

Chapter IR:III

III. Text Transformation

- ❑ Text Statistics
- ❑ Parsing Documents
- ❑ Information Extraction
- ❑ Link Analysis

Information Extraction

Overview

Information extraction is the task of extracting **structured** information from unstructured information sources.

Key goals:

- ❑ Rendering text amenable to structured processing and retrieval.
- ❑ Allowing for logical reasoning and inferences.

Basic analysis process:

- ❑ Input: Token stream
- ❑ Output: Token stream with **annotations / tags** on (subsequences of) tokens.

Methodology:

- ❑ Extraction rules
- ❑ Machine learning
- ❑ Sequence modeling

Information Extraction

Overview

Information extraction is employed to identify more complex index terms by means of natural language processing technology (computationally expensive):

- ❑ **Noun phrases**

Phrases which have a noun as its head word, i.e., a noun and any word that modifies it.
Examples: “*The yellow house* is for sale.”, “I want *a skate board*”.

- ❑ **Named entities**

Words or phrases that designate something (e.g., a place, a person, an organization, etc.).

- ❑ **Coreference resolution**

Coreferences, i.e., anaphora and cataphora, are expressions that refer backward or forward in a text, respectively. Resolving *them* is important for text understanding, yet, one of the most difficult problems of natural language processing

- ❑ **Relation detection**

Extraction of relations between named entities mentioned in the text. Example: “*Bill* lives in the *USA*”.

- ❑ **Semi-structured information extraction**

Extraction of tables, quotes, references, comments, etc.

Information Extraction

Part-of-Speech Tagging

Given a token sequence of text, markup each token with its part of speech (POS).

Common English (Western) parts of speech:

- ❑ **Noun** (names of abstract or concrete entities: persons, places, things, ideas, qualities)
- ❑ **Pronoun** (substitutes for nouns)
- ❑ **Verb** (actions, occurrences, or states of being)
- ❑ **Adjective** (modifiers of a noun or pronoun)
- ❑ **Adverb** (modifiers of verbs, adverbs, or adjectives)
- ❑ **Preposition** (words expressing relations in a phrase or sentence)
- ❑ **Conjunction** (connects words, phrases, or clauses)
- ❑ **Interjection** (expressions of feelings and emotions)
- ❑ **Determiner** (markers of definiteness or indefiniteness)

For practical purposes, these broad classes are insufficient. Typically some 30 to 150 parts of speech are distinguished.

Information Extraction

Part-of-Speech Tagging: Example

Original text:

A relevant document will describe marketing strategies carried out by U.S. companies for their agricultural chemicals, report predictions for market share of such chemicals, or report market statistics for agrochemicals, pesticide, herbicide, fungicide, insecticide, fertilizer, predicted sales, market share, stimulate demand, price cut, volume of sales.

Brill tagger:

A/**DT** relevant/**JJ** document/**NN** will/**MD** describe/**VB** marketing/**NN** strategies/**NNS** carried/**VBD** out/**IN** by/**IN** U.S./**NNP** companies/**NNS** for/**IN** their/**PRP\$** agricultural/**JJ** chemicals/**NNS** ,/, report/**NN** predictions/**NNS** for/**IN** market/**NN** share/**NN** of/**IN** such/**JJ** chemicals/**NNS** ,/, or/**CC** report/**NN** market/**NN** statistics/**NNS** for/**IN** agrochemicals/**NNS** ,/, pesticide/**NN** ,/, herbicide/**NN** ,/, fungicide/**NN** ,/, insecticide/**NN** ,/, fertilizer/**NN** ,/, predicted/**VBN** sales/**NNS** ,/, market/**NN** share/**NN** ,/, stimulate/**VB** demand/**NN** ,/, price/**NN** cut/**NN** ,/, volume/**NN** of/**IN** sales/**NNS** ./.

