

Hypergraphs and Information Fusion for Term Representation Enrichment. Applications to Named Entity Recognition and Word Sense Disambiguation

Ph.D. Thesis Defense

Pavel Soriano-Morales

Supervised by Sabine Loudcher and Julien Ah-Pine

February 7th, 2018



UNIVERSITÉ
LUMIÈRE
LYON 2



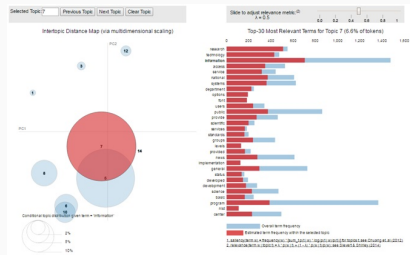
UNIVERSITÉ
DE LYON

INSTITUT
DES SCIENCES
DE L'HOMME



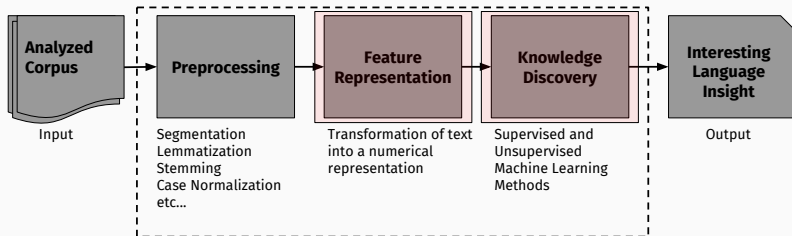
Introduction

Why is it useful to us to understand text?

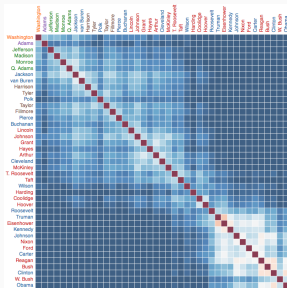


How do we extract meaning from text?

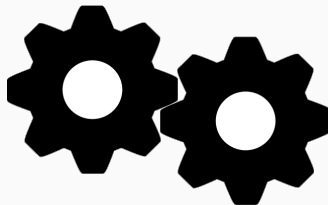
We use **Natural Language Processing** (NLP), a field of computer science interested in making computers extract useful information from text



How do we represent text for the machine to understand?



What techniques do we use to discover meaning from text?



Representing Text

- **Three common ways to represent text**

Representing Text

- **Three common ways to represent text**
 - Lexical

Representing Text

- **Three common ways to represent text**
 - Lexical
 - Syntactic

Representing Text

- **Three common ways to represent text**
 - Lexical
 - Syntactic
 - Constituency Tree

Representing Text

- **Three common ways to represent text**
 - Lexical
 - Syntactic
 - Constituency Tree
 - **Dependency Tree**

Representing Text

- **Three common ways to represent text**
 - Lexical
 - Syntactic
 - Constituency Tree
 - Dependency Tree
- **Working Example**

