

A Game – Theoretic Approach to Word Sense Disambiguation

Rocco Tripodi
Marcello Pelillo



context

+ WSD

Introduction

- **New model** for Word Sense Disambiguation (WSD) formulated in terms of evolutionary game theory
- **Word Sense Disambiguation:** the task of identifying the intended meaning of a word based on the context in which it appears
- “There is a *financial institution* near the river **bank**”
 - Financial sense
 - Naturalistic sense

Background

- Approaches based on
 - 1. Learning models
 - Supervised
 - Unsupervised
 - Semi-supervised
 - 2. Techniques
 - Heuristic
 - 3. Algorithms
 - Graph-based
 - Knowledge-based
- ✓ Remain on the surface of the word, compromising the coherence of the data to be analyzed

Game Theory

- Mathematical approach to study the interaction between two or more individuals
 - **Outcome** → benefits and costs depend on the strategies of each other
 - **Players** $I = \{1, \dots, n\} \rightarrow N$: total number of players
 - **Pure player strategy set** $S_i = \{s_1, s_2, \dots, s_M\}$, M : total number of strategies
 - **Mixed player strategy set** $x_i = \{x_{i1}, x_{i2}, \dots, x_{iM}\}$, x_h : probability of player i to choose its $h - th$ pure strategy
 - **Utility function** $u_i : S_1 \times S_2 \times \dots \times S_N \rightarrow \mathbb{R}$
 - strategies → payoffs
 - combination of strategies played by all players
- **Classical game theory VS Evolutionary game theory**
 - **Static VS Dynamic** strategies
 - **Evolutionary game theory** → dynamics of strategy change through repeated games
- **Nash equilibrium** → set of strategy profiles in which each strategy is a best response to the strategy of the co-players and no player has the incentive to deviate from their decision (changing strategy $\not\rightarrow$ payoff increase)
- **Payoff** → value associated with a possible outcome of a game (set of strategies)

Proposed System

Evolutionary game theory framework

- **First attempt** in the specific NLP task of WSD

WSD

- Sense - **labelling** task → sense assignment to words
- **Constraint** satisfaction problem

Game-theoretical approach

- **Players** → words (to be disambiguated)
- **Strategies** → senses (evolving population)
- **Payoff** matrices → sense similarity
- **Interactions** → weighted graph
- **Nash equilibrium** → consistent word-sense assignment
- **Selection process**
 - Iterative process of **fitness increase** (candidates / senses with certain features)
 - **Best candidates** (senses with higher fitness) in the population

Proposed System

Consistent final labelling of the data

The solution of the problem is always found

- System convergence to the nearest Nash equilibrium
(Nash Theorem 1951)