CC	coordinating conjunction	NN	singular or mass noun	VBD	verb, past tense
DT	singular determiner/quantifier	NNP	proper noun, singular	VBN	verb, past participle
IN	preposition	NNS	plural noun	,	comma
JJ	adjective	PRP\$	possessive pronoun	.	dot
MD	modal auxiliary	VB	verb, base form		other tags

Information Extraction

Part-of-Speech Tagging: Brill Tagger [\[Brill 1992\]](#)

Principle: “error-driven transformation-based tagging”

1. Assign each token its most likely part of speech tag. Stemming rules are applied to match inflected tokens with word stems stored in a dictionary.
2. Apply a list of transformation rules to correct tagging errors.
3. Repeat Step 2 until no rules can be applied, anymore, or after a pre-sepcified number of repetitions.

Concepts:

- ❑ Initial tag probabilities are trained on a large pre-tagged corpus.
- ❑ Rules are learned from errors made on a pre-tagged corpus, and applied in the order listed.
- ❑ Rules are defined as follows: $T1 \ T2 \ \langle \text{Premise} \rangle$

Semantics:

For each token currently tagged with $T1$ which fulfills the $\langle \text{Premise} \rangle$, replace $T1$ with $T2$.

Information Extraction

Part-of-Speech Tagging: Brill Tagger [\[Brill 1992\]](#)

Premises:

context x	A word in context is tagged x.
property	The word has a certain property.
context property	A word in context has a certain property.
context	One or any of $i \in [1, 3]$ preceding or following word(s).
property TRUE FALSE	Capitalized word.

Example rules:

TO	IN	next-tag AT	TO: to, IN: preposition, AT: article
VTB	VBD	prev-word-is-cap TRUE	VTB: verb past participle,
VBD	VTB	prev-1-or-2-or-3-tag HVD	VBD: verb past tense, HVD: auxiliary had
VB	NN	prev-1-or-2-tag AT	VB: verb, base form
NN	VB	prev-tag TO	NN: noun, singular or mass
TO	IN	next-word-is-cap TRUE	
NN	VB	prev-tag MD	MD: modal

Rules are learned starting with the initial tagging on a training dataset by instantiating rules from the above templates, keeping those that minimize tagging errors the most in each iteration, until some termination criterion is reached.

Information Extraction

Part-of-Speech Tagging: Brill Tagger [\[Brill 1994\]](#)

Problem: The tagger cannot tag words not occurring in the training data.

An unknown word tagger can be trained based on the same principles but with different premises as templates for rules. T1 may be UNK for unknown.

Premises:

affix x constraint
context word
char x

Token fulfills constraint regarding affix of at most 4 chars.
A word appears in context.
Character x occurs in word.

constraint

When deleting or adding affix x, word found in dictionary.
Else, affix x occurs in token.

Example rules:

NN	NNS	suffix -s occurs
NN	CD	char .
NN	JJ	char -
NN	VBN	suffix -ed occurs
NN	VBG	suffix -in occurs
UNK	JJ	suffix -ly addition
UNK	RB	suffix -ly occurs

NN: noun, singular or mass, NNS: noun, plural
CD: cardinal number
JJ: adjective
VBN: verb, past participle
VBG: verb, gerund or present participle
UNK: unknown
RB: adverb

Remarks:

- ❑ Large corpora for part of speech tagging have been painstakingly manually annotated, starting with the 1 million word Brown corpus in the 1960s, later superseded by the 100 million word British National Corpus, and others.
- ❑ Tag sets: [Brown](#) (87 tags), [Penn TreeBank II](#) (41 tags), [British National Corpus](#) (61 tags), [Oxford English part-of-speech tagset](#) (109 tags), [British National Corpus Sampler](#) (146 tags).
- ❑ Assigning the most probable tag to each known word and proper noun (NNP) to all unknown words already yields 90% accuracy. [[Charniak 1997](#)]
- ❑ The state of the art in part of speech tagging can be reviewed at [aclweb.org](#). Most taggers reported are based on statistical sequence models rather than rules. However, many taggers proposed are not included, including the Brill tagger.
- ❑ Nevertheless, the Brill tagger frequently serves as baseline for comparison, and as a last step in tagging pipelines.

