

Semantic Class Induction

Some slides were adapted from the ones created by Ellen
Riloff

Motivation

- A semantic lexicon assigns semantic categories to words.

politician → *human*
truck → *vehicle*
grenade → *weapon*

- Domain-specific vocabulary is often not found in general purpose resources, such as WordNet.
- Automatic methods could be used to enhance these resources or create domain-specific lexicons.

Syntactic Heuristics for Learning Semantic Labels

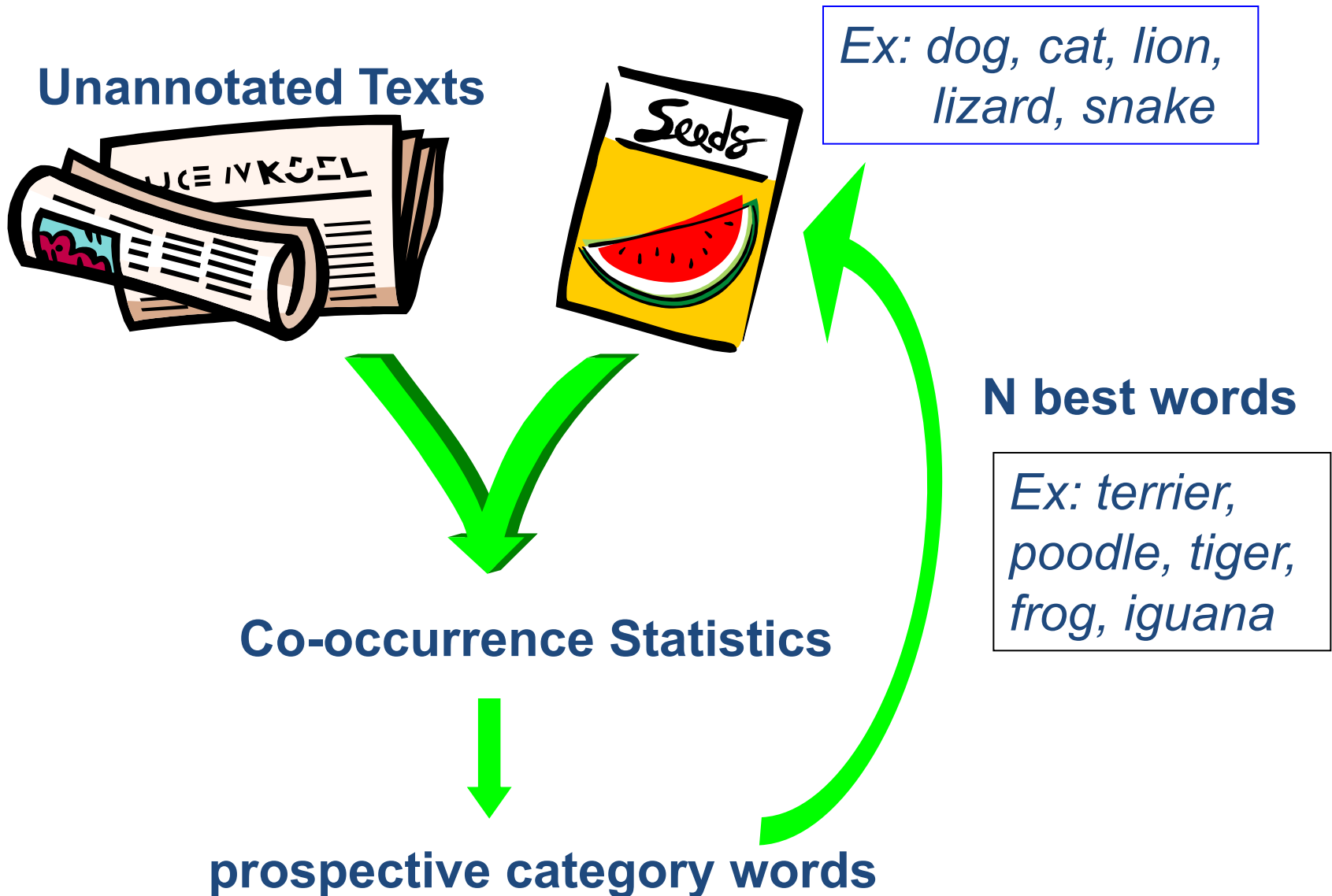
Conjunctions	lions and tigers and bears
Lists	lions, tigers, bears
Appositives	the horse, a stallion
Predicate Nominals	the wolf is a mammal
Compound nouns	tuna fish Honda Sedan

[Riloff & Shepherd 97; Roark & Charniak 98; Phillips & Riloff 02; etc.]

Hyponym patterns	dogs such as beagles and boxers dogs, including beagles and boxers
------------------	---

[Hearst 92; KnowItAll (U.Washington), Kozareva et al. 2008; etc.]

Bootstrapping Semantic Lexicons



Extraction Patterns

- Represent syntactic context that often reveals the semantic class of a word.
- AutoSlog: each pattern extracts an NP from one of 3 syntactic positions: *subject*, *direct object*, *pp obj*.