Most appropriate sense association

- Target word

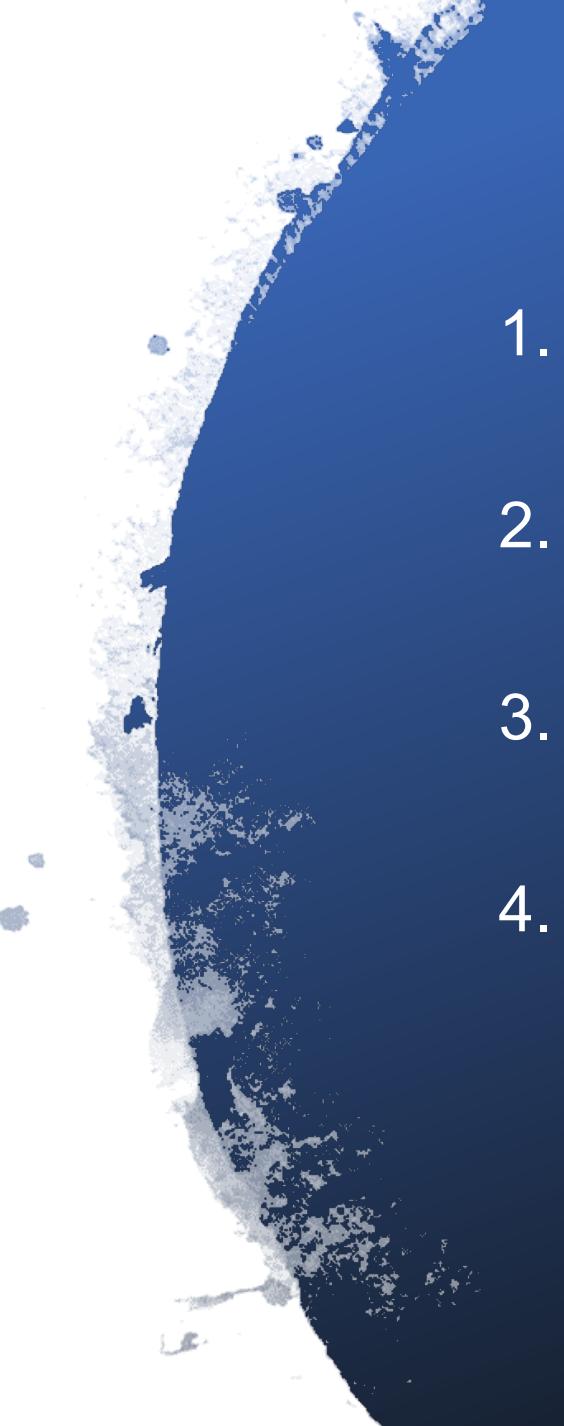
WSD

- Continuous optimization problem
- Exploitation of contextual information in a dynamic way
(evolutionary game theoretic framework)

Versatile approach

- Adaptive to different scenarios and tasks
- Unsupervised / Semi-supervised

Data Modelling

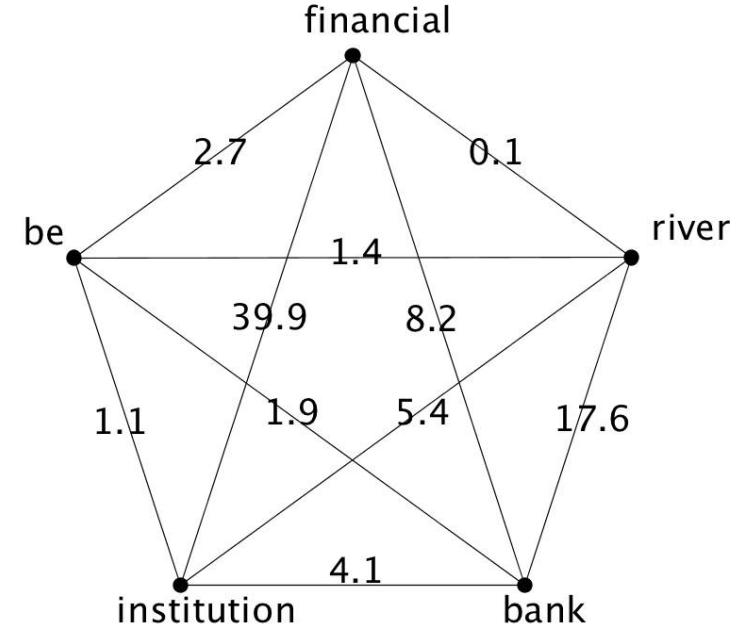
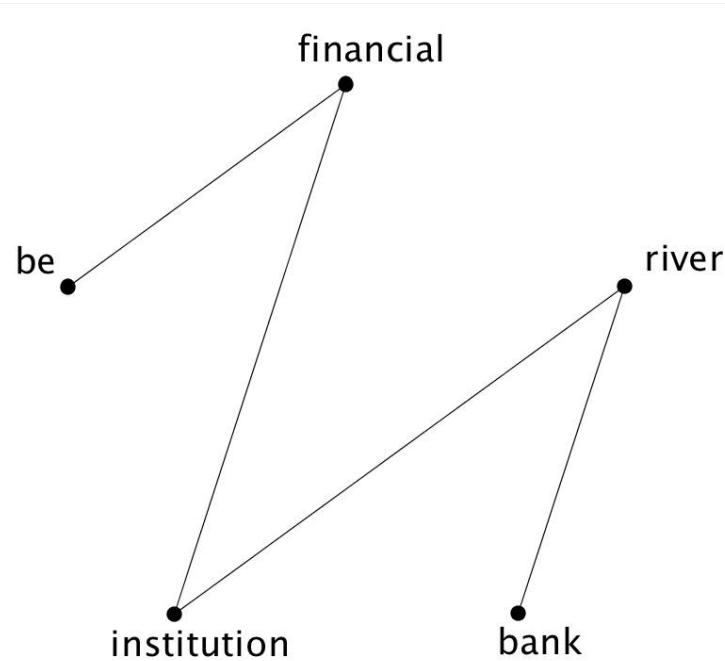
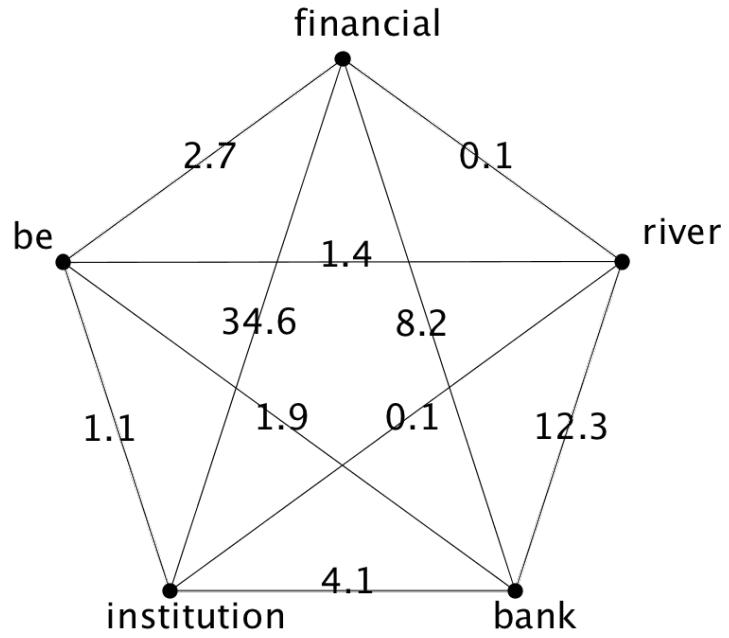
- 
1. Data Geometry
 2. Strategy Space
 3. Payoff
 4. System Dynamics & Sense Classification