Information Extraction

Noun Phrase

A phrase is a group of words (or possibly a single word) that functions as a constituent in the syntax of a sentence. A noun phrase (NP) is a phrase which has a noun (or indefinite pronoun) as its head word, preceded and/or followed by modifying words. Hence, a noun phrase is a phrase which acts like a noun.

Common noun modifiers:

- ❑ Determiners / articles (the inclusion of determiners is disputed)
- ❑ Adjectives
- ❑ Prepositional phrases

Examples: **What are noun phrases in the examples?**

- ❑ Cats sleep a lot.
- ❑ A cat is sleeping.
- ❑ The fluffy, long-haired cat sleeps.
- ❑ The cat on top of the stool is sleeping.
- ❑ Most big cats sleep at daytime.

Information Extraction

Noun Phrase

A phrase is a group of words (or possibly a single word) that functions as a constituent in the syntax of a sentence. A noun phrase (NP) is a phrase which has a noun (or indefinite pronoun) as its head word, preceded and/or followed by modifying words. Hence, a noun phrase is a phrase which acts like a noun.

Common noun modifiers:

- ❑ **Determiners / articles** (the inclusion of determiners is disputed)
- ❑ **Adjectives**
- ❑ **Prepositional phrases**

Examples:

- ❑ *Cats* sleep a lot.
- ❑ **A** *cat* is sleeping.
- ❑ **The fluffy, long-haired** *cat* sleeps.
- ❑ **The** *cat* **on top of the stool** is sleeping.
- ❑ **Most big** *cats* sleep at daytime.

Information Extraction

Noun Phrase Extraction

Noun phrases can be extracted from POS tagged text using regular expressions.

To extract noun phrases as index terms, simplified patterns are typically used:

- ❑ sequences of nouns
- ❑ sequences of adjectives followed by nouns

Example (noun phrases bold, tagging errors **fixed**):

A/DT **relevant/JJ document/NN** will/MD describe/VB **marketing/NN strategies/NNS** carried/VBD out/IN by/IN **U.S./NNP companies/NNS** for/IN their/PRP\$ **agricultural/JJ chemicals/NNS** ,/, report/VB predictions/NNS for/IN **market/NN share/NN** of/IN such/DT chemicals/NNS ,/, or/CC report/VB **market/NN statistics/NNS** for/IN agrochemicals/NNS ,/, pesticide/NN ,/, herbicide/NN ,/, fungicide/NN ,/, insecticide/NN ,/, fertilizer/NN ,/, **predicted/JJ sales/NNS** ,/, **market/NN share/NN** ,/, stimulate/VB demand/NN ,/, **price/NN cut/NN** ,/, volume/NN of/IN sales/NNS ./.

Compared to simply indexing n -grams, extracting noun phrases incurs high processing costs for POS tagging while offering the potential to save storage space for n -grams that would never be requested.

Information Extraction

Named Entity

A named entity refers to a real-world object that can be denoted with a proper name.

Named entities are phrases that form a subset of noun phrases, adjective phrases, or take the role of a noun phrase but have different internal structure.

Three basic object categories:

- ❑ Names of people, organizations, locations, facilities, products, events, natural objects, diseases, colors
- ❑ Times
- ❑ Numbers

Example: **What are the named entities?**

Fred Smith, living at 10 Water Street, Springfield, MA, has been breeding five species of tropical fish for the past 15 years.

Information Extraction

Named Entity

A named entity refers to a real-world object that can be denoted with a proper name.

Named entities are phrases that form a subset of noun phrases, adjective phrases, or take the role of a noun phrase but have different internal structure.

Three basic object categories:

- ❑ **Names** of people, organizations, locations, facilities, products, events, natural objects, diseases, colors
- ❑ **Times**
- ❑ **Numbers**

Example:

Fred Smith, living at **10 Water Street, Springfield, MA**, has been breeding **five species of tropical fish** for **the past 15 years**.

