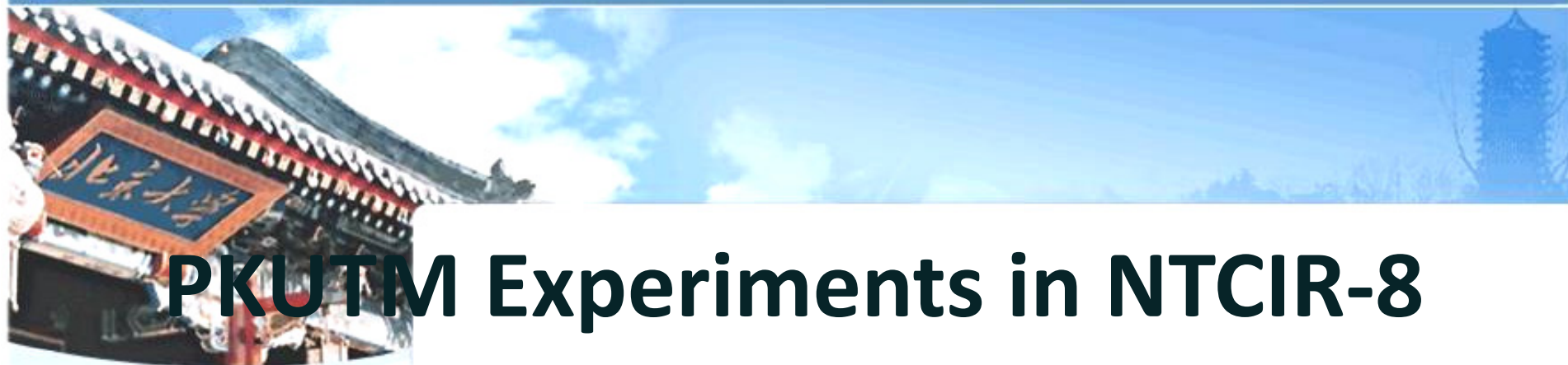




北京大学计算机科学技术研究所
INSTITUTE OF COMPUTER SCIENCE & TECHNOLOGY OF PEKING UNIVERSITY



PKUTM Experiments in NTCIR-8 MOAT Task

Author: Chenfeng Wang*, Tengfei Ma*, Liqiang Guo, Xiaojun Wan
and Jianwu Yang

Affiliation: Institute of Computer Science & Technology of Peking
University

Speaker: Tengfei Ma



Background of Opinion Analysis

- Aspects of Opinion Analysis
 - Is it opinionated?
 - Is the opinion positive or negative?
 - What is the opinion?
 - Who gives the opinion and who does the opinion point to?
 - How to summarize all the opinions?
 - ...



Background of Opinion Analysis



NTCIR-8 MOAT

- Aspects of Opinion Analysis
 - Is it opinionated?
 - Is the opinion positive or negative?
 - What is the opinion?
 - Who gives the opinion and who does the opinion point to?
 - How to summarize all the opinions?
 - ...



Background of Opinion Analysis

- The trend of opinion analysis
 - Coarse-grain to fine-grain
 - Holder/target extraction
 - General to domain-specific and domain-transfer
 - Opinion analysis in news, product reviews, movie reviews
 - Cross-Lingual, transfer learning
 - Publisher-predominate to interactive



Background of Opinion Analysis



- The trend of opinion analysis
 - Coarse-grain to fine-grain
 - Holder/target extraction
 - General to domain-specific and domain-transfer
 - Opinion analysis in news, product reviews, movie reviews
 - Cross-Lingual, transfer learning
 - Publisher-predominate to interactive

