

Clarification of TF-IDF score

- Some of the slides show this formula:

$$tfidf(t, d) = (1 + \log(tf_{t,d})) * \log\left(\frac{|corpus|}{df_{t,d}}\right)$$

- Precisely it should be:

$$tfidf(t, d) = WTF(t, d) * \log\left(\frac{|corpus|}{df_{t,d}}\right)$$

- The difference is just the special case when $tf = 0$

$WTF(t, d)$

1 **if** $tf_{t,d} = 0$

2 **then** $return(0)$

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



Queries in the vector space model

- Central idea: the query is a vector
- We regard the query as a short document
- We return the documents ranked by the closeness of their vectors to the query (also a vector)

$$\text{sim}(q, d_i) = \frac{\vec{V}(q) \cdot \vec{V}(d_i)}{|\vec{V}(q)| |\vec{V}(d_i)|}$$

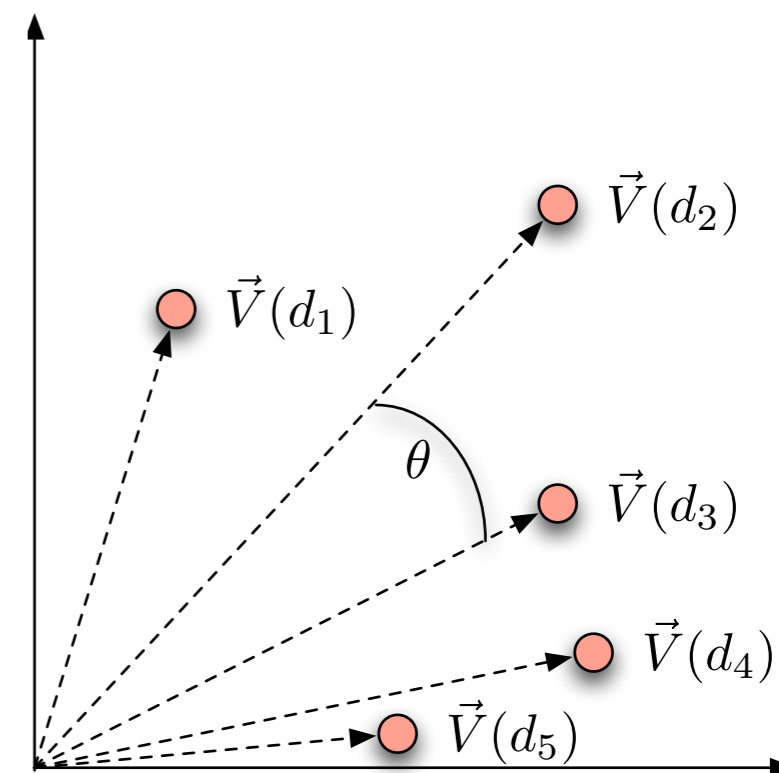
- Note that q is very sparse!



Cosine Similarity Score

- Also called **cosine similarity**

$$\begin{aligned}\vec{V}(d_1) \cdot \vec{V}(d_2) &= \frac{|\vec{V}(d_1)| |\vec{V}(d_2)|}{\cos(\theta)} \\ \cos(\theta) &= \frac{\vec{V}(d_1) \cdot \vec{V}(d_2)}{|\vec{V}(d_1)| |\vec{V}(d_2)|} \\ \text{sim}(d_1, d_2) &= \frac{\vec{V}(d_1) \cdot \vec{V}(d_2)}{|\vec{V}(d_1)| |\vec{V}(d_2)|}\end{aligned}$$

