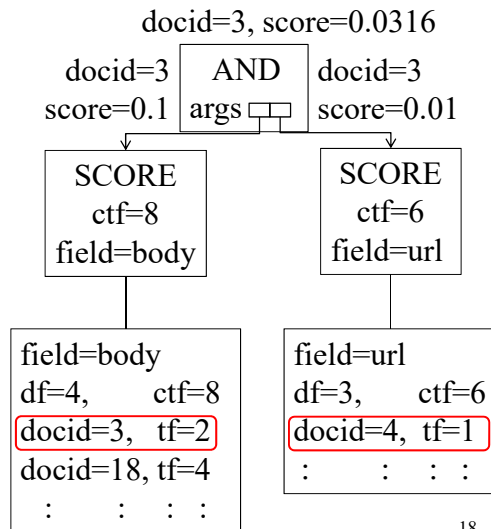
Min document is 3.

HasMatch caches the docid (3)
and returns True

17

---

## Call to docIteratorHasMatch & getScore
### (First Time)

docid=3, score=0.0316

docid=3 | AND | docid=3
score=0.1 | args | score=0.01

SCORE
ctf=8
field=body

SCORE
ctf=6
field=url

field=body
df=4,       ctf=8
docid=3,   tf=2
docid=18, tf=4
:      :    :   :

field=url
df=3,       ctf=6
docid=4,  tf=1
:      :    :   :

**AND Operator Evaluation**

getScore uses the cached docid (3).

args[0] has match, so call
  args[0].getScore ().
Suppose the result is 0.1.
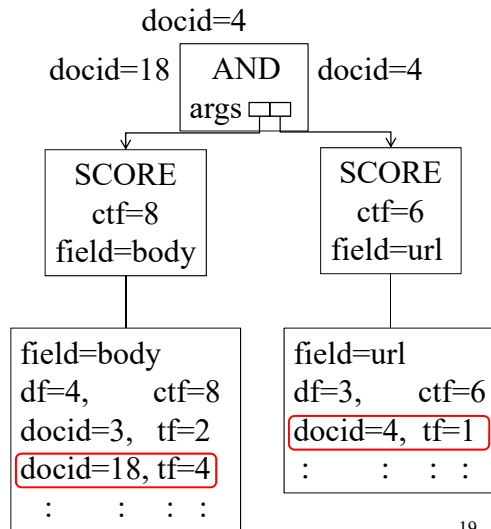
args[1] does not match, so call
  args[1].getDefaultScore(3).
Suppose the result is 0.01.

$Score_{AND}(3) = (0.1^{0.5} \times 0.01^{0.5})$

18

Page 9

## Call to docIteratorHasMatch & getScore
### (Second Time)

docid=4

docid=18 | AND | docid=4
args

SCORE
ctf=8
field=body

SCORE
ctf=6
field=url

field=body
df=4,      ctf=8
docid=3,   tf=2
docid=18, tf=4
:      :    :  :

field=url
df=3,      ctf=6
docid=4,  tf=1
:      :    :  :

**AND Operator Evaluation**
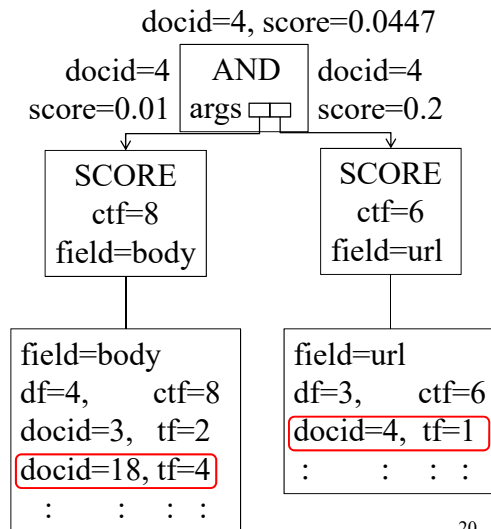Min document is 4.

HasMatch caches the docid (4)
and returns True

19

---

## Call to docIteratorHasMatch & getScore
### (Second Time)

docid=4, score=0.0447

docid=4 | AND | docid=4
score=0.01 | args | score=0.2

SCORE
ctf=8
field=body

SCORE
ctf=6
field=url

field=body
df=4,      ctf=8
docid=3,   tf=2
docid=18, tf=4
:      :    :  :

field=url
df=3,      ctf=6
docid=4,  tf=1
:      :    :  :