Representing Text

- **Three common ways to represent text**

- Lexical
- Syntactic
 - Constituency Tree
 - Dependency Tree

- **Working Example**

The report contains copies of the minutes of these meetings

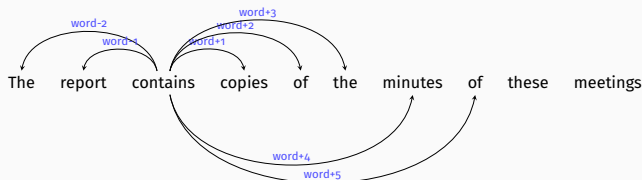
Representing Text

- **Three common ways to represent text**

- Lexical
- Syntactic
 - Constituency Tree
 - Dependency Tree

- **Working Example**

The report contains copies of the minutes of these meetings



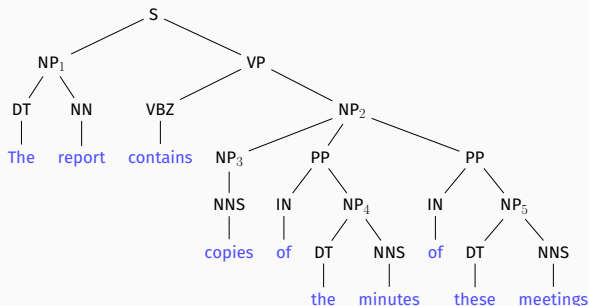
Representing Text

- **Three common ways to represent text**

- Lexical
- Syntactic
 - Constituency Tree
 - Dependency Tree

- **Working Example**

The report contains copies of the minutes of these meetings



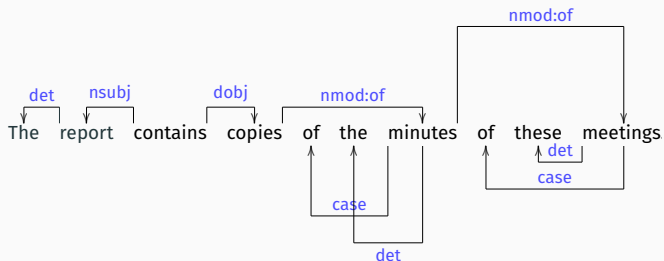
Representing Text

- **Three common ways to represent text**

- Lexical
- Syntactic
 - Constituency Tree
 - Dependency Tree

- **Working Example**

The report contains copies of the minutes of these meetings



- **Text Representation Models**

- Words and features can be represented by means of graph-based models matrices
- Or directly with (sparse) matrices

- **Leveraging the Network Structure**

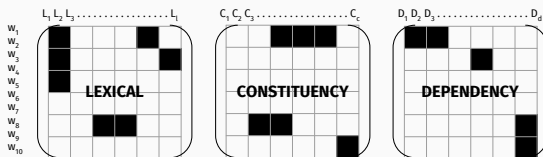
- We can find communities of similar words according to their meaning

- **Text Representation Models**

- Words and features can be represented by means of graph-based models matrices
- Or directly with (sparse) matrices

- **Leveraging the Network Structure**

- We can find communities of similar words according to their meaning



Main Challenges and Contributions

1. What type of model can we employ to represent a corpus using heterogeneous features?

Main Challenges and Contributions

1. What type of model can we employ to represent a corpus using heterogeneous features?
 - *Hypergraph linguistic model to hold different types of linguistic information*

Main Challenges and Contributions

1. What type of model can we employ to represent a corpus using heterogeneous features?
 - *Hypergraph linguistic model to hold different types of linguistic information*
2. How can we combine these features while dealing with feature sparsity?

Main Challenges and Contributions

1. What type of model can we employ to represent a corpus using heterogeneous features?
 - *Hypergraph linguistic model to hold different types of linguistic information*
2. How can we combine these features while dealing with feature sparsity?
 - *Multimedia fusion techniques to combine and densify representation spaces*

Main Challenges and Contributions

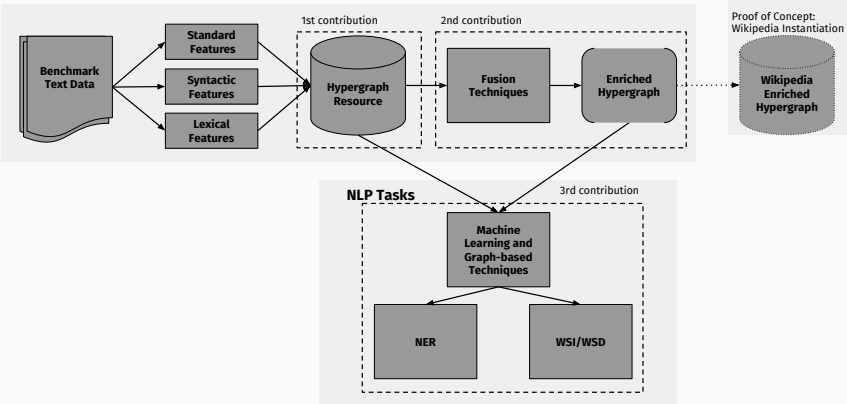
1. What type of model can we employ to represent a corpus using heterogeneous features?
 - *Hypergraph linguistic model to hold different types of linguistic information*
2. How can we combine these features while dealing with feature sparsity?
 - *Multimedia fusion techniques to combine and densify representation spaces*
3. How can we find and employ communities existing within the language networks?

