

# Vector Space Scoring

Introduction to Information Retrieval  
CS 221  
Donald J. Patterson

Content adapted from Hinrich Schütze  
<http://www.informationretrieval.org>

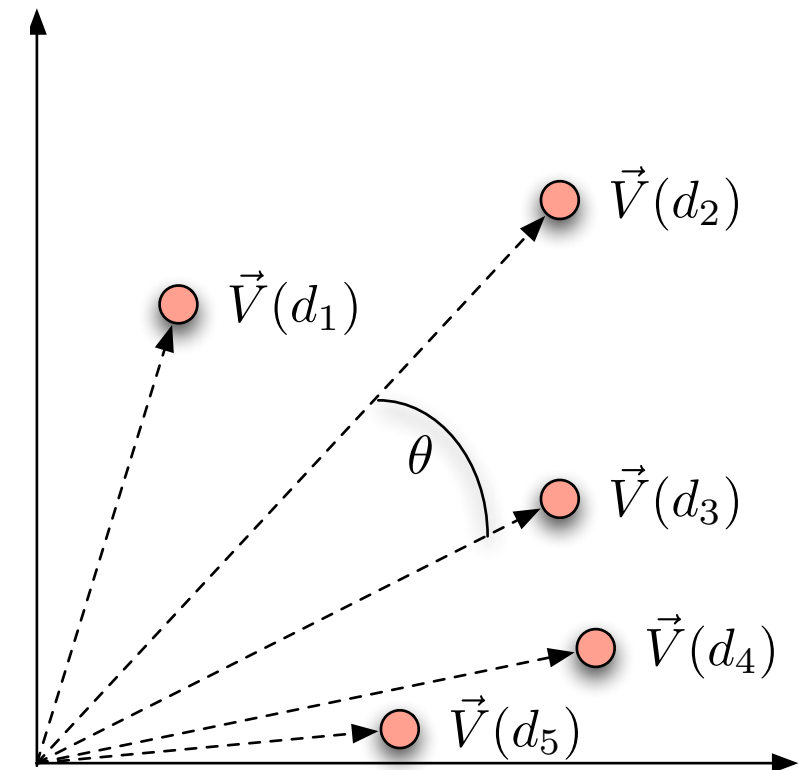


## Cosine Similarity Score

- Define: Euclidean Length

$$|\vec{V}(d_1)| = \sqrt{\sum_{i=t_1}^{t_n} (\vec{V}(d_1)_i \vec{V}(d_1)_i)}$$

|                  | <i>Antony and Cleopatra</i> | <i>Julius Caesar</i> | <i>The Tempest</i> | <i>Hamlet</i> | <i>Othello</i> | <i>Macbeth</i> |
|------------------|-----------------------------|----------------------|--------------------|---------------|----------------|----------------|
| <i>Antony</i>    | 13.1                        | 11.4                 | 0.0                | 0.0           | 0.0            | 0.0            |
| <i>Brutus</i>    | 3.0                         | 8.3                  | 0.0                | 1.0           | 0.0            | 0.0            |
| <i>Caesar</i>    | 2.3                         | 2.3                  | 0.0                | 0.5           | 0.3            | 0.3            |
| <i>Calpurnia</i> | 0.0                         | 11.2                 | 0.0                | 0.0           | 0.0            | 0.0            |
| <i>Cleopatra</i> | 17.7                        | 0.0                  | 0.0                | 0.0           | 0.0            | 0.0            |
| <i>mercy</i>     | 0.5                         | 0.0                  | 0.7                | 0.9           | 0.9            | 0.3            |
| <i>worser</i>    | 1.2                         | 0.0                  | 0.6                | 0.6           | 0.6            | 0.0            |



$$\begin{aligned} |\vec{V}(d_1)| &= \sqrt{(13.1 * 13.1) + (3.0 * 3.0) + (2.3 * 2.3) + (17.7 * 17.7) + (0.5 * 0.5) + (1.2 * 1.2)} \\ &= 22.38 \end{aligned}$$