**AND Operator Evaluation**
getScore uses the cached docid (4).

args[0] does not match, so call
  args[0].getDefaultScore (4).
Suppose the result is 0.01.

args[1] has match, so call
  args[1].getScore().
Suppose the result is 0.2.

$Score_{AND}(4) = (0.01^{0.5} \times 0.2^{0.5})$

20

Page 10

## Default Belief Scores
## Are Only a Small Complication

**When evaluating a query argument**

• If it has a match for the current document
  – Ask the <u>query argument</u> to calculate the <u>document score</u> for the current document
  – Else ask the <u>query argument</u> to calculate a <u>default score</u> for the current document
    » E.g., a SCORE operator (the example given)
    » E.g., an OR operator (similar logic)

21

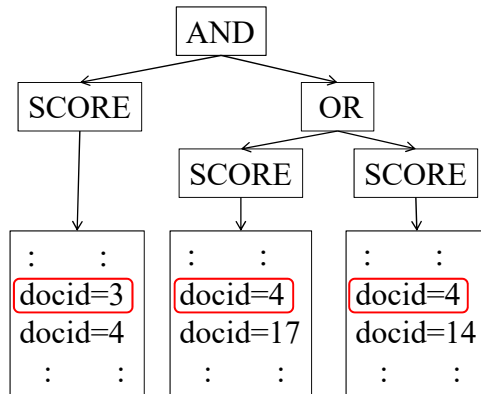## Default Belief Scores
## Are Only a Small Complication

**Which types of query operators calculate default scores?**

• If an operator calculates scores
  … it also calculates default scores
• QrySop operators calculate default scores
• QryIop operators do not calculate default scores

22

## Call to docIteratorHasMatch & getScore
## (First Time)

AND

SCORE     OR

    SCORE     SCORE

```
  :    :       :    :       :    :
docid=3     docid=4     docid=4
docid=4     docid=17    docid=14
  :    :       :    :       :    :
```

**AND Operator Evaluation**
Min document is 3.

args[0] has match, so call
  args[0].getScore ().
Suppose the result is 0.3.

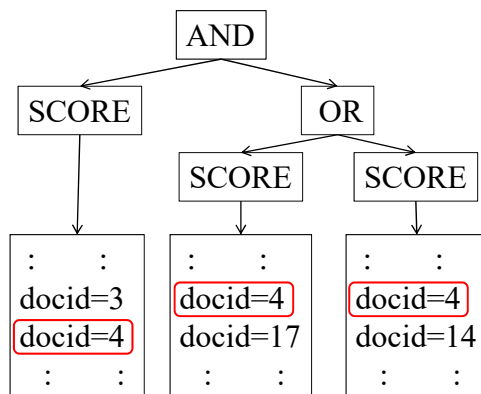args[1] does not match, so call
  args[1].getDefaultScore(3).

OR calls getDefaultScore(3)
  for all of its args and computes
  a score. Suppose it is 0.01.

$Score_{AND}(3) = (0.3^{0.5} \times 0.01^{0.5})$

23

---

## Call to docIteratorHasMatch & getScore
## (Second Time)

AND

SCORE     OR

    SCORE     SCORE

```
  :    :       :    :       :    :
docid=3     docid=4     docid=4
docid=4     docid=17    docid=14
  :    :       :    :       :    :
```

**AND Operator Evaluation**
Min document is 4.

args[0] matches, so call
  args[0].getScore ().
Suppose the result is 0.3.

args[1] matches, so call
  args[1].getScore().
Suppose the score is 0.2.