NTCIR-8 MOAT



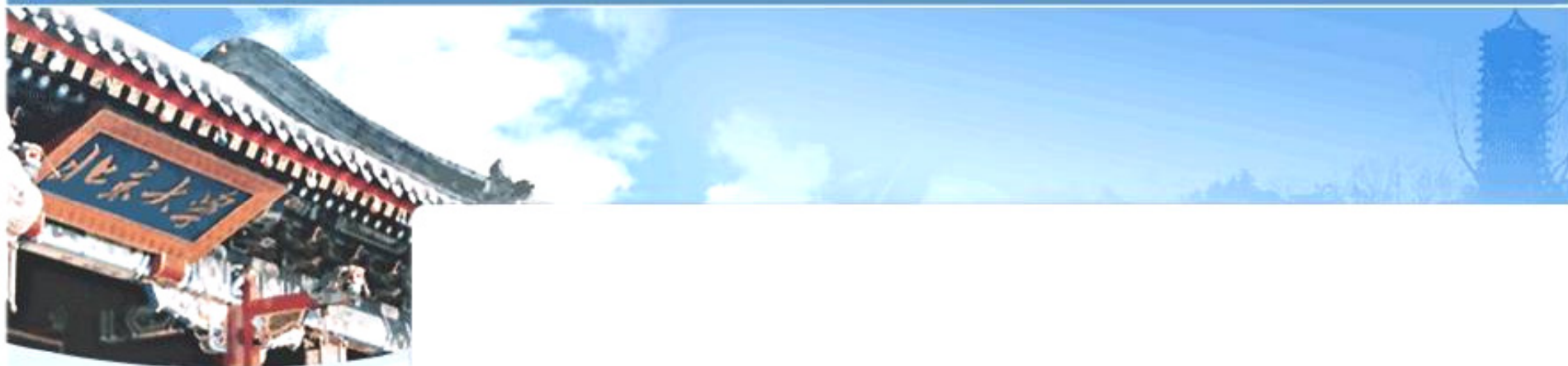
Our tasks in NTCIR-8 Moat



- Opinionated subtask.
- Opinion holder extraction.
- Opinion target extraction.



北京大学计算机科学技术研究所
INSTITUTE OF COMPUTER SCIENCE & TECHNOLOGY OF PEKING UNIVERSITY



(Opinionated task)

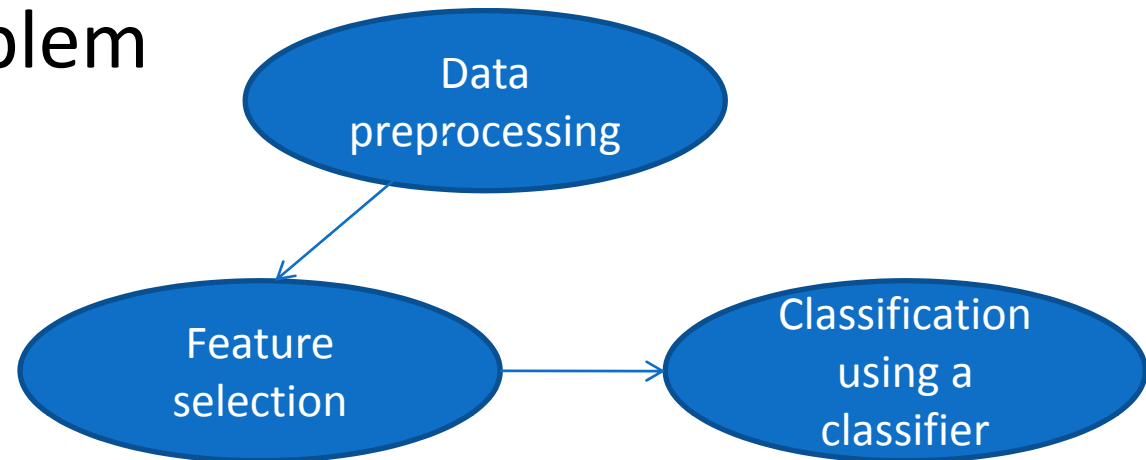
I. DETECTION OF SUBJECTIVE SENTENCES



Detection of subjective sentences



- Equivalent to a classification problem



- Our method:
 - Some combined datasets
 - Some special opinion features
 - A general classifier and an improvement