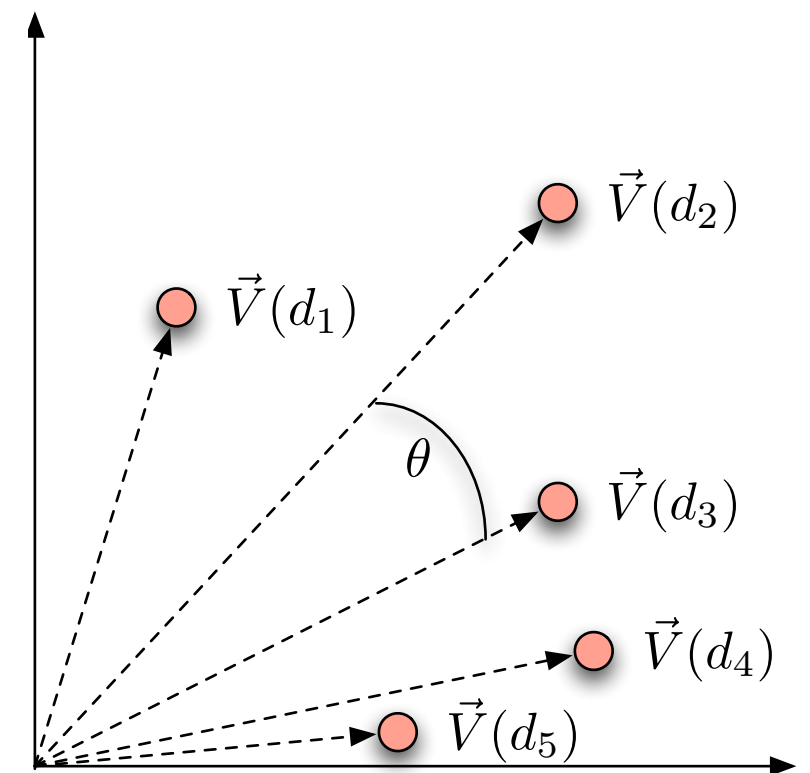


## Cosine Similarity Score

- Define: Euclidean Length

$$|\vec{V}(d_1)| = \sqrt{\sum_{i=t_1}^{t_n} (\vec{V}(d_1)_i \vec{V}(d_1)_i)}$$

|                  | <i>Antony and Cleopatra</i> | <i>Julius Caesar</i> | <i>The Tempest</i> | <i>Hamlet</i> | <i>Othello</i> | <i>Macbeth</i> |
|------------------|-----------------------------|----------------------|--------------------|---------------|----------------|----------------|
| <i>Antony</i>    | 13.1                        | 11.4                 | 0.0                | 0.0           | 0.0            | 0.0            |
| <i>Brutus</i>    | 3.0                         | 8.3                  | 0.0                | 1.0           | 0.0            | 0.0            |
| <i>Caesar</i>    | 2.3                         | 2.3                  | 0.0                | 0.5           | 0.3            | 0.3            |
| <i>Calpurnia</i> | 0.0                         | 11.2                 | 0.0                | 0.0           | 0.0            | 0.0            |
| <i>Cleopatra</i> | 17.7                        | 0.0                  | 0.0                | 0.0           | 0.0            | 0.0            |
| <i>mercy</i>     | 0.5                         | 0.0                  | 0.7                | 0.9           | 0.9            | 0.3            |
| <i>worser</i>    | 1.2                         | 0.0                  | 0.6                | 0.6           | 0.6            | 0.0            |



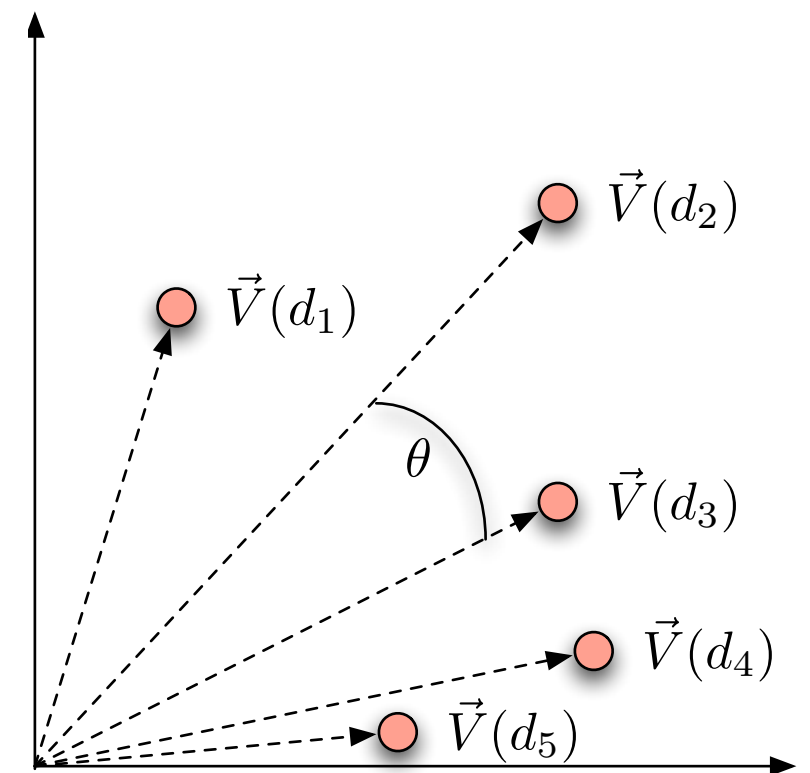
$$\begin{aligned} |\vec{V}(d_1)| &= \sqrt{(11.4 * 11.4) + (8.3 * 8.3) + (2.3 * 2.3) + (11.2 * 11.2)} \\ &= 18.15 \end{aligned}$$