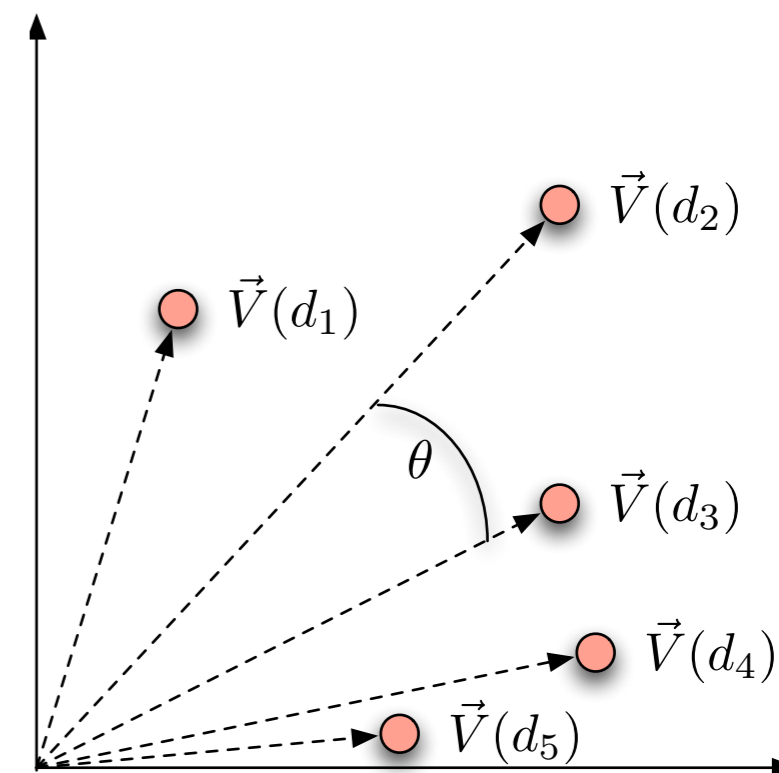


Cosine Similarity Score

- Define: dot product

$$\vec{V}(d_1) \cdot \vec{V}(d_2) = \sum_{i=t_1}^{t_n} (\vec{V}(d_1)_i \vec{V}(d_2)_i)$$

| | <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 |



$$\begin{aligned} \vec{V}(d_1) \cdot \vec{V}(d_2) &= (13.1 * 11.4) + (3.0 * 8.3) + (2.3 * 2.3) + (0 * 11.2) + (17.7 * 0) + (0.5 * 0) + (1.2 * 0) \\ &= 179.53 \end{aligned}$$