Main Challenges and Contributions

1. What type of model can we employ to represent a corpus using heterogeneous features?
 - *Hypergraph linguistic model to hold different types of linguistic information*
2. How can we combine these features while dealing with feature sparsity?
 - *Multimedia fusion techniques to combine and densify representation spaces*
3. How can we find and employ communities existing within the language networks?
 - *An alternative network-based algorithm to discover semantically related words within a text*

Introduction

Work Overview



Contributions in Detail

Hypergraph Linguistic Model

Introduction

Based on the distributional hypothesis, a word is defined by its surroundings, we can extract useful information from a text.

- **How do we represent textual data?**

Based on the distributional hypothesis, a word is defined by its surroundings, we can extract useful information from a text.

- **How do we represent textual data?**
 - Network Models [MTFo4]

Based on the distributional hypothesis, a word is defined by its surroundings, we can extract useful information from a text.

- **How do we represent textual data?**
 - Network Models [MTFo4]
 - Vector Space Models [MS+99]

Based on the distributional hypothesis, a word is defined by its surroundings, we can extract useful information from a text.

- **How do we represent textual data?**

- Network Models [MTFo4]
- Vector Space Models [MS+99]

- **We choose network models**

Based on the distributional hypothesis, a word is defined by its surroundings, we can extract useful information from a text.

- **How do we represent textual data?**
 - Network Models [MTFo4]
 - Vector Space Models [MS+99]
- **We choose network models**
 - Used in a large quantity of NLP tasks [MR11]

Based on the distributional hypothesis, a word is defined by its surroundings, we can extract useful information from a text.

- **How do we represent textual data?**

- Network Models [MTFo4]
- Vector Space Models [MS+99]

- **We choose network models**

- Used in a large quantity of NLP tasks [MR11]
- Graphs structures can give us a clearer view into the relations of words within a text [CM09]

Based on the distributional hypothesis, a word is defined by its surroundings, we can extract useful information from a text.

- **How do we represent textual data?**

- Network Models [MTFo4]
- Vector Space Models [MS+99]

- **We choose network models**

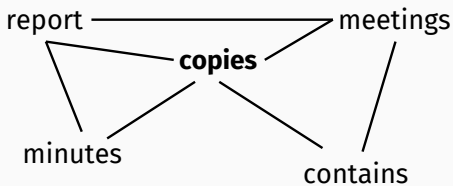
- Used in a large quantity of NLP tasks [MR11]
- Graphs structures can give us a clearer view into the relations of words within a text [CM09]
- Ultimately graphs are transformed to a vectorial representation through the adjacency/incidence matrices

The report contains copies of the minutes of these meetings

The report contains copies of the minutes of these meetings

Lexical Networks

Sentence Level

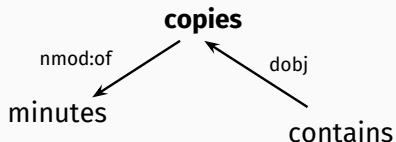


[KMo8]

The report contains copies of the minutes of these meetings

Syntactic Networks

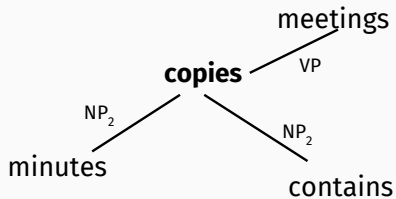
Dependency Tree



The report contains copies of the minutes of these meetings

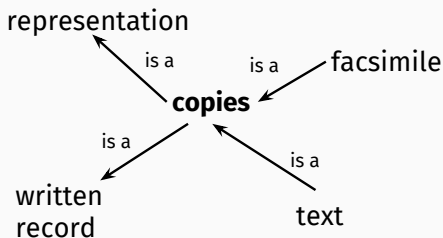
Syntactic Networks

Constituency Tree



The report contains copies of the minutes of these meetings

Semantic Networks



Limitations and Proposition

- **Limitations of existing representations**

Limitations and Proposition

- **Limitations of existing representations**
 - Language networks generally employ a single type of textual information

Limitations and Proposition

- **Limitations of existing representations**
 - Language networks generally employ a single type of textual information
 - The edges of the network may relate maximum two words at each time