## Cosine Similarity Score

- Example

$$\begin{aligned} \text{sim}(d_1, d_2) &= \frac{\vec{V}(d_1) \cdot \vec{V}(d_2)}{|\vec{V}(d_1)| |\vec{V}(d_2)|} \\ &= \frac{179.53}{22.38 * 18.15} \\ &= 0.442 \end{aligned}$$



## Exercise

- Rank the following by decreasing cosine similarity.
  - Assume tf-idf weighting:
    - Two docs that have only frequent words in common
      - (the, a , an, of)
    - Two docs that have no words in common
    - Two docs that have many rare words in common
      - (mocha, volatile, organic, shade-grown)



# Vector Space Scoring

## Exercise

tf =

|    |    |   |    |    |    |    |    |    |    |    |    |    |    |    |    |
|----|----|---|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 24 | 24 | 0 | 24 | 24 | 24 | 24 | 24 | 24 | 24 | 24 | 24 | 24 | 24 | 24 | 24 |
| 10 | 10 | 0 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 |
| 24 | 24 | 0 | 24 | 24 | 24 | 24 | 24 | 24 | 24 | 24 | 24 | 24 | 24 | 24 | 24 |
| 12 | 12 | 0 | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 11 |
| 10 | 10 | 0 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 |
| 0  | 0  | 0 | 0  | 1  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  |
| 0  | 0  | 1 | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  |
| 0  | 0  | 1 | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  |
| 0  | 0  | 0 | 0  | 1  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  |
| 0  | 0  | 0 | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  |
| 0  | 0  | 0 | 0  | 1  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  |
| 0  | 0  | 1 | 1  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  |
| 0  | 0  | 1 | 1  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  |
| 0  | 0  | 1 | 1  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  |

df =

14  
14  
14  
14  
14  
1  
1  
1  
1  
1  
1  
2  
2  
2



## Exercise

```
c = 15;
tf = load('tf.txt', '-ASCII');
df = load('df.txt', '-ASCII');
tfidf = zeros(size(tf));

for i = 1:size(tf,1)
    for j = 1:size(tf,2)
        if tf(i,j) == 0
            tfidf(i,j) = (0) * log2(c/df(i));
        else
            tfidf(i,j) = (1+log2(tf(i,j))) * log2(c/df(i));
        end
    end
end
```



# Vector Space Scoring

## Exercise

tfidf =

|        |        |        |        |        |        |        |        |
|--------|--------|--------|--------|--------|--------|--------|--------|
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |
| 0.4564 | 0.4564 | 0      | 0.4439 | 0.4439 | 0.4439 | 0.4439 | 0.4439 |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |
| 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      |
| 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |

More of the same →





# Vector Space Scoring

## Exercise

tfidf =

|        |        |        |        |        |        |        |        |                          |
|--------|--------|--------|--------|--------|--------|--------|--------|--------------------------|
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 | More of<br>the same<br>→ |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |                          |
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |                          |
| 0.4564 | 0.4564 | 0      | 0.4439 | 0.4439 | 0.4439 | 0.4439 | 0.4439 |                          |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |                          |
| 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      |                          |
| 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      |                          |
| 0      | 0      | 0      | 0      | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      |                          |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |                          |



# Exercise

[illegible]

```
>> num = tfidf(:,1)'*tfidf(:,2)

num =

    1.1964
```

num =

1.1964

# More of the same



# Exercise

|        |        |
|--------|--------|
| 0.5559 | 0.5559 |
| 0.4302 | 0.4302 |
| 0.5559 | 0.5559 |
| 0.4564 | 0.4564 |
| 0.4302 | 0.4302 |
| 0      | 0      |
| 0      | 0      |
| 0      | 0      |
| 0      | 0      |
| 0      | 0      |
| 0      | 0      |
| 0      | 0      |
| 0      | >> der |
| 0      | .      |

```
>> num = tfidf(:,1)'*tfidf(:,2)

num =

    1.1964
```

