









CS 510 Winter 2007















Local Clustering

- Use term co-occurrence data to calculate correlation between terms u, v: C_{uv}
- Add the *n* query terms with the highest $C_{\mu\nu}$ where *u* is a term in the original query

CS 510 Winter 2007

14

16

Reflects term clustering









Weighted by similarity to original query

CS 510 Winter 2007

19

Creating a statistical thesaurus

- · Based on clustering documents
 - Algorithm produces hierarchical clusters
 Uses cosine formula from vector space model
- Classes derived from clusters and 3 parameters
 - Similarity threshold (want tight clusters)
 - Similarity theshold (want tight clus
 - Cluster size (want small clusters)
 Minimum inverse document frequency
 - Only want to use low frequency terms
- Weight the thesaurus classes
 - average term weight/number of terms in class

CS 510 Winter 2007

20











specification rampin

4

.

and an inclusion

×

1.11

Specify

a biologic

concept

Expand

from

category

to instances