Detection of subjective sentences

- Data Preprocessing
 - Choosing the training Datasets
 - NTCIR6/NTCIR7 corpora and NTCIR8's samples
 - Containing both simplified and traditional Chinese
 - Translate traditional Chinese to Simplified Chinese
 - POS, NER
 - Building Lexicons



Building Lexicons

– Source:

- expanded Hownet by using the Synonymy Thesaurus + MPQA(English--->Chinese) + NTU + our in-house labeled corpora

– Types:

- **Opinion Operators** e.g.声称(claim)
- **Opinion Indicators** e.g.但是(but)
- **Degree Adverbs** e.g.非常(very), 缺乏(lack of)
- **Opinion Words** (28421 opinion words)
- **Strong Opinion Words** (6471 words)



Detection of subjective sentences



- Feature Selection

Punctuations Features
Presence of quotation marks like “, [,] and ”
Presence of colon followed by quotation marks
Percentage of punctuations in sentences
Words and Entities Features
The percentage of numeral words
The presence of pronoun
The presence of a named entity
The presence of a word which indicates a sequence
Lexical Subjective Clues
The presence of opinion operator
The presence of opinion indicator
The logarithm of percentage of opinion words
The logarithm of percentage of strong opinion words
The presence of degree verb
Collocation Features
The presence of collocations between named entities and opinion operators
The presence of collocations between pronouns or nouns and opinion operators
The presence of collocations between opinion operators and opinion words
The presence of collocations between pronouns and opinion words
The presence of collocations between nouns or pronouns and opinion words
The presence of collocations between degree adverbs and opinion operators
The presence of collocations between degree adverbs and opinion words
The presence of collocations between nouns or named entities and opinion words



Detection of subjective sentences



- Classifier
 - Basic classifiers
 - such as SVM, Naive Bayes, Max Entropy and Decision Tree
 - The comparison is shown in the following section
 - Improved classifier
 - Iterative classifier using former results of detecting subjective sentences



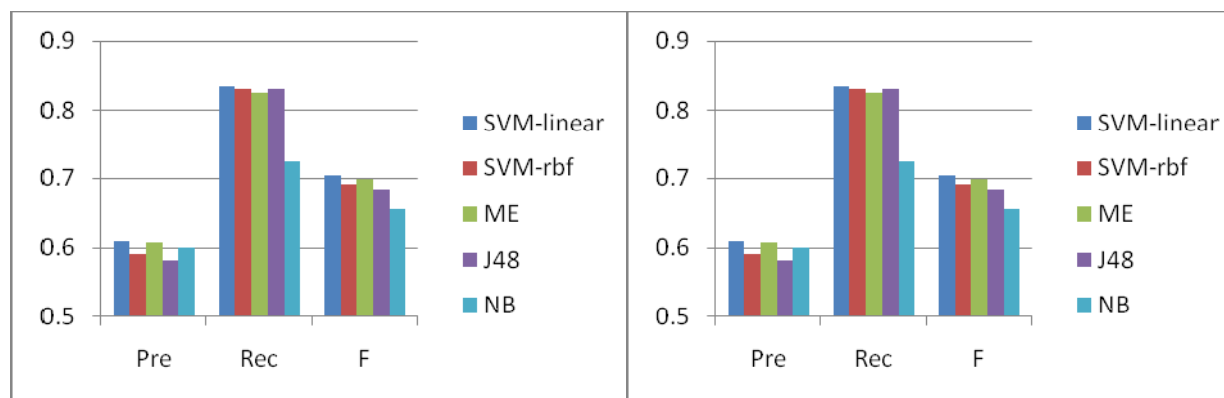
Detection of subjective sentences



- Results in NTCIR8

		Precision	Recall	F-measure
Run1	All datasets +iterative classifier	0.3721	0.8370	0.5152
Run2	NTCIR7 + NTCIR8 simplified Chinese + basic classifier	0.4134	0.8335	0.5527
Run3	Run2 + NTCIR7 traditional dataset	0.3405	0.9062	0.4950