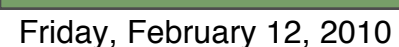
num =

1.1964

denom =

1.1964

# More of the same



# Exercise

```
>> num = tfidf(:,1)'*tfidf(:,2)

num =

    1.1964
```

num =

```
>> denom = sqrt(tfi df(:,1)'*tfi df(:,1)).*sqrt(tfi df(:,2)'*tfi df(:,2))
```

1.1964

score =

1.0000

# More of the same

# Vector Space Scoring

## Exercise

tfidf =

|        |        |        |        |        |        |        |        |        |                          |
|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------------------------|
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 | More of<br>the same<br>→ |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |                          |
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |                          |
| 0.4564 | 0.4564 | 0      | 0.4439 | 0.4439 | 0.4439 | 0.4439 | 0.4439 | 0.4439 |                          |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |                          |
| 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      |                          |
| 0      | 0      | 0      | 0      | 0      | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      | 0      |                          |



# Vector Space Scoring

## Exercise

tfidf =

|        |        |        |
|--------|--------|--------|
| 0.5559 | 0.5559 | 0      |
| 0.4302 | 0.4302 | 0      |
| 0.5559 | 0.5559 | 0      |
| 0.4564 | 0.4564 | 0      |
| 0.4302 | 0.4302 | 0      |
| 0      | 0      | 0      |
| 0      | 0      | 3.9069 |
| 0      | 0      | 3.9069 |
| 0      | 0      | 0      |
| 0      | 0      | 0      |
| 0      | 0      | 0      |
| 0      | 0      | 2.9069 |
| 0      | 0      | 2.9069 |
| 0      | 0      | 2.9069 |

|        |        |        |        |        |
|--------|--------|--------|--------|--------|
| 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |
| 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |
| 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |
| 0.4439 | 0.4439 | 0.4439 | 0.4439 | 0.4439 |
| 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |
| 0      | 3.9069 | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      |

```
>> num = tfidf(:,2)'*tfidf(:,3)
```

```
num =
```

```
2.9069  
2.9069 0  
2.9069 0 0
```

More of  
the same



# Vector Space Scoring

## Exercise

tfidf =

|        |        |        |        |        |        |        |        |
|--------|--------|--------|--------|--------|--------|--------|--------|
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |
| 0.4564 | 0.4564 | 0      | 0.4439 | 0.4439 | 0.4439 | 0.4439 | 0.4439 |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |
| 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      |
| 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 3.9069 |        |        |        |        |        |
| 0      | 0      | 0      |        |        |        |        |        |
| 0      | 0      | 0      |        |        |        |        |        |
| 0      | 0      | 0      |        |        |        |        |        |
| 0      | 0      | 2.9069 | 2.9069 |        |        |        |        |
| 0      | 0      | 2.9069 | 2.9069 | 0      |        |        |        |
| 0      | 0      | 2.9069 | 2.9069 |        |        |        |        |

More of the same →

```
>> num = tfidf(:,2)'*tfidf(:,3)
```

```
num =
```

```
2.9069
```

```
>> denom = sqrt(tfidf(:,2)'*tfidf(:,2)).*sqrt(tfidf(:,3)'*tfidf(:,3))
```

```
denom =
```

```
8.1765
```





# Vector Space Scoring

## Exercise

tfidf =

|        |        |        |        |        |        |        |        |
|--------|--------|--------|--------|--------|--------|--------|--------|
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |
| 0.4564 | 0.4564 | 0      | 0.4439 | 0.4439 | 0.4439 | 0.4439 | 0.4439 |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |
| 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      |
| 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 3.9069 |        |        |        |        |        |
| 0      | 0      | 0      |        |        |        |        |        |
| 0      | 0      | 0      |        |        |        |        |        |
| 0      | 0      | 0      |        |        |        |        |        |
| 0      | 0      | 2.9069 | 2.9069 |        |        |        |        |
| 0      | 0      | 2.9069 | 2.9069 | 0      |        |        |        |
| 0      | 0      | 2.9069 | 2.9069 |        |        |        |        |

More of the same →

```
>> num = tfidf(:,2)'*tfidf(:,3)
```

```
num =
```

```
2.9069  
2.9069  
2.9069
```

```
>> denom = sqrt(tfidf(:,2)'*tfidf(:,2)).*sqrt(tfidf(:,3)'*tfidf(:,3))
```

```
denom =
```

```
8.1765
```

```
>> score = num/denom
```

```
score =
```

```
0
```