$Score_{AND}(4) = (0.3^{0.5} \times 0.2^{0.5})$

24

**Using Default Belief Scores Properly**
**Requires Two Components**

**Add a new method to all QrySop operators**

double getDefaultScore (RetrievalModel r, long docid)

  – QrySopScore.getDefaultScore <u>calculates</u> a score for a term

  – QrySop <other>. getDefaultScore <u>combines</u> scores for $n$ terms

**When any QrySop operator calculates scores**

If the $i^{th}$ query argument contains document d

then read its score from the $i^{th}$ score list

else call the $i^{th}$ query argument's getDefaultScore method

 **It may sound complicated now, but actually it is very easy**

---

**Outline**

**HW2 implementation**
- Indri default beliefs
- Window operator

**Proximity Operators:**
**The Window (or Unordered Window) Operator**

**The WINDOW/n operator is used to match related concepts**

- Arguments can be <u>in any order</u>
- n specifies <u>the maximum distance</u> between any pair of terms

**Examples**

- WINDOW/100 (obama merkel putin)
  - We don't care which order they occur in

**Typically proximity operators have complexity O(|C|)**

- A single pass down each inverted list, similar to NEAR/n

---

**Proximity Operators:**
**The Window (or Unordered Window) Operator**

| a | | b | |
|---|---|---|---|
| df: | 47 | df: | 95 |
| doc: | 19 | doc: | 23 |
| tf: | 1 | tf: | 1 |
| locs: | 7 | locs: | 99 |
| doc: | 27 | doc: | 27 |
| tf: | 3 | tf: | 4 |
| locs: | 47 | locs: | 48 |
| | 98 | | 49 |
| | 132 | | 133 |
| doc: | 92 | | 134 |
| … | | doc: | 148 |
| | | ... | |

Query:  # WINDOW/20 (a b)

**Proximity Operators:**
**The Window (or Unordered Window) Operator**

| a | | b | |
|---|---|---|---|
| **a** | | **b** | |
| df: | 47 | df: | 95 |
| doc: | 19 | doc: | 23 |
| tf: | 1 | tf: | 1 |
| locs: | 7 | locs: | 99 |
| doc: | 27 | doc: | 27 |
| tf: | 3 | tf: | 4 |
| locs: | 47 | locs: | 48 |
| | 98 | | 49 |
| | 132 | | 133 |
| doc: | 92 | | 134 |
| … | | doc: | 148 |
| | | ... | |

Query: # WINDOW/20 (a b)

Initialize doc iterators

29                    © 2021 Jamie Callan

---

**Proximity Operators:**
**The Window (or Unordered Window) Operator**

| a | | b | |
|---|---|---|---|
| **a** | | **b** | |
| df: | 47 | df: | 95 |
| doc: | 19 | doc: | 23 |
| tf: | 1 | tf: | 1 |
| locs: | 7 | locs: | 99 |
| doc: | 27 | doc: | 27 |
| tf: | 3 | tf: | 4 |
| locs: | 47 | locs: | 48 |
| | 98 | | 49 |
| | 132 | | 133 |
| doc: | 92 | | 134 |
| … | | doc: | 148 |
| | | ... | |

Query: # WINDOW/20 (a b)

Advance all doc iterators
until they point to the
same document
- This is a simple nested
  loop

30                    © 2021 Jamie Callan

Page 15

## Proximity Operators:
## The Window (or Unordered Window) Operator

| a | | b | |
|---|---|---|---|
| df: | 47 | df: | 95 |
| doc: | 19 | doc: | 23 |
| tf: | 1 | tf: | 1 |
| locs: | 7 | locs: | 99 |
| doc: | 27 | doc: | 27 |
| tf: | 3 | tf: | 4 |
| locs: | 47 | locs: | 48 |
| | 98 | | 49 |
| | 132 | | 133 |
| doc: | 92 | | 134 |
| … | | doc: | 148 |
| | | ... | |

Query: # WINDOW/20 (a b)