Data Geometry

- **List I of N words** from the text
 - $I = \{i_1, i_2, \dots, i_N\}$
- Word **similarity matrix $W_{N \times N}$**
 - Pairwise similarities among words → players' interactions
 - $w_{kj} = \text{sim}(i_k, i_j), \forall k, j \in I : k \neq j$
 - Weighted adjacency matrix of the graph
- **Similarity measure:** strength of co-occurrence between words i, j
 - Semantically correlated words
 - 8 association measures: Dice coefficient (dice), modified Dice coefficient (mDice), pointwise mutual information (pmi), t-score measure (t-score), z-score measure (z-score), odds ratio (odds-r), chi-squared test (chi-s), chi-squared correct (chi-s-c).
- **Proximity** relations with n - neighbours (similarity **augmentation**)
 - Sentence structure
 - Size of n : fixed expressions / semantic concepts



Data Geometry



Strategy Space



- **Sense inventories** of each word $i \in I$
 - $M_i = \{ m_{i1}, m_{i2}, \dots, m_{ik} \}$, where $k = |M_i|$ the number of senses associated with word i
 - WordNet and BabelNet as knowledge base
- **Player strategy space**
 - $\Delta_p = \{ x_p \in \mathbb{R}^n : \sum_{i=1}^n x_{pi} = 1, \text{where } x_{pi} \geq 0 \text{ for } i = 0, 1, \dots, n \}$
 - n : number of pure strategies of player i and each component
 - x_{pi} : probability that player p chooses its $i - th$ pure strategy
 - Graphical representation: regular polygon of radius 1, where the distance from the centre to any vertex (mixed strategy) represents the probability associated with a particular word sense (pure strategy) →
- **Unique concepts** in sense inventories
 - Game space $C = \{ c_1, c_2, \dots, c_u \}$, where u the total number of unique concepts
 - Assignment of **probability distribution** over the senses in C , for each word $i \in I$

Strategy Space



- **System** strategy space S
 - s_{ih} : **probability** of player to choose its $h - th$ pure strategy to play
 - **Uniform** distribution
 - $s_{ih} = \begin{cases} |M_i|^{-1} & , \text{ if sense } h \text{ is in } M_i \\ 0, & \text{otherwise} \end{cases}$
 - unsupervised learning: no prior knowledge
- **Geometric** distribution
 - semi-supervised learning: exploits information from prior knowledge
 - $s_{ih} = \begin{cases} p(1 - p)^{r_h}, & \text{if sense } h \text{ is in } M_i \\ 0, & \text{otherwise} \end{cases}$
 - ✓ higher probability on senses that have a high frequency
 - $s_{ih} = \begin{cases} p(1 - p)^{r_o}, & \text{if sense } h \text{ is in cluster } o \\ 0, & \text{otherwise} \end{cases}$
 - ✓ equal probability to the senses belonging to a determined cluster