# Vector Space Scoring

## Exercise

tfidf =

|        |        |        |        |        |        |        |        |                          |
|--------|--------|--------|--------|--------|--------|--------|--------|--------------------------|
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 | More of<br>the same<br>→ |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |                          |
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |                          |
| 0.4564 | 0.4564 | 0      | 0.4439 | 0.4439 | 0.4439 | 0.4439 | 0.4439 |                          |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |                          |
| 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      |                          |
| 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      |                          |
| 0      | 0      | 0      | 0      | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      |                          |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |                          |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |                          |



# Vector Space Scoring

## Exercise

tfidf =

|        |        |        |        |        |        |        |        |
|--------|--------|--------|--------|--------|--------|--------|--------|
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |
| 0.4564 | 0.4564 | 0      | 0.4439 | 0.4439 | 0.4439 | 0.4439 | 0.4439 |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |
| 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      |
| 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |

```
>> num = tfidf(:,3)'*tfidf(:,4)
```

```
num =
```

```
25.3500
```

More of  
the same



# Vector Space Scoring

## Exercise

tfidf =

|        |        |        |        |        |        |        |        |
|--------|--------|--------|--------|--------|--------|--------|--------|
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |
| 0.4564 | 0.4564 | 0      | 0.4439 | 0.4439 | 0.4439 | 0.4439 | 0.4439 |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |
| 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      |
| 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |

More of the same →

```
>> num = tfidf(:,3)'*tfidf(:,4)
```

num =

25.3500

```
>> denom = sqrt(tfidf(:,3)'*tfidf(:,3)).*sqrt(tfidf(:,4)'*tfidf(:,4))
```

denom =

38.5062



# Vector Space Scoring

## Exercise

tfidf =

|        |        |        |        |        |        |        |        |
|--------|--------|--------|--------|--------|--------|--------|--------|
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |
| 0.5559 | 0.5559 | 0      | 0.5559 | 0.5559 | 0.5559 | 0.5559 | 0.5559 |
| 0.4564 | 0.4564 | 0      | 0.4439 | 0.4439 | 0.4439 | 0.4439 | 0.4439 |
| 0.4302 | 0.4302 | 0      | 0.4302 | 0.4302 | 0.4302 | 0.4302 | 0.4302 |
| 0      | 0      | 0      | 0      | 3.9069 | 0      | 0      | 0      |
| 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 3.9069 | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 0      | 0      | 0      | 0      | 0      | 0      |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |
| 0      | 0      | 2.9069 | 2.9069 | 0      | 0      | 0      | 0      |

More of the same →

```
>> num = tfidf(:,3)'*tfidf(:,4)
```

```
num =
```

```
25.3500
```

```
>> denom = sqrt(tfidf(:,3)'*tfidf(:,3)).*sqrt(tfidf(:,4)'*tfidf(:,4))
```

```
denom =
```

```
38.5062
```

```
>> score = num/denom
```

```
score =
```

```
0.6583
```

## Exercise

- Rank the following by decreasing cosine similarity.
  - Assume tf-idf weighting:
    - Two docs that have only frequent words in common
      - (the, a , an, of)
    - Two docs that have no words in common
    - Two docs that have many rare words in common
      - (mocha, volatile, organic, shade-grown)



## Exercise

- Rank the following by decreasing cosine similarity.
  - Assume tf-idf weighting:
    - Two docs that have only frequent words in
      - (the, a , an, of)
    - Two docs that have no words in common
    - Two docs that have many rare words in common
      - (mocha, volatile, organic, shade-grown)

```
>> score = num/denom  
score =  
1.0000
```



## Exercise

- Rank the following by decreasing cosine similarity.
- Assume tf-idf weighting:
  - Two docs that have only frequent words in common
    - (the, a , an, of)
  - Two docs that have no words in common
  - Two docs that have many rare words in common
    - (mocha, volatile, organic, shade-grown)

```
>> score = num/denom  
score =  
1.0000
```

```
>> score = num/denom  
score =  
0
```





## Exercise

- Rank the following by decreasing cosine similarity.
- Assume tf-idf weighting:
  - Two docs that have only frequent words in common
    - (the, a , an, of)
  - Two docs that have no words in common
  - Two docs that have many rare words in common
    - (mocha, volatile, organic, shade-grown)

```
>> score = num/denom  
score =  
1.0000
```

```
>> score = num/denom  
score =  
0
```

```
>> score = num/denom  
score =  
0.6583
```





## Spamming indices

- This was invented before spam
- Consider:
  - Indexing a sensible passive document collection
  - vs.
  - Indexing an active document collection, where people, companies, bots are shaping documents to maximize scores
- Vector space scoring may not be as useful in this context.