- Additional Tests (Comparison of different classifiers)



lenient

strict

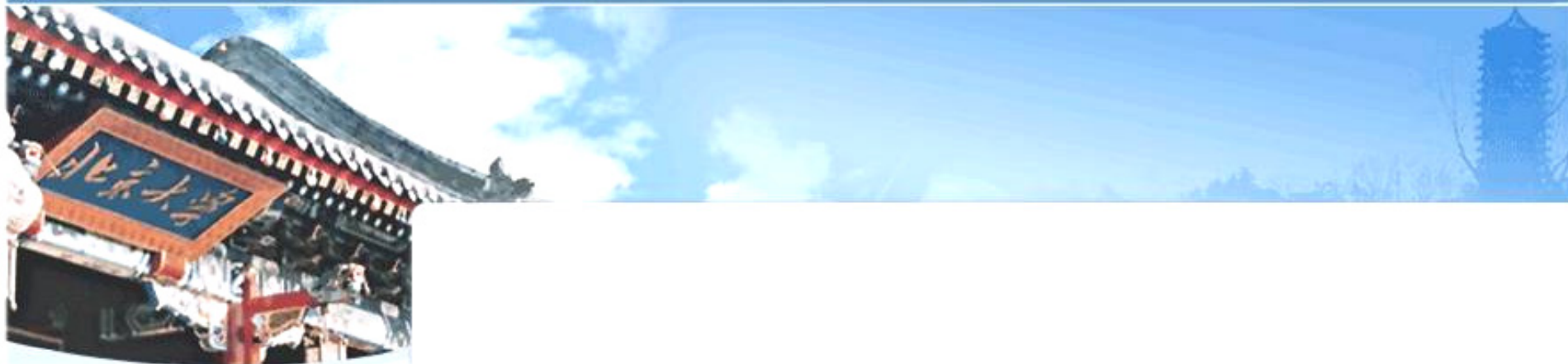


Detection of subjective sentences

- Discussion of the results
 - Training data
 - More \neq Better
 - When and how to leverage translated datasets
 - Classifier
 - Iterative \rightarrow risk
 - Problem
 - Ambiguous definition
 - Ambiguous words



北京大学计算机科学技术研究所
INSTITUTE OF COMPUTER SCIENCE & TECHNOLOGY OF PEKING UNIVERSITY



Holder/target task

EXTRACTION OF OPINION HOLDERS AND TARGETS



Extracting opinion holders/targets



- Common methods
 - Parsing and direct training (Bethard)
 - Maximum Entropy ranking (Kim and Hovy)
 - Labeling
- Our method
 - Chunking and heuristic rules



Extracting opinion holders/targets



- Advantage of Chunking
 - Better than parsing in Chinese
 - Easier to control and modify than shallow parsing
- Process:
 - Training data: proposition bank
 - Modifying training data
 - Training and labeling by CRF



Extracting opinion holders/targets



- Heuristic rules for opinion holder extraction
 - before an opinion operator (include a colon) or following a quotes.
 - not governed by a preposition
 - in other sentences sometimes
 - using nouns or pronouns as candidates to complement the upper missing cases
 - author