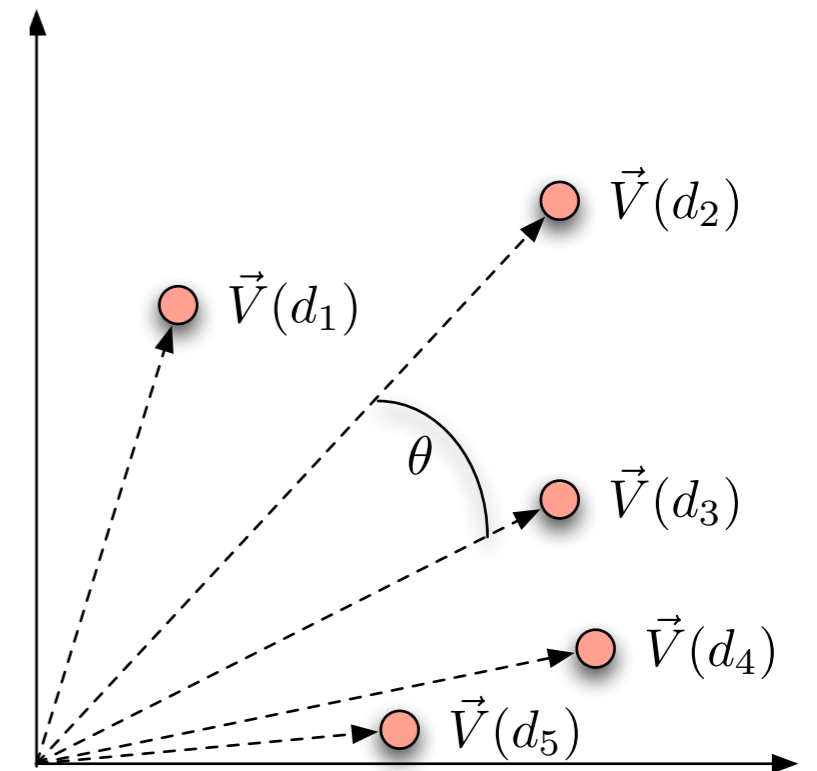


Cosine Similarity Score

- 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 |



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

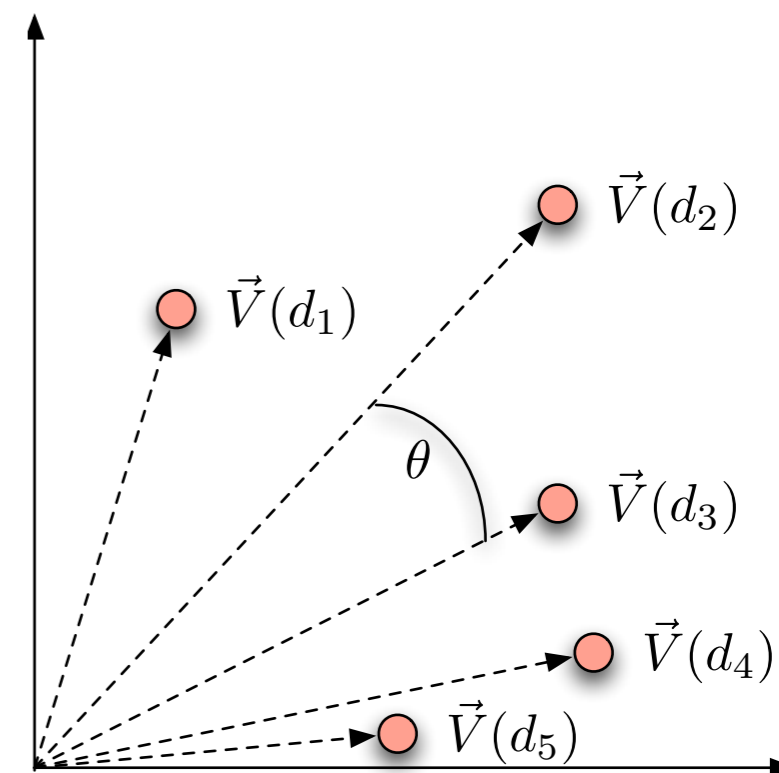


Cosine Similarity Score

- 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 |



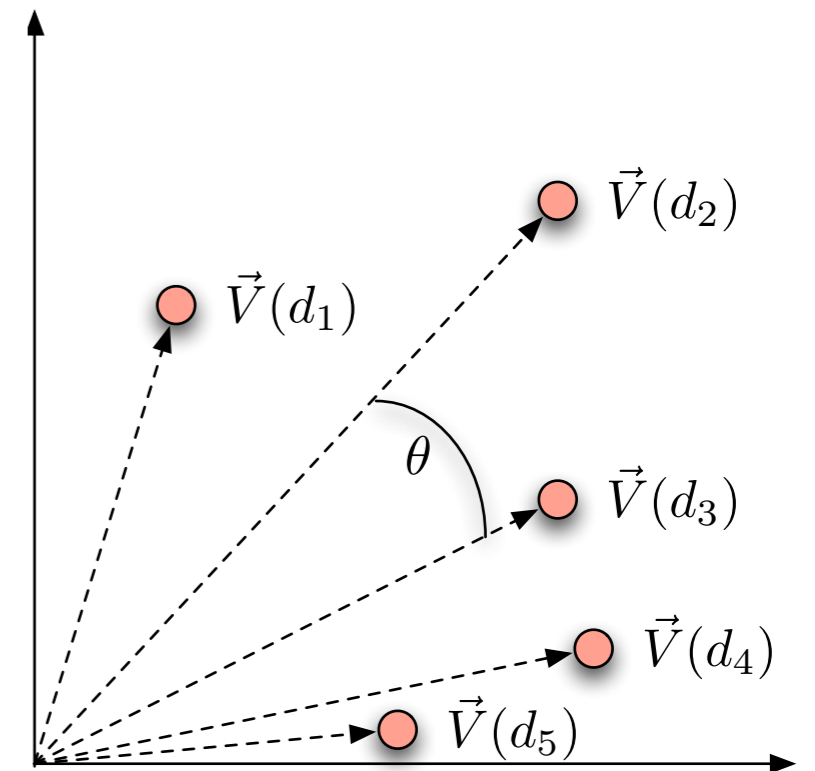
$$|\vec{V}(d_1)| = \sqrt{(11.4 * 11.4) + (8.3 * 8.3) + (2.3 * 2.3) + (11.2 * 11.2)}$$
$$= 18.15$$



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)



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 more effectively discriminate results.