## Interaction: vectors and phrases

- Scoring phrases doesn't naturally fit into the vector space world:
  - How do we get beyond the "bag of words"?
  - "dark roast" and "pot roast"
  - There is no information on "dark roast" as a phrase in our indices.
- Biword index can treat some phrases as terms
  - postings for phrases
  - document wide statistics for phrases



## Interaction: vectors and phrases

- Theoretical problem:
  - Axes of our term space are now correlated
    - There is a lot of shared information in “light roast” and “dark roast” rows of our index
- End-user problem:
  - A user doesn’t know which phrases are indexed and can’t effectively discriminate results.



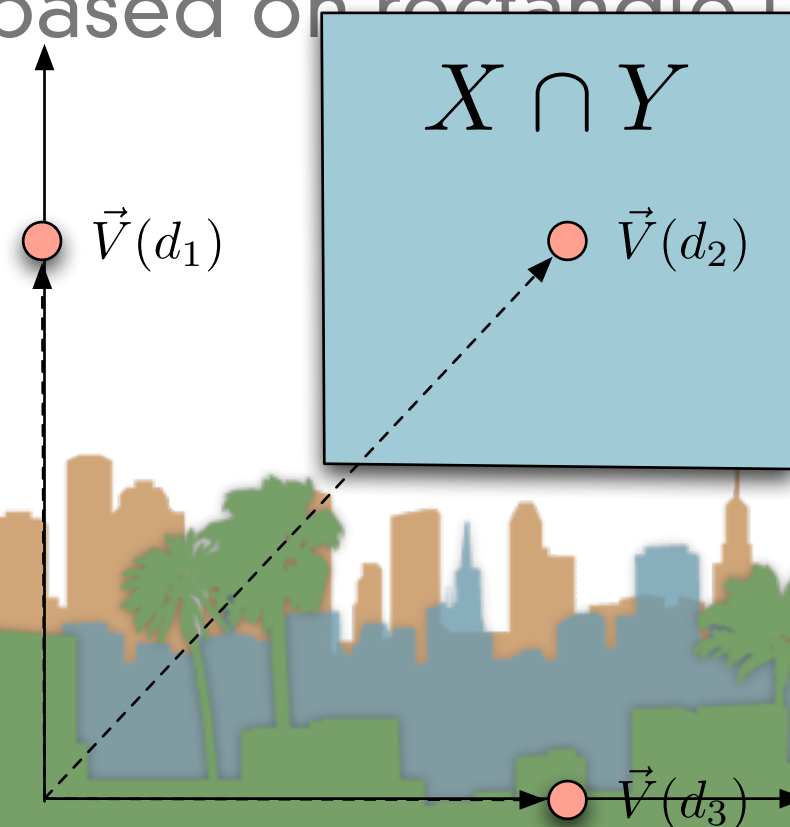
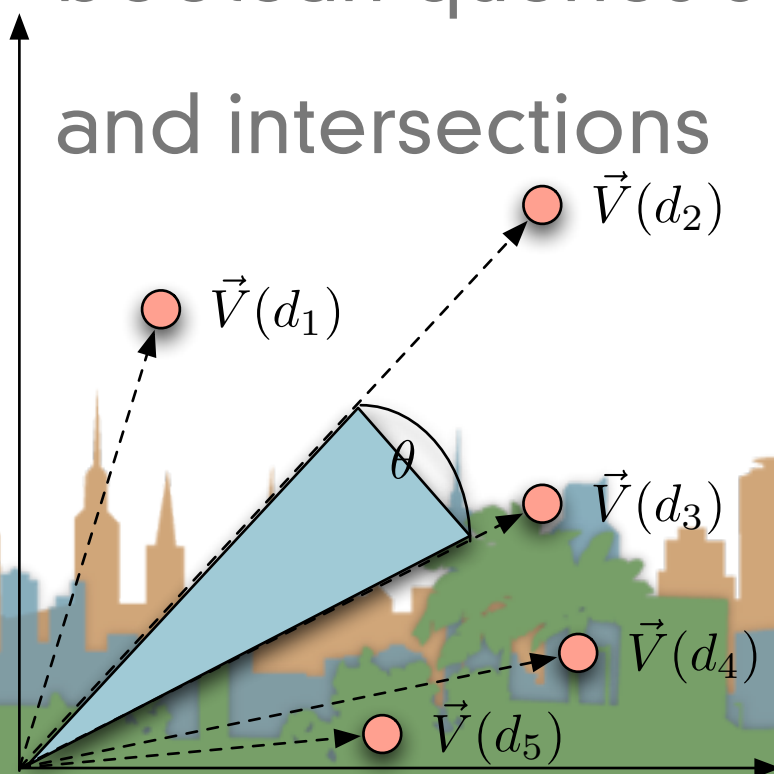
# Multiple queries for phrases and vectors

- Query: “rising interest rates”
- Iterative refinement:
  - Run the phrase query vector with 3 words as a term.
  - If not enough results, run 2-phrase queries and fold into results: “rising interest” “interest rates”
  - If still not enough results run query with three words as separate terms.



