

COLLOCATION EXTRACTION

- *A collocation is an expression consisting of two or more words that correspond to some conventional way of saying things (Manning and Schütze 1999: p151)*
- The difference between comprehensible and natural-sounding language usage
- CAUTION: “collocation” is an ill-defined term!!!

APPLICATIONS

- Information retrieval (query expansion, query segmentation)
- Language modelling in Speech processing (N-grams)
- Parsing (symbolic, statistical)
- Generation (symbolic, statistical)
- Word sense disambiguation (“one sense per collocation” principle)
- Lexicography (e.g. COBUILD)
- Terminology
- Text simplification
- Machine translation (multi-word translation pairs)

EXTRACTION PARADIGMS

- **Segment-based knowledge-driven/statistical extraction:** extract multi-segments as part of segmentation process
- **Word-based, knowledge-driven extraction:** extract word sequences of pre-defined type (e.g. nominal compounds)

POS-based regular expressions, structural analysis

- **Word-based, statistical extraction:** extract statistically idiosyncratic word sequences

STATISTICAL TESTS USED IN COLLOCATION EXTRACTION

- Simple frequency: $f(XY)$
- Pointwise/specific mutual information: $\log \frac{P(x,y)}{P(x)P(y)}$
- Dice's coefficient: $\frac{2 f(x,y)}{f(x)f(y)}$
- (Student's) t score
- (Pearson's) chi-square (χ^2)
- Z score
- Log likelihood
- Selectional association

⋮

BIGRAM RESULTS FROM THE WSJ

Rank	Frequency	Mutual information	χ^2	<i>t</i> test
1	<i>of the</i>	<i>Quadi Doum</i>	<i>Posse Comitatus</i>	<i>of t</i>
2	<i>in the</i>	<i>Wrongful Discharge</i>	<i>LORIMAR TELEPICTURES</i>	<i>in t</i>
3	<i>to NUMB</i>	<i>Seh Jik</i>	<i>Petits Riens</i>	<i>to NU</i>
4	<i>for the</i>	<i>Noo Yawk</i>	<i>Wrongful Discharge</i>	<i>on t</i>
5	<i>to the</i>	<i>WESTDEUTSCHE LANDESBANK</i>	<i>Tupac Amaru</i>	<i>the cor</i>
6	<i>of NUMB</i>	<i>Naamloze Vennootschap</i>	<i>Sary Shagan</i>	<i>about M</i>
7	<i>on the</i>	<i>Caisses Regionales</i>	<i>Outlaw Biker</i>	<i>said</i>
8	<i>NUMB to</i>	<i>Centenaire Blanzly</i>	<i>GEMINI SOGETI</i>	<i>for r</i>
9	<i>that the</i>	<i>Guillen Landrau</i>	<i>Centenaire Blanzly</i>	<i>to r</i>
10	<i>the company</i>	<i>Ea Matsekha</i>	<i>Smith-Corona Typewriters</i>	<i>a sh</i>
⋮				

WHY STATISTICS?

- Pick up on word combinations which occur with “significantly” high relative frequency when compared to the frequencies of the individual words (i.e. $f(x, y)$ as compared to $f(x)$ and $f(y)$)
- BUT WHY SO MANY #\$\$%! STATISTICAL TESTS?
 - complications in evaluation (hard to say which is the “best” test, conflicting results from different researchers)
 - different corpora have different distributional idiosyncracies
 - different tests have different statistical idiosyncracies
- AND WHERE’S THE #\$\$%! LINGUISTICS!!!
 - bear with me!

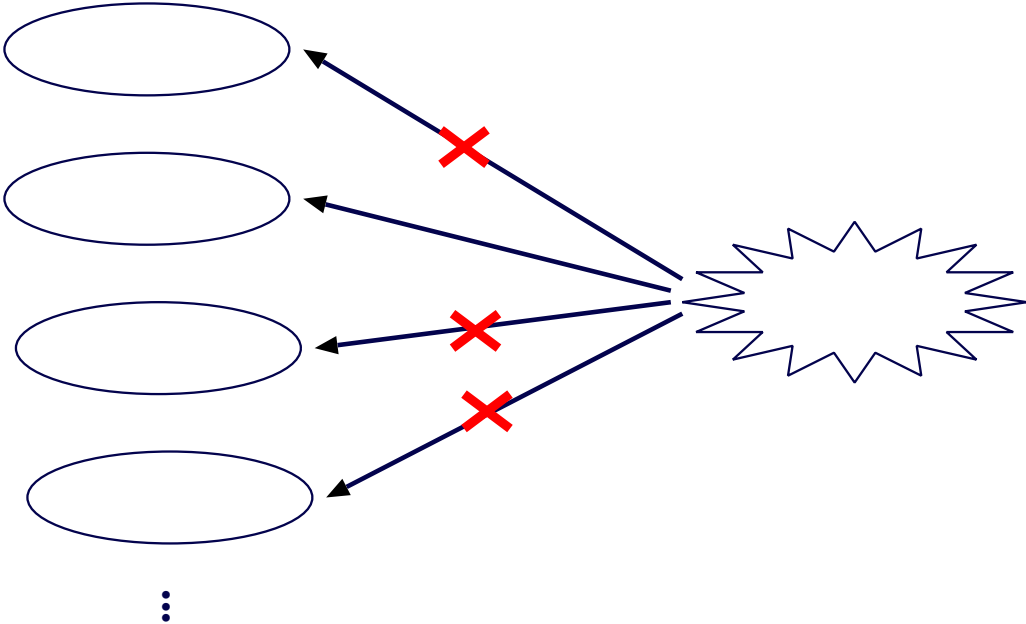
LINGUISTICS IN COLLOCATION EXTRACTION

- Apply statistical measures to (head) bigrams in a given dependency relation (e.g. subject-verb)
 - filters out stop words, produces “collocations” of pre-defined type for direct use in parsing, etc
- Look beyond contiguous bigrams, to bigrams occurring within a “collocational window” of fixed size (e.g. within 3-4 words of each other)
- Utilise linguistic qualities of collocations:
 - limited internal modifiability (applicable as a post-filter)
 - limited substitutability (contrast with anti-collocations, e.g. (*strong* **powerful*) *coffee*)
 - non-compositional semantics

SUBSTITUTABILITY

Lexicalisation

Concept



SUBSTITUTABILITY

- Most immediate means of testing substitutability via synonyms
- Synonyms accessible from thesauri, but word sense disambiguation is generally needed to isolate which synset(s) over which to apply substitution test
- Possibilities of getting at synonyms via distributional analysis (based on dependency pairs)???

CONCEPTUAL IDENTIFIABILITY

Lexicalisation

Concept