Information Extraction

Named Entity Recognition

Approaches to construct a named entity recognizer:

- ❑ **Dictionary-based recognition**

Compile lists of words for each named entity category. If a (sequence of) token in a text is found in a text, or if it can be matched with a word from the dictionary, it is tagged as a named entity.

- ❑ **Rule-based recognition**

Handcrafted or learned pattern rules are defined to recognize named entities. Some aspects of learning rules can be transferred from the Brill tagger.

- ❑ **Statistical sequence modeling**

Application of probabilistic sequence models, such as hidden Markov models, to obtain named entities from sentences.

For every named entity category (or combinations thereof), a sufficiently large set of manually annotated training documents is required.

Domain transfer of trained recognizers between entity categories is difficult.

Information Extraction

Named Entity Recognition: Hidden Markov Models (informal)

Hypothesis: The decision whether a token in a sequence is part of a named entity can be made from tokens in its context. Example: `marathon` is part of a named entity when preceded by `Boston`.

A Markov model describes a process as a collection of states with transition probabilities between them. If state transitions depend only on the current state, the process modeled is said to have the Markov **property**.

Markov **assumption** for text: Given a token sequence, the next token depends only on the last one of the sequence. Generalization: the last n tokens.

Problem: We want to know to what category of named entities a token belongs, not what the next most likely token is.

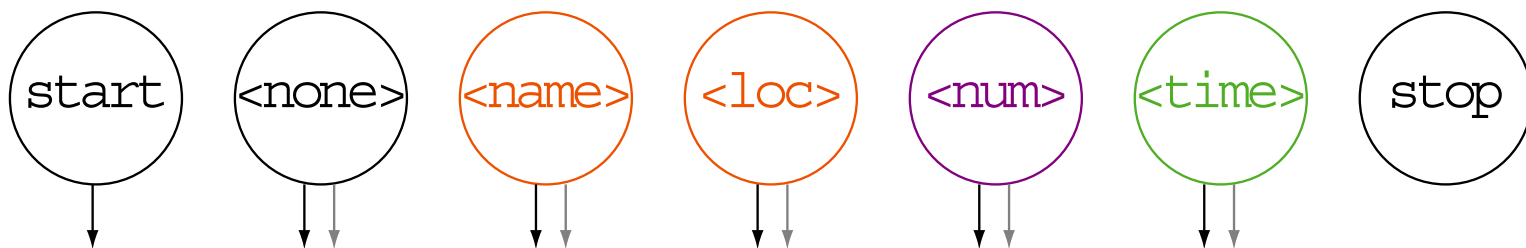
Information Extraction

Named Entity Recognition: Hidden Markov Models (informal)

In a hidden Markov model, the state transitions cannot be observed. But every state may produce an output with a certain probability.

Example hidden state transitions:

```
start <name> <name> <none> <none> <none> <loc> <loc> <loc> <loc> <loc>
<loc> <loc> <none> <none> <none> <none> <num> <num> <num> <num> <num>
<none> <time> <time> <time> <time> stop
```



to any state except start with prior probability

output of word with probability

Example output:

Fred/**NNP** Smith/**NNP**,/, living/**VBG** at/**IN** 10/**CD** Water/**NNP** Street/**NNP**,/,
Springfield/**NNP**,/, MA/**NNP**,/, has/**VBZ** been/**VBN** breeding/**VBG** five/**CD**
species/**NNS** of/**IN** tropical/**JJ** fish/**NN** for/**IN** the/**DT** past/**JJ** 15/**CD**
years/**NNS**./.

Information Extraction

Named Entity Recognition: Hidden Markov Models (informal)

Given a trained hidden Markov model and a token sequence, the question is which sequence of hidden state transitions most likely produces the given token sequence as an output.

Based on large-scale training data, transition and output probabilities are trained.

Using dynamic programming algorithms, such as the Viterbi algorithm or the Forward-backward algorithm, the above question can then be answered.

Rule-based approaches are competitive to sequence modeling approaches, whereas performance depends on the category of named entity.

>90% accuracy for names, locations, organizations require about 1 million token of training data (1500 news articles).