Payoff

- Sense similarity matrix Z
 - Pairwise similarity among senses in $C \rightarrow$ partial payoff matrices of each game
 - $z_{ij} = ssim(s_i, s_j), \forall i, j \in C, i \neq j$
- Partial payoff matrix:
 - Games played between two words i, j
 - Dimension $|M_i| \times |M_j|, |M_k|$: cardinality of senses set of word k
- Semantic similarity
 - Relations of likeness (“is-a”)
 - wup: path length among two senses s_i, s_j
 - $ssim_{wup}(s_i, s_j) = 2 * \frac{depth(msa)}{depth(s_i) + depth(s_j)}$
 - msa : “most specific ancestor”
 - jcn measure: corpus statistics and structural properties
 - $ssim_{jcn}(s_i, s_j) = IC(s_1) + IC(s_2) - 2 IC(msa) * \frac{depth(msa)}{depth(s_i) + depth(s_j)}$
 - $IC(s_i) = -\log(p(s_i))$

Payoff

- Semantic relatedness
 - Similarity among the definitions of two concepts
 - Wider range of relations (“is-a-part-of”, “is-the-opposite-of”)
 - Definitions derived from glosses of the synsets in WordNet
 - Co-occurrence vector $v_i = (w_{1i}, w_{2i}, \dots, w_{ni})$, i : concept, w : word gloss occurrences, n : total words
 - Cosine similarity
 - $sim(v_i, v_j) = \cos \theta = \frac{v_i \cdot v_j}{\|v_i\| \|v_j\|}$, i, j : concepts
 - $\|v_i\| = \sqrt{\sum_{j=1}^n w_{ji}}$
 - cosine of the angle θ between the two co-occurrence vectors v_i, v_j
 - 4 variations (construction of vectors)
 - Difference in co-occurrence calculation, corpus use and relation source
 - tf-idf, tf-idf_ext, vec and vec_ext