Same document
Initialize location iterators

31

---

## Proximity Operators:
## The Window (or Unordered Window) Operator

| a | | b | |
|---|---|---|---|
| df: | 47 | df: | 95 |
| doc: | 19 | doc: | 23 |
| tf: | 1 | tf: | 1 |
| locs: | 7 | locs: | 99 |
| doc: | 27 | doc: | 27 |
| tf: | 3 | tf: | 4 |
| locs: | 47 | locs: | 48 |
| | 98 | | 49 |
| | 132 | | 133 |
| doc: | 92 | | 134 |
| … | | doc: | 148 |
| | | ... | |

Query: # WINDOW/20 (a b)

Find the min (47) and
max (48) locations

32

## Proximity Operators:
## The Window (or Unordered Window) Operator

**a**

df:      47
doc:   19
tf:        1
locs:    7
doc:   27
tf:        3
locs:  47
          98
        132
doc:   92
…

**b**

df:      95
doc:   23
tf:        1
locs:  99
doc:   27
tf:        4
locs:  48
          49
        133
        134
doc: 148
...

Query:  # WINDOW/20 (a b)

Is (max – min) < window?

48 – 47 < 20 (match)

Record match
• max location (48)

33

© 2021 Jamie Callan

---

## Proximity Operators:
## The Window (or Unordered Window) Operator

**a**

df:      47
doc:   19
tf:        1
locs:    7
doc:   27
tf:        3
locs:  47
          98
        132
doc:   92
…

**b**

df:      95
doc:   23
tf:        1
locs:  99
doc:   27
tf:        4
locs:  48
          49
        133
        134
doc: 148
...

Query:  # WINDOW/20 (a b)

Increment all loc iterators

34

© 2021 Jamie Callan

Page 17

## Proximity Operators:
## The Window (or Unordered Window) Operator

**a**

| | |
|---|---|
| df: | 47 |
| doc: | 19 |
| tf: | 1 |
| locs: | 7 |
| doc: | 27 |
| tf: | 3 |
| locs: | 47 |
| | 98 |
| | 132 |
| doc: | 92 |
| … | |

**b**

| | |
|---|---|
| df: | 95 |
| doc: | 23 |
| tf: | 1 |
| locs: | 99 |
| doc: | 27 |
| tf: | 4 |
| locs: | 48 |
| | 49 |
| | 133 |
| | 134 |
| doc: | 148 |
| ... | |

Query: # WINDOW/20 (a b)

Find the min (49) and
max (98) locations

35

---

## Proximity Operators:
## The Window (or Unordered Window) Operator

**a**

| | |
|---|---|
| df: | 47 |
| doc: | 19 |
| tf: | 1 |
| locs: | 7 |
| doc: | 27 |
| tf: | 3 |
| locs: | 47 |
| | 98 |
| | 132 |
| doc: | 92 |
| … | |

**b**

| | |
|---|---|
| df: | 95 |
| doc: | 23 |
| tf: | 1 |
| locs: | 99 |
| doc: | 27 |
| tf: | 4 |
| locs: | 48 |
| | 49 |
| | 133 |
| | 134 |
| doc: | 148 |
| ... | |

Query: # WINDOW/20 (a b)

Is (max – min) < window?

$98 – 49 \geq 20$ (no match)

36

Page 18

## Proximity Operators:
## The Window (or Unordered Window) Operator

**a**        **b**
df:    47    df:    95
doc:  19    doc:  23
tf:     1    tf:     1
locs:  7    locs:  99
doc:  27    doc:  27
tf:     3    tf:     4
locs: 47    locs:  48
        98            49
       132          133
doc:  92          134
…          doc: 148
            ...

Query:  # WINDOW/20 (a b)

Increment the iterator for
the min location

37                          © 2021 Jamie Callan

---

## Proximity Operators:
## The Window (or Unordered Window) Operator

**a**        **b**
df:    47    df:    95
doc:  19    doc:  23
tf:     1    tf:     1
locs:  7    locs:  99
doc:  27    doc:  27
tf:     3    tf:     4
locs: 47    locs:  48
        98            49
       132          133
doc:  92          134
…          doc: 148
            ...