Limitations and Proposition

- **Limitations of existing representations**
 - Language networks generally employ a single type of textual information
 - The edges of the network may relate maximum two words at each time
- **Proposition**

Limitations and Proposition

- **Limitations of existing representations**
 - Language networks generally employ a single type of textual information
 - The edges of the network may relate maximum two words at each time
- **Proposition**
 - Represent together linguistic co-occurrences through a hypergraph model

Limitations and Proposition

- **Limitations of existing representations**

- Language networks generally employ a single type of textual information
- The edges of the network may relate maximum two words at each time

- **Proposition**

- Represent together linguistic co-occurrences through a hypergraph model
 - Link together three different types of networks, using lexical and syntactic data

Limitations and Proposition

- **Limitations of existing representations**

- Language networks generally employ a single type of textual information
- The edges of the network may relate maximum two words at each time

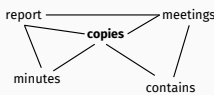
- **Proposition**

- Represent together linguistic co-occurrences through a hypergraph model
 - Link together three different types of networks, using lexical and syntactic data
 - Get a semantic overview at three different levels: short range (with dependency functions), medium range (phrase constituency membership), and long range (lexical co-occurrence)

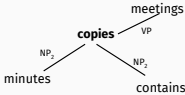
Hypergraph Linguistic Model

Proposed Model

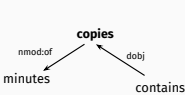
Lexical Networks
Sentence Level



Syntactic Networks
Constituency Tree



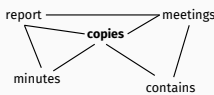
Syntactic Networks
Dependency Tree



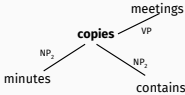
Hypergraph Linguistic Model

Proposed Model

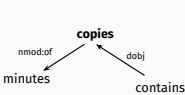
Lexical Networks
Sentence Level



Syntactic Networks
Constituency Tree

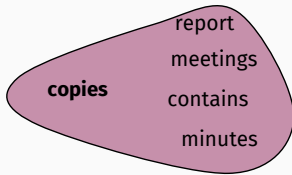


Syntactic Networks
Dependency Tree



Hypergraph Model

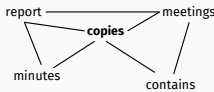
■ Lexical



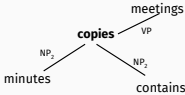
Hypergraph Linguistic Model

Proposed Model

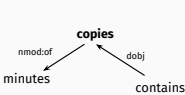
Lexical Networks
Sentence Level



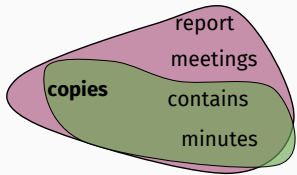
Syntactic Networks
Constituency Tree



Syntactic Networks
Dependency Tree



Hypergraph Model

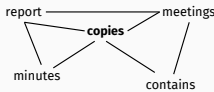


- Lexical
- Constituency (NP₂)

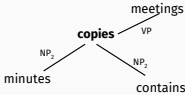
Hypergraph Linguistic Model

Proposed Model

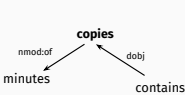
Lexical Networks
Sentence Level



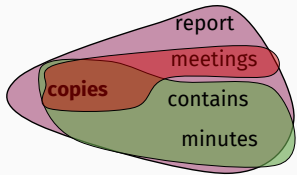
Syntactic Networks
Constituency Tree



Syntactic Networks
Dependency Tree



Hypergraph Model

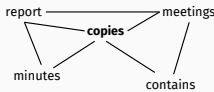


- Lexical
- Constituency (NP₂)
- Constituency (VP)

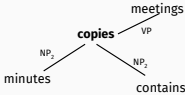
Hypergraph Linguistic Model

Proposed Model

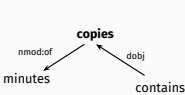
Lexical Networks
Sentence Level



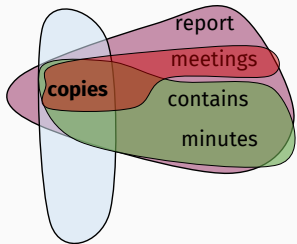
Syntactic Networks
Constituency Tree



Syntactic Networks
Dependency Tree



Hypergraph Model

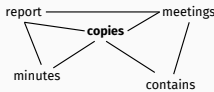


- Lexical
- Constituency (NP₂)
- Constituency (VP)
- Dependency (dobj:contains)

