

Dependency Analysis of Japanese Spoken Language via SVM

Kiyonori Ohtake

ATR Spoken Language Translation Research Laboratories
August 6, 2003

Dependency Analysis of Japanese Spoken Language via SVM - p. 1/24

Background

- Dependency analysis for Japanese is important.
- Many works for **text** have been done.
- Few analyzers for **spoken** Japanese.
 - difficult task?
 - hard to develop analyzed spoken language corpora?

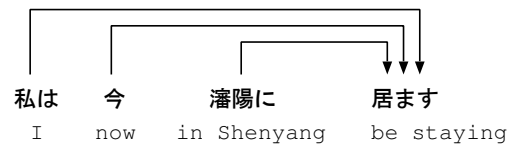
Dependency Analysis of Japanese Spoken Language via SVM - p. 2/24

Objective

- To build a dependency analyzer for spoken language with corpus-based method.
- To prepare a consistent dependency analyzed corpus of spoken language with currently available resources.

Dependency Analysis of Japanese Spoken Language via SVM - p. 3/24

Japanese Dependency Structure



1. From left to right.
2. Do not cross each other.
3. One segment depends on only one segment.

Dependency Analysis of Japanese Spoken Language via SVM - p. 4/24

Dependency Analysis Approach

- Rule-based
- Corpus-based using ME, **SVM**, etc.

Why SVM?

very high generalization ability



it works well if the learning corpus is small.

Dependency Analysis of Japanese Spoken Language via SVM - p. 5/24

Resources

- Analyzers
 - KNP: rule-based analyzer
 - CaboCha: analyzer with SVM
- Analyzed corpora
 - KUTC: newspaper articles, directly applicable to CaboCha
 - JDEP: conversations, **not** applicable to CaboCha

Dependency Analysis of Japanese Spoken Language via SVM - p. 6/24

CaboCha: analyzer via SVM

Characteristics:

- Cascaded chunking model [Kudo and Matsumoto 2002]
- SVM tells whether a segment depends on the next one or not
- Available analyzed corpus is **written text** (KUTC)

Accuracy = 89% (by 2-fold cross validation for newspapers)

Dependency Analysis of Japanese Spoken Language via SVM - p. 7/24

KNP: rule-based analyzer

Characteristics:

Detecting conjunctive structures, and then analyze with manually built rules.

Accuracy = 91% (for newspapers)

Dependency Analysis of Japanese Spoken Language via SVM - p. 8/24

Corpus - KUTC

KUTC: Kyoto University Text Corpus

- newspaper articles, 40,000 sentences
- analyzed with KNP and manually corrected
- dependency segment: *bunsetu*

bunsetu: a phrasal unit that consists of one or more morphemes.

Dependency Analysis of Japanese Spoken Language via SVM - p. 1024

Corpus - JDEP

JDEP= Japanese Dependency structures of ATR SLDB

- travel-type conversations, 21,761 utterances
- automatically analyzed and manually corrected
- dependency segment: morpheme

Converted into KUTC format

Dependency Analysis of Japanese Spoken Language via SVM - p. 1024

Comparing two analyzers and two corpora

	KUTC	
	Dep. acc.	Sent. acc.
KNP	91.23%	57.17%
CaboCha	97.36%	83.40%
	JDEP	
	Dep. acc.	Sent. acc.
KNP	84.27%	69.54%
CaboCha	88.43%	76.01%

Dependency Analysis of Japanese Spoken Language via SVM - p. 1124

Corpus Cleaning

- Objective: to remove **inconsistencies**
- How: repeating closed test with CaboCha and correct manually. very **SIMPLE**

Standard guideline: based on KUTC annotation standard, and arranged for spoken language.

Dependency Analysis of Japanese Spoken Language via SVM - p. 1024

Gain by Cleaning

Evaluated by 2-fold cross validation

	Accuracy	
	KUTC	JDEP
Before	88.85%	92.78%
After	89.92%	93.03%

Dependency Analysis of Japanese Spoken Language via SVM - p. 1324

Discussions

- Why JDEP results were higher?
- Whether *bunsetu* is good segment for spoken language dependency analyzing?

Dependency Analysis of Japanese Spoken Language via SVM - p. 1424

Discussions

- Why JDEP results were higher? Sentences in JDEP are shorter.
- Whether *bunsetu* is good segment for spoken language dependency analyzing?

Dependency Analysis of Japanese Spoken Language via SVM - p. 1424

Discussions

- Why JDEP results were higher? Sentences in JDEP are shorter.
- Whether *bunsetu* is good segment for spoken language dependency analyzing?
⇒ No. Need to investigate further.

Dependency Analysis of Japanese Spoken Language via SVM - p. 1424

Open problems

- Much learning time.
- To prepare properly segmented corpus.
- How to analyze real spoken language.

Dependency Analysis of Japanese Spoken Language via SVM - p. 15/24

Conclusion

- Introduced available resources, and built an analyzer via SVM.
- Carried out the corpus cleaning.
- Cleaning raised the acc. (1.07% for KUTC, 0.25% for JDEP).
- The analyzer resulted 93% acc. for spoken language.

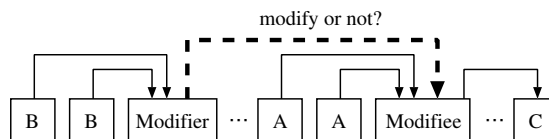
Dependency Analysis of Japanese Spoken Language via SVM - p. 16/24

Thank you!

Dependency Analysis of Japanese Spoken Language via SVM - p. 17/24

Appendix - Features for SVM

Dynamic features:



Static features: head words, parts-of-speech, functional words and inflection forms, and so on.

Dependency Analysis of Japanese Spoken Language via SVM - p. 18/24

Appendix - cascaded chunking

私は 今 中国の 瀋陽に 居ます
O O O O O

Step 1. Initialization

O tag means the dependency relation is undecided.

Dependency Analysis of Japanese Spoken Language via SVM - p. 19/24

Appendix - cascaded chunking

私は 今 中国の 瀋陽に 居ます
O O D D O

Step 2. Each **O** tag depends on immediate right hand side?
⇒ replace with **D**

Dependency Analysis of Japanese Spoken Language via SVM - p. 20/24

Appendix - cascaded chunking

私は 今 中国の 瀋陽に 居ます
O O D D O
delete

Step 3. Delete segments with **D** tag that immediately follow **O** tag.

Dependency Analysis of Japanese Spoken Language via SVM - p. 21/24

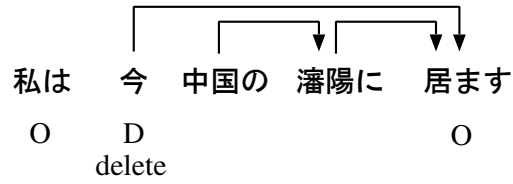
Appendix - cascaded chunking

私は 今 中国の 瀋陽に 居ます
O O D O
delete

Repeating Step 2 and 3.

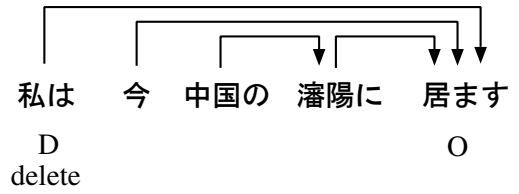
Dependency Analysis of Japanese Spoken Language via SVM - p. 22/24

Appendix - cascaded chunking



Repeating Step 2 and 3.

Appendix - cascaded chunking



Repeating Step 2 and 3.
And finished.