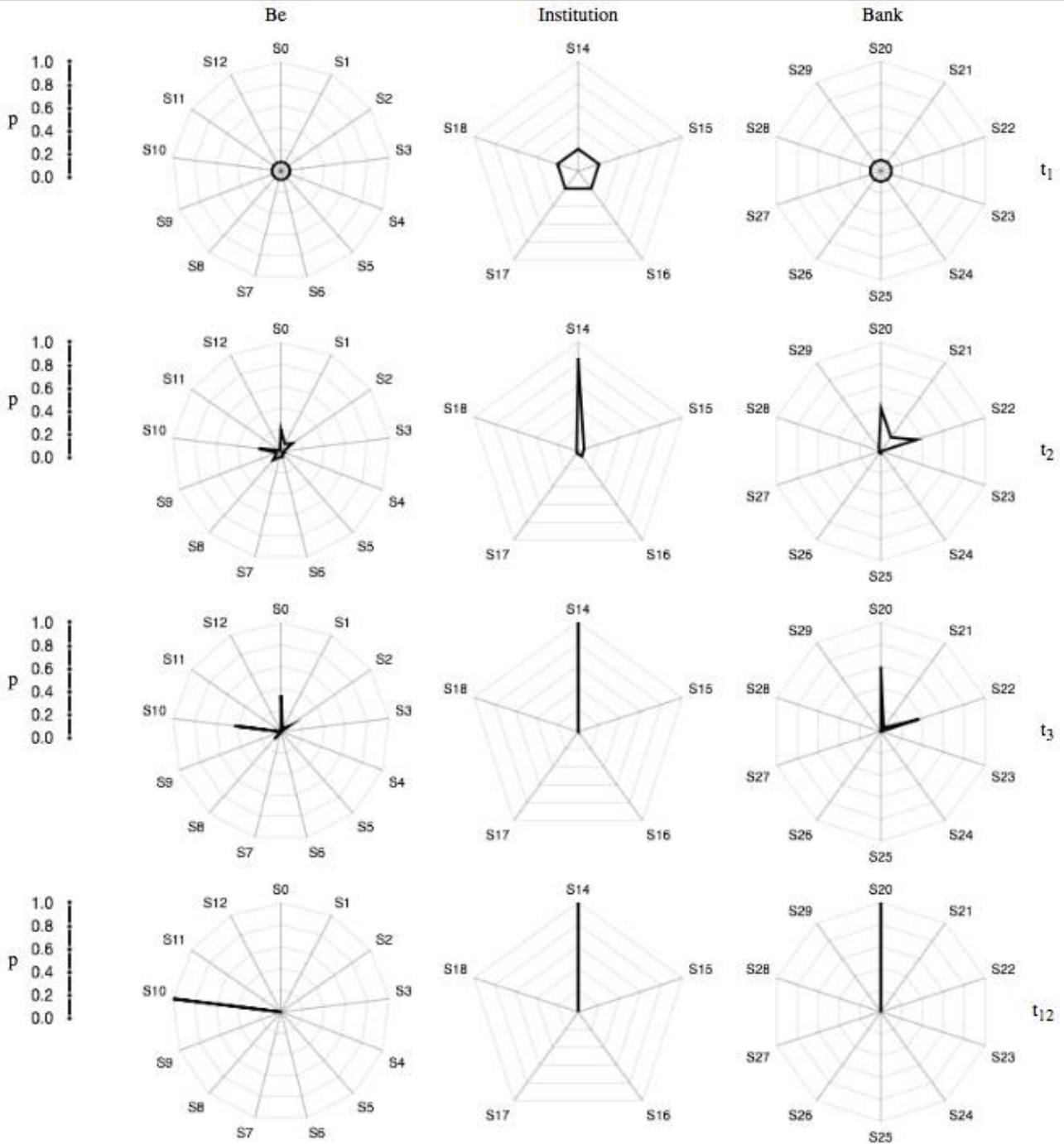
System Dynamics & Sense Classification

- Computation of the Nash equilibrium
 - Expected **fitness of sense** h (word i)
 - $u_i(e^h) = \sum_{j \in N_i} (w_{ij}, Z_{ij}, x_j)_h$, where N_i the neighbours (context) of word i
 - Pure strategy h : e^h
 - Similarity with word j : w_{ij}
 - Similarities among senses of i, j : Z_{ij}
 - Sense preference of word j : x_j
 - Average **fitness of population**
 - $u_i(x) = \sum_{j \in N_i} x_i^T (w_{ij}, Z_{ij}, x_j)$
 - **Replicator Dynamics Equation**
 - $x_i^h(t+1) = x_i^h(t) \frac{u_i(e^h)}{u_i(x)}$, $\forall h \in S_i$
 - Models the change in sense frequency
 - Depends on performance → measured relatively to the average fitness of population

System Dynamics & Sense Classification

- Classification
 - Unique strategy $s \in S \rightarrow$ each word $i \in I$
 - The strategy with the highest probability is chosen
 - $\varphi_i = \text{argmax}(x_{ih})$, where $h = 1, 2, \dots, |C|$, $|C|$ total number of senses
 - Each word $i \in I$ is mapped to exactly one sense $c \in C$
 - Word fails to be disambiguated
 - Unable to update its strategy space
 - Strategy space initialized with a uniform distribution
 - Zero entry payoff (no similar senses with co-players)
 - Not connected with other nodes in the graph

System Dynamics & Sense Classification



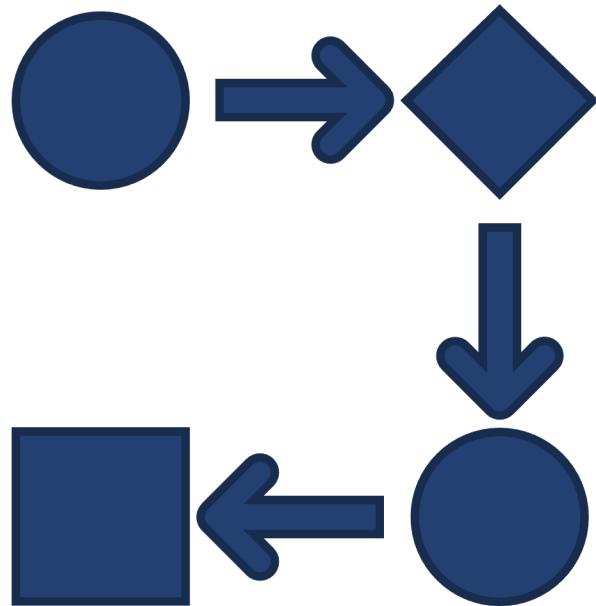
Parameter Tuning

- Two data sets to evaluate our algorithm in different scenarios
- From each data set
 - 50 different data sets to simulate a situation in which the system must work on texts of different sizes and on different domains
- The parameters that will be tuned
 - Association, similarity and relatedness measures to weight the similarity among word and senses
 - The n -gram graph to increase the weights of near words (proximity)
 - The p of the geometric distribution used by their semi-supervised system