## Vectors and Boolean queries

- Ranked queries and Boolean queries don't work very well together
- In term space
  - ranked queries select based on sector containment - cosine similarity
  - boolean queries select based on rectangle unions



# Vectors and wild cards



# Vectors and wild cards

- How could we work with the query, “quick\* print\*” ?



## Vectors and wild cards

- How could we work with the query, “quick\* print\*” ?
- Can we view this as a bag of words?





## Vectors and wild cards

- How could we work with the query, “quick\* print\*” ?
- Can we view this as a bag of words?
- What about expanding each wild-card into the matching set of dictionary terms?



## Vectors and wild cards

- How could we work with the query, “quick\* print\*” ?
  - Can we view this as a bag of words?
  - What about expanding each wild-card into the matching set of dictionary terms?
- Danger: Unlike the boolean case, we now have tf's and idfs to deal with



## Vectors and wild cards

- How could we work with the query, “quick\* print\*” ?
  - Can we view this as a bag of words?
  - What about expanding each wild-card into the matching set of dictionary terms?
- Danger: Unlike the boolean case, we now have tf's and idfs to deal with
- Overall, not a great idea



# Vectors and other operators

- Vector space queries are good for no-syntax, bag-of-words queries
  - Nice mathematical formalism
  - Clear metaphor for similar document queries
  - Doesn't work well with Boolean, wild-card or positional query operators
- But ...



# Query language vs. Scoring

- Interfaces to the rescue
  - Free text queries are often separated from operator query language
  - Default is free text query
  - Advanced query operators are available in “advanced query” section of interface
  - Or embedded in free text query with special syntax
    - aka -term -“terma termb”



## Alternatives to tf-idf

- Sublinear tf scaling
  - 20 occurrences of “mole” does not indicate 20 times the relevance

- This motivated the WTF score.

$WTF(t, d)$

1    **if**  $tf_{t,d} = 0$

2        **then**  $return(0)$

3        **else**  $return(1 + \log(tf_{t,d}))$

- There are other variants for reducing the impact of repeated terms



## TF Normalization

- Normalize tf weights by maximum tf in that document

$$ntf_{t,d} = \alpha + (1 - \alpha) \frac{tf_{t,d}}{tf_{max}(d)}$$

- alpha is a smoothing term from (0 - 1.0 ) ~0.4 in practice
- This addresses a length bias.
- Take one document, repeat it, WTF goes up



## TF Normalization

- Normalize tf weights by maximum tf in that document

$$ntf_{t,d} = \alpha + (1 - \alpha) \frac{tf_{t,d}}{tf_{max}(d)}$$

- a change in the **stop word list** can change weights drastically - hard to tune
- still based on bag of words model
  - one outlier word, repeated many times might throw off the algorithmic understanding of the content





## Laundry List

| <i>Term Frequency</i> |  | <i>Document Frequency</i> |  | <i>Normalization</i> |  |
|-----------------------|--|---------------------------|--|----------------------|--|
| <i>(n)atural</i>      | $tf_{t,d}$   | <i>(n)o</i>               | 1  | <i>(n)one</i>        | 1  |
| <i>(l)ogarithm</i>    | $1 + \log(tf_{t,d})$   | <i>(t)idf</i>             | $\log \frac{ \text{corpus} }{df_t}$                    | <i>(c)osine</i>      | $\frac{1}{\sqrt{w_1^2 + w_2^2 + \dots + w_m^2}}$ |
| <i>(a)ugmented</i>    | $\alpha + (1 - \alpha) \frac{tf_{t,d}}{tf_{max}(d)}$           | <i>(p)robidf</i>          | $\max\{0, \log(\frac{ \text{corpus}  - df_t}{df_t})\}$ | <i>(u)pivoted</i>    | $1/u$  |
| <i>(b)oolean</i>      | $tf_{t,d} > 0 ? 1 : 0$   |                           |  | <i>(b)yte</i>        | $1/CharLength^\alpha, \alpha < 1$                |
| <i>(L)ogaverage</i>   | $\frac{1 + \log(tf_{t,d})}{1 + \log(ave_{t \in d}(tf_{t,d}))}$ |                           |  |                      |  |

- SMART system of describing your IR vector algorithm
  - ddd.qqq (ddd = document weighting) (qqq = query weighting)
  - first is term weighting, second is document, then normalization
  - Inc.ltc is what?



