Words vs. Terms

Words vs. Terms

- Information Retrieval cares about "terms"
- You search for 'em, Google indexes 'em
- Query:
 - What kind of monkeys live in Costa Rica?

Words vs. Terms

- What kind of monkeys live in Costa Rica?
 - words?
 - content words?
 - word stems?
 - word clusters?
 - multi-word phrases?
 - thematic content? (this is a "habitat question")

Finding Phrases ("collocations")

- kick the bucket
- directed graph
- iambic pentameter
- Osama bin Laden
- United Nations
- real estate
- quality control
- international best practice
- ... have their own meanings, translations, etc.





Finding Phrases ("collocations")

- Still want to filter out "new companies"
- These words occur together reasonably often but only because both are frequent
- Do they occur more often [among A N pairs?] than you would expect by chance?
 - Expect by chance: p(new) p(companies)

 - Actually observed: p(new companies)
 mutual information = p(new) p(companies | new) mutual information
 - binomial significance test

data from Manning & Schütze textbook (14 million words of NY Times)			
(Pointwise) Mutual Information			
			TOTAL
	new	_new	TOTAL
companies	8	4,667 ("old companies")	4,675
¬companies	15,820	14,287,181 ("old machines")	14,303,001
TOTAL	15,828	14,291,848	14,307,676
 p(new companies) = p(new) p(companies) ? MI = log₂ p(new companies) / p(new)p(companies) = log₂ (8/N) /((15828/N)(4675/N)) = log₂ 1.55 = 0.63 			
 MI > 0 if and only if p(co's new) > p(co's) > p(co's ¬new) Here MI is positive but small. Would be larger for stronger collocations. 			

data from Manning & Schütze textbook (14 million words of NY Times)			
	new	_new	TOTAL
companies	1	583 ("old companies")	584
¬companies	1978	1,785,898 ("old machines")	1,787,876
TOTAL	1979	1,786,481	1,788,460
 Sparse data. In fact, suppose we divided all counts by 8: Would MI change? No, yet we should be less confident it's a real collocation. Extreme case: what happens if 2 novel words next to each other? So do a significance test! Takes sample size into account. 			

data from Manning & Schütze textbook (14 million words of NY Times)			
Binomial Significance ("Coin Flips")			
	new	−new	ΤΟΤΑΙ
companies	8	4,667	4,675
¬companies	15,820	14,287,181	14,303,001
TOTAL	15,828	14,291,848	14,307,676
 Assume we have 2 coins that were used when generating the text. Following new, we flip coin A to decide whether companies is next. Following ¬new, we flip coin B to decide whether companies is next. We can see that A was flipped 15828 times and got 8 heads. Probability of this: p⁸ (1-p)¹⁵⁸²⁰ * ¹⁵⁸²⁸¹/₈₁₁₅₈₂₀₁ We can see that B was flipped 14291848 times and got 4667 heads. Our question: Do the two coins have different weights? (equivalently, are there really two separate coins or just one?) 			

data from Manning & Schütze textbook (14 million w Binomial Significance ("Coin F	rords of NY Times)		
Binomial Significance ("Coin F	lips")		
Binomial Significance ("Coin F	lips")		
newnew	TOTAL		
companies 8 4,667	4,675		
¬companies 15,820 14,287,181 1	14,303,001		
TOTAL 15,828 14,291,848 1	14,307,676		
Null hypothesis: same coin			
• assume $p_{null}(co's new) = p_{null}(co's \neg new) = p_{null}(co's) = 4675/14307676$			
• $p_{null}(data) = p_{null}(8 \text{ out of } 15828) * p_{null}(4667 \text{ out of } 14291848) = .00042$			
Collocation hypothesis: different coins			
assume $p_{coll}(co's new) = 8/15828$, $p_{coll}(co's -new) = 466//14291848$			
• $p_{coll}(data) = p_{coll}(8 \text{ out of } 13828)^{\circ}p_{coll}(4667 \text{ out of } 14291848) = .00081$			
 So collocation hypothesis doubles p(data). 			
 We can sort bigrams by the log-likelihood ratio: log p_{coll}(data)/p_{null}(data) 			
I.e., how sure are we that "companies" is more likely after "new"?			

data from Manning & Schütze textbook (14 million words of NY Times)			
Binomial Significance ("Coin Flips")			
	new	-new	TOTAL
companies	1	583	584
–companies	1978	1,785,898	1,787,876
TOTAL	1979	1,786,481	1,788,460
Null hypothesis: same coin			
assume p _{null} (cors n(data) = n(new) = p _{null} (cos 1 out of 1979)*n	s ¬new) = p _{null} (cos) = (583 out of 1786481)	= 584/1/88460 = 0056
Collocation hypothesis: different coins			
 assume p_{coll}(co's new) = 1/1979, p_{coll}(co's ¬new) = 583/1786481 			
• $p_{coll}(data) = p_{coll}(1 \text{ out of } 1979)*p_{coll}(583 \text{ out of } 1786418) = .0061$			
 Collocation hypothesis still increases p(data), but only slightly now. 			
If we don't have much data, 2-coin model can't be much better at explaining it.			
 Pointwise mutual information as strong as before, but based on much less data. So it's now reasonable to believe the null hypothesis that it's a coincidence 			
600 465 into the Pisner			

data from Manning & Schütze textbook (14 million words of NY Times) Binomial Significance ("Coin Flips")			
	new	-new	TOTAL
companies	8	4,667	4,675
¬companies	15,820	14,287,181	14,303,001
TOTAL	15,828	14,291,848	14,307,676
 Null hypothesis: same coin assume p_{null}(co's new) = p_{null}(co's ¬new) = p_{null}(co's) = 4675/14307676 p_{null}(data) = p_{null}(8 out of 15828)*p_{null}(4667 out of 14291848) = .00042 Collocation hypothesis: different coins assume p_{coll}(co's new) = 8/15828, p_{coll}(co's ¬new) = 4667/14291848 p_{coll}(data) = p_{coll}(8 out of 15828)*p_{coll}(4667 out of 14291848) = .00081 Does this mean that collocation hypothesis is twice as likely? No, as it's far less probable <i>a priorl</i> ! (most bigrams ain't collocations) Bayes: p(coll data) = p(coll) * p(data coll) / p(data) sint twice p(null data) 			















Latent Semantic Analysis

- Themes extracted for IR might help sense disambiguation
- Each word is like a tiny document: (0,0,0,1,0,0,...)
- Express word as a linear combination of themes
- Each theme corresponds to a sense?
 - E.g., "Jordan" has Mideast and Sports themes
 - (plus Advertising theme, alas, which is same sense as Sports) Word's sense in a document: which of its themes are strongest in the document?
- Groups senses as well as splitting them
 - One word has several themes and many words have same theme