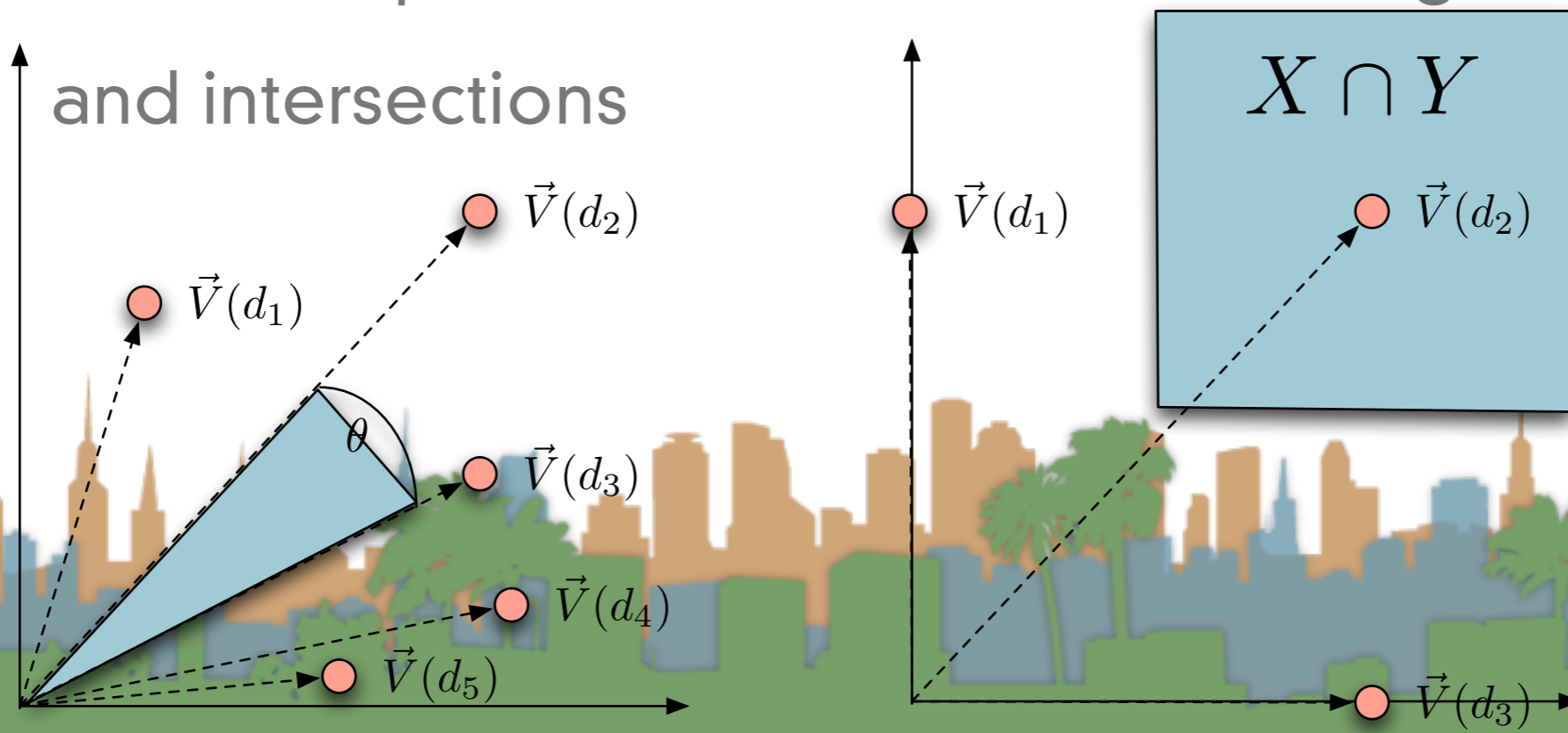
Multiple queries for phrases and vectors

- 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

- 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 idf's 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> | $\log \frac{ corpus }{df_t}$ | <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> | $\max\{0, \log(\frac{ corpus - df_t}{df_t})\}$ | <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?