# Efficient Cosine Ranking

- Find the  $k$  docs in the corpus “nearest” to the query
  - the  $k$  largest query-doc cosines
- Efficient ranking means:
  - Computing a single cosine efficiently
  - Computing the  $k$  largest cosine values efficiently
    - Can we do this without computing all  $n$  cosines?
      - $n$  = number of documents in corpus



## Efficient Cosine Ranking

- Computing a single cosine
- Use inverted index
- At query time use an array of accumulators  $A_j$  to accumulate component-wise sum
- Accumulate scores as postings lists are being processed (numerator of similarity score)

$$A_j = \sum_t (w_{q,t} w_{d,t})$$



# Efficient Cosine Ranking

- For the web
  - an array of accumulators in memory is infeasible
  - so only create accumulators for docs that occur in postings list
    - dynamically create accumulators
  - put the `tf_d` scores in the postings lists themselves
  - limit docs to non-zero cosines on rare words
    - or non-zero cosines on all words
  - reduces number of accumulators



## Efficient Cosine Ranking

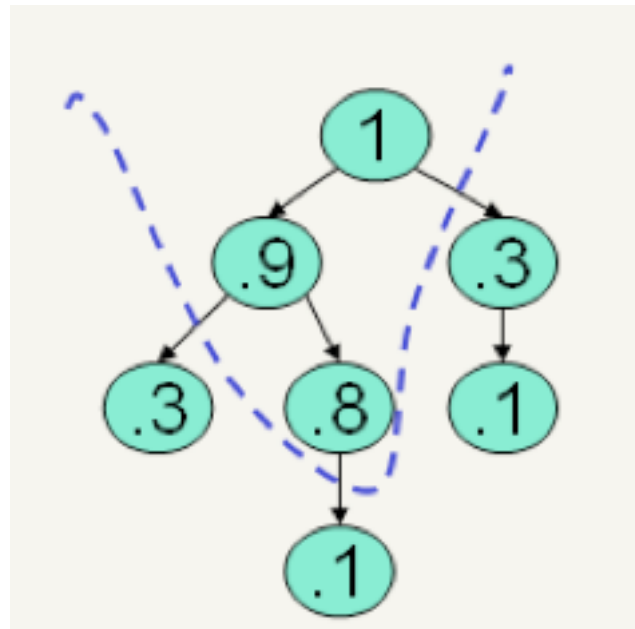
COSINESCORE( $q$ )

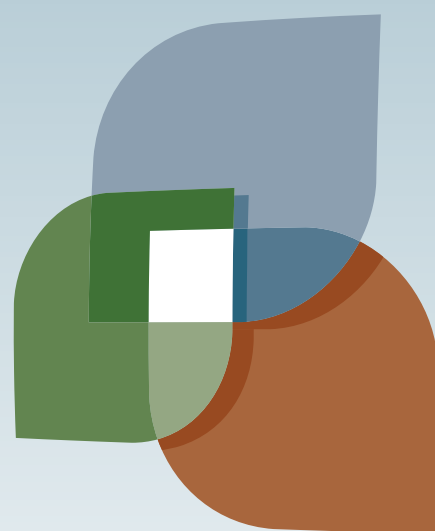
```
1  INITIALIZE( $Scores[d \in D]$ )
2  INITIALIZE( $Magnitude[d \in D]$ )
3  for each term ( $t \in q$ )
4      do  $p \leftarrow \text{FETCHPOSTINGSLIST}(t)$ 
5           $df_t \leftarrow \text{GETCORPUSWIDESTATS}(p)$ 
6           $\alpha_{t,q} \leftarrow \text{WEIGHTINQUERY}(t, q, df_t)$ 
7          for each  $\{d, tf_{t,d}\} \in p$ 
8              do  $Scores[d] += \alpha_{t,q} \cdot \text{WEIGHTINDOCUMENT}(t, q, df_t)$ 
9  for  $d \in Scores$ 
10     do  $\text{NORMALIZE}(Scores[d], Magnitude[d])$ 
11  return top  $K \in Scores$ 
```



## Use heap for selecting the top K Scores

- Binary tree in which each node's value  $>$  the values of children
- Takes  $2N$  operations to construct
  - then each of  $k$  "winners" read off in  $2\log n$  steps
  - For  $n = 1M$ ,  $k = 100$  this is about 10% of the cost of sorting





L U C I