Query:  # WINDOW/20 (a b)

Find the min (98) and
max (133) locations

38                          © 2021 Jamie Callan

Page 19

## Proximity Operators:
## The Window (or Unordered Window) Operator

**a**

df:  47
doc:  19
tf:  1
locs:  7
doc:  27
tf:  3
locs: 47
98
132
doc:  92
…

**b**

df:  95
doc:  23
tf:  1
locs:  99
doc:  27
tf:  4
locs:  48
49
133
134
doc: 148
...

Query:  # WINDOW/20 (a b)

Is (max – min) < window?

133 – 98 ≥ 20 (no match)

39

© 2021 Jamie Callan

## Proximity Operators:
## The Window (or Unordered Window) Operator

**a**

df:  47
doc:  19
tf:  1
locs:  7
doc:  27
tf:  3
locs: 47
98
132
doc:  92
…

**b**

df:  95
doc:  23
tf:  1
locs:  99
doc:  27
tf:  4
locs:  48
49
133
134
doc: 148
...

Query:  # WINDOW/20 (a b)

Increment the iterator for
the min location

40

© 2021 Jamie Callan

Page 20

## Proximity Operators:
## The Window (or Unordered Window) Operator

**a**

df: 47
doc: 19
tf: 1
locs: 7
doc: 27
tf: 3
locs: 47
 98
 132
doc: 92
…

**b**

df: 95
doc: 23
tf: 1
locs: 99
doc: 27
tf: 4
locs: 48
 49
 133
 134
doc: 148
...

Query: # WINDOW/20 (a b)

Find the min (132) and
max (133) locations

41

---

## Proximity Operators:
## The Window (or Unordered Window) Operator

**a**

df: 47
doc: 19
tf: 1
locs: 7
doc: 27
tf: 3
locs: 47
 98
 132
doc: 92
…

**b**

df: 95
doc: 23
tf: 1
locs: 99
doc: 27
tf: 4
locs: 48
 49
 133
 134
doc: 148
...

Query: # WINDOW/20 (a b)

Is (max – min) < window?

133 – 132 < 20 (match)

Record match
• max location (133)

42

Page 21

**Proximity Operators:**
**The Window (or Unordered Window) Operator**

| a | | b | |
|---|---|---|---|
| df: | 47 | df: | 95 |
| doc: | 19 | doc: | 23 |
| tf: | 1 | tf: | 1 |
| locs: | 7 | locs: | 99 |
| doc: | 27 | doc: | 27 |
| tf: | 3 | tf: | 4 |
| locs: | 47 | locs: | 48 |
| | 98 | | 49 |
| | 132 | | 133 |
| doc: | 92 | | 134 |
| … | | doc: | 148 |
| | | ... | |

Query: # WINDOW/20 (a b)

Increment all loc iterators

$q_0$ locs are exhausted.
No more matches are
  possible in this document

43

© 2021 Jamie Callan

---

| a | | b | |
|---|---|---|---|
| df: | 47 | df: | 95 |
| doc: | 19 | doc: | 23 |
| tf: | 1 | tf: | 1 |
| locs: | 7 | locs: | 99 |
| doc: | 27 | doc: | 27 |
| tf: | 3 | tf: | 4 |
| locs: | 47 | locs: | 48 |
| | 98 | | 49 |
| | 132 | | 133 |
| doc: | 92 | | 134 |
| … | | doc: | 148 |
| | | ... | |

Query: # WINDOW/20 (a b)

Increment all doc iterators
…

Continue until the inverted
lists are exhausted

44

© 2021 Jamie Callan

**Proximity Operators:**
**The Window (or Unordered Window) Operator**

**Implementation note**

- A document term can only match the query once
- **Query:** #WINDOW/100 (obama merkel putin)
- **Document:** obama … merkel … putin … merkel … obama
- There is just <u>one</u> match here

**One can imagine other implementations, but this is the norm**

- Usually more complicated matching doesn't improve accuracy