Parameter Tuning - Results

1. Association, Semantic & Relatedness Measures
 - The relatedness measures perform significantly better than the semantic similarity measures
 - Particularly suitable measures for the algorithm
 - mdice
 - Tf-idf
2. n - gram Graph
 - Highest results on both data sets using $n = 5$ (n nearest neighbours)
3. Geometric Distribution:
 - Best results obtained with $p = 0.4$

Evaluation set-up



- Results are provided as balanced F-score (F_1) measure
 - F_1 : determines the weighted harmonic mean of precision and recall
 - $$F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} \times 100 \quad (\%)$$
- Evaluation based on
 - Experiments with WordNet as Knowledge base
 - The best performance on nouns on all the data sets
 - Low results on verbs on all data sets
 - The use of prior knowledge is beneficial in general domain data sets
 - Experiments with BabelNet as Knowledge base
 - We used BabelNet to collect the sense inventories of each word to be disambiguated and NASARI to obtain the semantic representation of each sense
 - This data set contains highly ambiguous mentions

Comparison to state-of-the-art algorithms

WordNet:

- Our unsupervised system performs better than any other unsupervised algorithm in all datasets
- The performance of our system is more stable on the four datasets
- The comparison with semi-supervised systems shows that our system always performs better than the most frequent sense heuristic when we use information from sense-labeled corpora

BabelNet:

- The performance of our system is close to the results obtained with Babelfy on S13 and substantially higher on KORE50
- It is also difficult to exploit distributional information on this data set because the sentences are short and, in many cases, cryptic.

Comparison to state-of-the-art algorithms

	S7CG	S7CG (N)	S7	S3	S2
unsup.	<i>Nav10</i>	—	—	43.1	52.9
	<i>PPR</i> _{w2w}	80.1	83.6	41.7	57.9
	<i>WSD</i> _{games}	80.4*	85.5	43.3	59.1
semi sup.	<i>IRST-DDD-00</i>	—	—	—	58.3
	<i>MFS</i>	76.3	77.4	54.7	62.8
	<i>MRF-LP</i>	—	—	50.6*	58.6
	<i>Nav05</i>	83.2	84.1	—	60.4
	<i>PPR</i> _{w2w}	81.4	82.1	48.6	63.0
	<i>WSD</i> _{games}	82.8	85.4	56.5	64.7*
	<i>Best</i>	82.5	82.3*	59.1	65.2
sup.	<i>Zhong10</i>	82.6	—	58.3	67.6

	S13	KORE50
<i>WSD</i> _{games}	70.8	75.7
<i>Babelfy</i>	69.2	71.5
<i>SUDOKU</i>	66.3	—
<i>MFS</i>	66.5*	—
<i>PPR</i> _{w2w}	60.8	—
<i>KORE</i>	—	63.9*
<i>GETALP</i>	58.3	—

Conclusion



A new method for WSD

Evolutionary Game Theory

Similarity measures that perform better

Continuation of knowledge-based, graph-based approaches



WSD as a constraint-satisfaction problem

Consistency in the assignment of senses to related words

Development of contextual coherence on the assignment of senses
(characteristic missing in many state-of-the-art systems)



Replicator dynamics equation

Best labelling assignment



Versatile

Unsupervised

Semi-supervised



Competitive compared to state-of-the-art systems

Considers the influence of each word on the others

Imposes sense compatibility among each sense before assigning a meaning
The meaning of a word depends only on words that share a proximity relation
and on those that enjoy a high distributional similarity

Thank you

THEORY



M902 – “Βασικές Μαθηματικές Έννοιες στη Γλωσσική Τεχνολογία”

Kylafi Christina-Theano
Piriasi Juliana

M.Sc. in Language Technology

Department of Informatics and Telecommunications,
National and Kapodistrian University of Athens
&

Institute for Language and Speech Processing,
“Athena” Research Center