# An Integrated Approach to Automatic Synonym Detection in Turkish Corpus

Dr. Tuğba YILDIZ Assist. Prof. Dr. Savaş YILDIRIM Assoc. Prof. Dr. Banu DİRİ

ISTANBUL BILGI UNIVERSITY YILDIZ TECHNICAL UNIVERSITY Department of Computer Engineering

September 17, 2014

- 2 Experimental Setup
- 3 Methodology
- 4 Results and Evaluation
- 5 Future Work



- Semantic relations are underlying relations betwen two concepts expressed by words or phrases (Moldovan 2004)
- Contradictory views about number of semantic relations
- Uncountability (Murphy 2003)
  - 13 classes (Rosario 2001)
  - 17 classes (Stephens 2001)
  - 5 classes at the top and 30 classes at the bottom (Nastase 2003)
  - 35 classes (Moldovan 2004)
  - 7 classes (SemEval 2007 Task4)
  - etc.

- Five main families of relations:
  - Hyponym/Hypernym (Class inclusion)
  - Meronym/Holonym (Part-Whole)
  - Synonym (Similars)
  - Antonym (Contrast)
  - Case

#### Five main families of relations:

- Hyponym/Hypernym (Class inclusion)
- Meronym/Holonym (Part-Whole)
- Synonym (Similars)
- Antonym (Contrast)
- Case

#### In this study:

- Synonym (Similars)
  - Synonyms are words with identical or similar meanings
  - "car is synonymous with automobile"
- Hyponym/Hypernym (Class inclusion)
  - A relation of inclusion (IS-A / kind of)
  - "A dog is an animal"
  - dog: hyponym / animal: hypernym
- Meronym/Holonym (Part-Whole)
  - A relationship between terms that respect to the significant parts of a whole

伺 ト く ヨ ト く ヨ ト

- "the eye is part of the face"
- eye: part / face: whole

## **Related Works**

- Methods are mostly based on distributional similarity.
- Distributional hypothesis which adopts that "semantically similar words share similar contexts" (Harris 1954).
- However insufficient for synonymy!
  - covers near-synonyms (Eg. top-5 for orange is: yellow, lemon, peach, pink, lime) (Lin et. al. 2003)
  - does not distinguish between synonyms and other relations
- Different strategies:
  - integrating two independent approaches such as distributional similarity and pattern-based approaches

・ 同 ト ・ ヨ ト ・ ヨ ト …

utilizing external features such as dependency relations.

## **Related Works**

- In Turkish, recent studies on synonym relations are based on dictionary definition TDK<sup>1</sup> and Wiktionary.
- Defined rules and phrasal patterns that occur in the dictionary definitions with.
- Objective of this study: to determine synonymy in a Turkish Corpus
- The main contributions of our work:
  - its corpus-driven characteristics
  - relies on both dependency and semantic relations.

<sup>1</sup>TDK:Turkish Language Association

#### Language Resources

- BOUN Web corpus with 500M tokens (Sak et.al. 2008)
  - four sub-corpora:
    - Three of them named NewsCor are from three major Turkish news portals
    - GenCor is a general sampling of web pages in the Turkish Language
  - morphological parser based on two-level morphology
  - an averaged perceptron-based morphological disambiguator
- Preprocessing
  - surface+root+POS-tag+[and all other markers] Turkish: araba English: car
    "arabanın+araba+noun+a3sg+pnon+gen"

伺下 イヨト イヨト

## Methodology

- No particular LSPs to extract synonymy
- Our main assumption:
  - Synonym pairs mostly show similar dependency and semantic characteristics in corpus.
  - They share the same meronym/holonym relations, same particular list of governing verbs, adjective modification profile and so on.
  - Examples:
    - Automobile is a vehicle
    - Car is a vehicle
    - Wheel is part of automobile
    - Wheel is part of car

# Methodology

- **Model:** determine if a given word pair is synonym or not.
- Data: Manually and randomly selected 200 synonym pairs(SYN) and 200 non-synonym pairs(NONSYN) to build a training data set.
  - Non-synonym pairs are especially selected from associated (relevant) pairs such as tree-leaf, student-school, computer-game, etc.
- Features from Corpus: For each synonym pair, 15 different features.
  - Co-occurrence (1)
  - Semantic relations based on LSPs (4)
  - Dependency relations based on syntactic patterns (10)
- Target Class: SYN and NONSYN

・ 同 ト ・ ヨ ト ・ ヨ ト

- Feature 1- Co-occurrence: The first feature is the co-occurrence of word pairs within a broad context
- Features 2/3- Meronym/Holonym: A big matrix in which rows depict <u>Whole</u> candidates, columns depict <u>Part</u> candidates
- Example:

Whole/Part	Part1	Part2	 Partm	
Whole1	val.	val.	 val.	
Whole2	val.	val.	 val.	
	val.	val.	 val.	
Wholen	val.	val.	 val.	