Some patterns to extract locations:

<subject> was inhabited
patrolling <direct object>
lives in <pp obj>

the locality was inhabited...
...patrolling **Zacamil neighborhood**
...lives in **Argentina**

EXTRACTION PATTERN TYPES

<subject> passive-vp
<subject> active-vp
<subject> active-vp dobj
<subject> active-vp infinitive
<subject> passive-vp infinitive
<subject> auxiliary dobj

<target> was bombed
<perpetrator> bombed
<perpetrator> threw dynamite
<perpetrator> tried to kill
<perpetrator> was hired to kill
<victim> was fatality

active-vp <dobj>
infinitive <dobj>
active-vp infinitive <dobj>
passive-vp infinitive <dobj>
subject auxiliary <dobj>

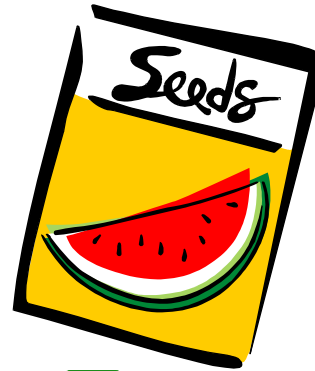
bombed <target>
to kill <victim>
tried to kill <victim>
was hired to kill <victim>
fatality was <victim>

passive-vp prep <np>
active-vp prep <np>
infinitive prep <np>
noun prep <np>

was killed by <perpetrator>
exploded in <target>
to kill with <weapon>
assassination of <victim>

Mutual Bootstrapping [Riloff & Jones 99]

Unannotated Texts



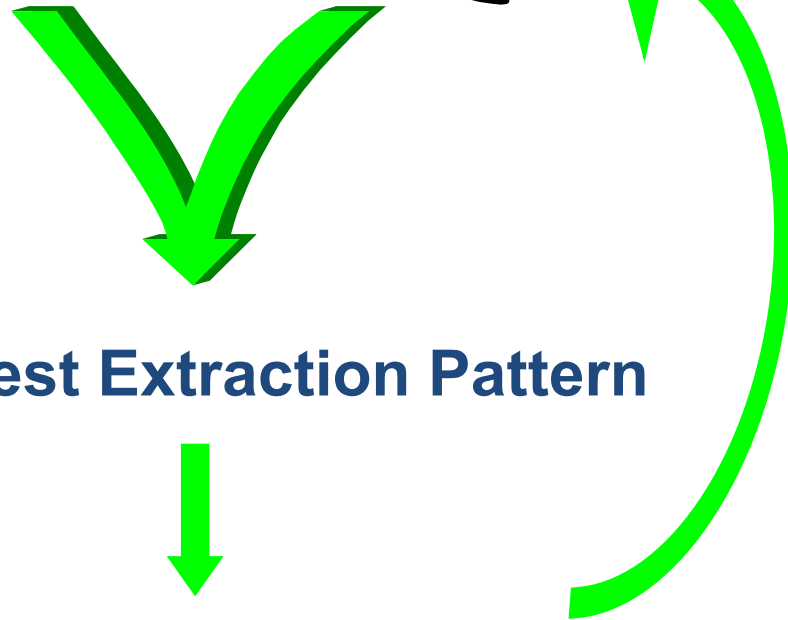
*Ex: dog, cat, lion,
lizard, snake*

*Ex:
<NP> growled*

Best Extraction Pattern

*Ex: Rottweiler,
terrier, cougar*

Extractions (Nouns)



Mutual Bootstrapping Example

SEEDS: Nicaragua, city, region, town

Best pattern: headquartered in <NP>

Extractions: Nicaragua, city, Chapare region, San Miguel

Best pattern: downed in <NP>

Extractions: Nicaragua, city, Usulután region,
San Miguel, area, Soyapango

Best pattern: to occupy <NP>

Extractions: Nicaragua, town, this northern area,
small country, San Sebastian neighborhood,
private property

Examples of Learned Patterns

Location Patterns (Web)

offices in <np>
facilities in <np>
operations in <np>
loans in <np>
operates in <np>
locations in <np>
producer in <np>
states of <np>
seminars in <np>
activities in <np>
consulting in <np>
countries of <np>

Location Patterns (Terrorism)

living in <np>
traveled in <np>
become in <np>
sought in <np>
presidents in <np>
parts of <np>
to enter <np>
condemned in <np>
relations between <np>
ministers of <np>
part in <np>
taken in <np>

Bootstrapping Procedure

- 1. Start from several seed words for a semantic class and an unlabeled text corpus.**
- 2. Score patterns and keep the top N patterns.**
- 3. Score pattern extractions (candidate lexicon words) and select the top M new words as new lexicon words.**
- 4. Increase N by 1 and go back to step 1.**

Pattern Scoring

Every extraction pattern is scored and the best patterns are kept.

The scoring function is:

$$\mathbf{RlogF}(\text{pattern}_i) = \frac{\mathbf{F}_i}{\mathbf{N}_i} * \log_2(\mathbf{F}_i)$$

where:

\mathbf{F}_i is the number of unique category members extracted by pattern_i

\mathbf{N}_i is the total number of unique nouns extracted by pattern_i

Selecting Words for the Lexicon

Score: the average number of category members extracted by each pattern (*while the original algorithm considers all patterns, we'll only consider patterns selected in the previous iteration*) that extracted the candidate word.

$$\text{score}(\text{word}_i) = \frac{\sum_{j=1}^{N_i} F_j}{N_i}$$

$$\text{AvgLog}(\text{word}_i) = \frac{\sum_{j=1}^{N_i} \log_2 (F_j + 1)}{N_i}$$

where:

F_j is the number of unique category members extracted by pattern j ,

N_i is the total number of patterns that extract word i .

More notes

- The mutual bootstrapping approach can be extended to learn semantic lexicon in multiple categories (the Basilisk system).
- Very often, manual review is still necessary to use the learned dictionaries.
- Performance for some categories is beginning to approach levels for which manual review may not be necessary.