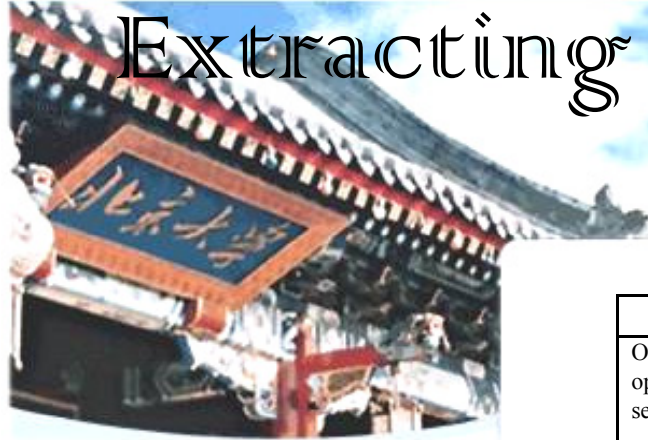
Extracting opinion holders/targets



- Heuristic rules for opinion target extraction
 - Similar to opinion holder extraction
 - Mainly existing in the opinion clause or as the object of an opinion operator
 - Coherent with neighbor sentences

Table 3. Evaluations Results for Opinion Holders

Table 4. Evaluations Results for Opinion Targets



Extracting opinion holders/targets

Holder Extraction

		Precision	Recall	F-measure
Only for opinionated sentences	Run1	0.550	0.434	0.485
	Run2	0.554	0.431	0.485
	Run3	0.548	0.473	0.508
For all sentences	Run1	0.204	0.434	0.277
	Run2	0.232	0.431	0.301
	Run3	0.186	0.473	0.267

Target Extraction

		Precision	Recall	F-measure
Only for opinionated sentences	Run1	0.892	0.736	0.806
	Run2	0.896	0.732	0.805
	Run3	0.877	0.792	0.832
For all sentences	Run1	0.339	0.736	0.464
	Run2	0.385	0.732	0.504
	Run3	0.307	0.792	0.442



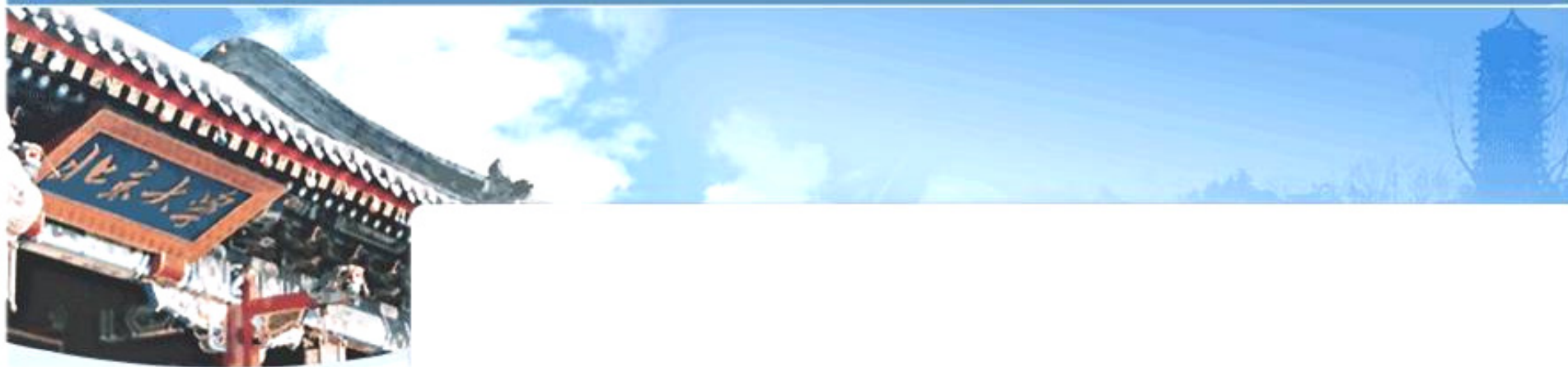
Extracting opinion holders/targets



- Discussion
 - Limited by the parsing technique
 - Features are complex for machine learning
 - Future research (See (Ma, Coling10))
 - Adding semantic information
 - Adding syntactic rules to leverage relevant information (e.g. reviews--news)



北京大学计算机科学技术研究所
INSTITUTE OF COMPUTER SCIENCE & TECHNOLOGY OF PEKING UNIVERSITY



Thank you~

Any questions?