Table: 1. Whole X Part Matrix

 Cells represent the possibility of that corresponding whole and part are in meronymy relation

伺下 イヨト イヨト

▲ 同 ▶ ▲ 国 ▶ ▲ 国 ▶

## Features from Corpus

Table: 2. General Patterns (GP), Dictionary-based Patterns(TDK-P) and Bootstrapped Patterns(BP)

GP	TDK-P	BP
NPx part of NPy	Group-of	NPy-gen NPx-pos
	(whole group all)	
	(set flock union)	
NPx member of NPy	Member-of	NPy-nom NPx-pos
	(class member team)	
	(from the family of Y)	
NPy constituted of NPx	Amount-of	NPy-Gen (N-ADJ)+ NPx-Pos
	(amount measure unit)	
NPy made of NPx	Has/Have (Y has I(H))	NPy of one-of NPx
NPy consist of NPx	Consist-of	NPx whose NPy
NPy has/have NPx	Made-of	NPxs with NPy
NPy with NPx		

 To measure the similarity of meronymy/holonymy profile of two given words, cosine function is applied on two rows/columns indexed by two given words

Whole/Part	Wheel	 Car	 Wiper	Automobile	
Car	Х	 	 Х		
School		 	 		
Person		 	 		
Automobile	Х	 	 Х		
Parking		 Х	 	Х	

Table: 3. Example of Whole X Part Matrix

Dr. Tuğba YILDIZ Assist. Prof. Dr. Savaş YILDIRIM Assoc. Pr An Integrated Approach to Automatic Synonym Detection in Tur

伺 ト く ヨ ト く ヨ ト

- Features 4/5 Hyponym/Hypernym: Another matrix is built for hypernymy/hyponymy by applying LSPs
- Same procedure is carried out as Meronym/Holonym.
  - "NPs gibi CLASS" (CLASS such as NPs)
  - "NPs ve diğer CLASS" (NPs and other CLASS)
  - "CLASS IArdAn NPs" (NPs from CLASS)
  - "NPs ve benzeri CLASS" (NPs and similar CLASS)

- Features 6–15 Dependency Relations: are obtained by syntactic patterns
- The more they are modified by same adjectives, the more likely they are synonym.
- Examples:
- Fast Car (+) / Fast Automobile (+) / Handsome Car (X) / Handsome Automobile (X)
- 36 different patterns were extracted, 8 were eliminated because of the poor results.
- Then we grouped them according to their syntactic structures.

・ 同 ト ・ ヨ ト ・ ヨ ト …

イロト イポト イヨト イヨト

# Features from Corpus

#### Table: 4. Dependency Features

Dependency relation	# of Patterns	Examples
direct object of verb	13	I drive a car
		araba sürüyorum
subject of verb	3	waiting car
		bekleyen araba
direct object/subject of verb	3	-
		-
modified by adjective+(with/without)	2	car with gasoline
		benzinli araba
modified by inf	1	swimming pool
10.11		yuzme havuzu
modified by noun	1	toy car
modified by adjustive	1	oyuncak araba
modified by adjective	1	lummum araba
modified by acronym locations	1	the cars in ABD
modified by actonym locations	1	ABD'deki arabalar
modified by proper noun locations	1	the cars in Istanbul
modified by proper nour locations	1	lite cars in istanbul Istanbul'daki arabalar
modified by locations	2	the car at narking lot
incurred by locations	-	otoparktaki araba
	Dependency relation direct object of verb subject of verb direct object/subject of verb modified by adjective+(with/without) modified by inf modified by noun modified by adjective modified by adjective modified by acronym locations modified by proper noun locations modified by locations	Dependency relation     # of Patterns       direct object of verb     13       subject of verb     3       direct object/subject of verb     3       modified by adjective+(with/without)     2       modified by inf     1       modified by noun     1       modified by adjective     1       modified by adjective     1       modified by acronym locations     1       modified by proper noun locations     2

# Results and Evaluation

- All features explained before are computed for all pairs/instances
- All computed scores of the pairs are kept as a training set
- Taking all features into account and applying machine learning algorithms
- The most suitable algorithm, logistic regression

#### Results and Evaluation

The first aim is to find out which feature is the most informative for detecting synonymy and contributes most to the overall success of the model.

Table: 5. F-Measure of Semantic Relations (SRs) Features

	co-occurrence	hyponym	hypernym	meronym	holonym
F-Measure	62.5	60.5	60	68.7	73.7

- The semantic features are notably better than dependency features.
- Among semantic relations, the most powerful attributes are meronymy and holonymy features.
- The possible reason: the sufficient number of cases matched by lexico-syntactic patterns.

# Results and Evaluation

 Among dependency relations, G1, G4 and G7 have better performance

Table: 6. F-measure of Dependency Relations Features

	G1	G2	G3	G4	G5	G6	G7	G8	G9	G10
F-Measure	64.7	58	60.5	65	61.6	58.8	63	49.4	48.3	62.6

- The poorest groups, G8 and G9, have low production capacity
- The successful features are linearly dependent on target class.
- On the aggregated data where all useful features are considered, the performance of logistic regression is Fmeasure of 80.3%.
- The achieved score is better than the individual performance of each feature.

化原因 化原因

#### Done

- In our paper: "Dictionary definitions, WordNet, and other useful resources could be used and evaluated in future work."
- In our review: "WordNet seems a good resource helping improve it even more."

直 ト イヨト イヨト

#### Done



## Conclusion

- Semantic relation is one of the problem in applications in NLP.
- Automatic acquisition of semantic relation has become more popular in recent years.
- It is obviously impossible to discover synonym pairs by applying LSP.
- Thanks to semantic and syntactic relations.
- The most powerful semantic relations: Meronym and Holonym.
- The most powerful syntactic relations: G4 modified by adjective (with/without).
- Aggregation of useful features gives promising result.

伺下 イヨト イヨト

Thank You For Listening...

Dr. Tuğba YILDIZ Assist. Prof. Dr. Savaş YILDIRIM Assoc. Pr An Integrated Approach to Automatic Synonym Detection in Tur

▲ 同 ▶ ▲ 国 ▶ ▲ 国 ▶