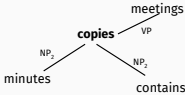
Hypergraph Linguistic Model

Proposed Model

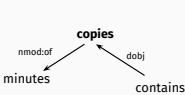
Lexical Networks
Sentence Level



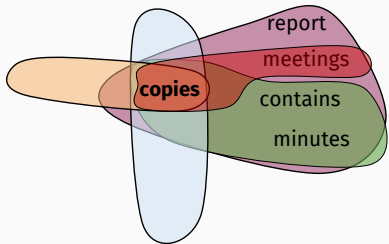
Syntactic Networks
Constituency Tree



Syntactic Networks
Dependency Tree



Hypergraph Model



- Lexical
- Constituency (NP₂)
- Constituency (VP)
- Dependency (dobj:contains)
- Dependency (nmod:of)

Contributions in Detail

**Combining Features and Dealing with
Sparsity**

- **Definition**

- **Definition**

- Set of techniques used in multimedia analysis tasks to integrate multiple media [Atr+10; ABL10]

- **Definition**

- Set of techniques used in multimedia analysis tasks to integrate multiple media [Atr+10; ABL10]
- The goal is to obtain rich insights about the data being treated

- **Definition**

- Set of techniques used in multimedia analysis tasks to integrate multiple media [Atr+10; ABL10]
- The goal is to obtain rich insights about the data being treated
- We adapt these techniques to our use case: textual information

- **Definition**

- Set of techniques used in multimedia analysis tasks to integrate multiple media [Atr+10; ABL10]
- The goal is to obtain rich insights about the data being treated
- We adapt these techniques to our use case: textual information

- **Main fusion operators:**

- **Definition**

- Set of techniques used in multimedia analysis tasks to integrate multiple media [Atr+10; ABL10]
- The goal is to obtain rich insights about the data being treated
- We adapt these techniques to our use case: textual information

- **Main fusion operators:**

- Early Fusion $E_{\alpha}(\cdot)$,

- **Definition**

- Set of techniques used in multimedia analysis tasks to integrate multiple media [Atr+10; ABL10]
- The goal is to obtain rich insights about the data being treated
- We adapt these techniques to our use case: textual information

- **Main fusion operators:**

- Early Fusion $E_{\alpha}(\cdot)$,
- Late Fusion $L_{\beta}(\cdot)$,

- **Definition**

- Set of techniques used in multimedia analysis tasks to integrate multiple media [Atr+10; ABL10]
- The goal is to obtain rich insights about the data being treated
- We adapt these techniques to our use case: textual information

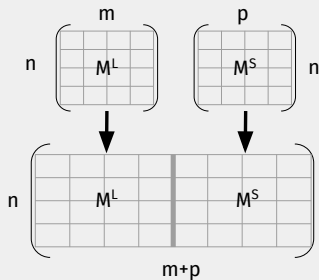
- **Main fusion operators:**

- Early Fusion $E_\alpha(\cdot)$,
- Late Fusion $L_\beta(\cdot)$,
- Cross Fusion $X_\gamma(\cdot), X_F(\cdot)$

Early and Late Fusion

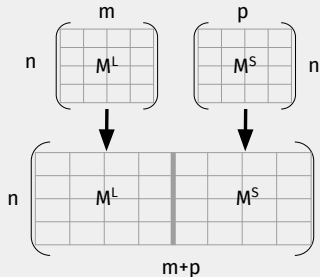
EARLY FUSION

Matrices M^L and M^S have the same number of rows



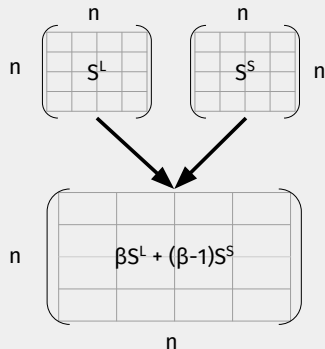
EARLY FUSION

Matrices M^L and M^S have the same number of rows

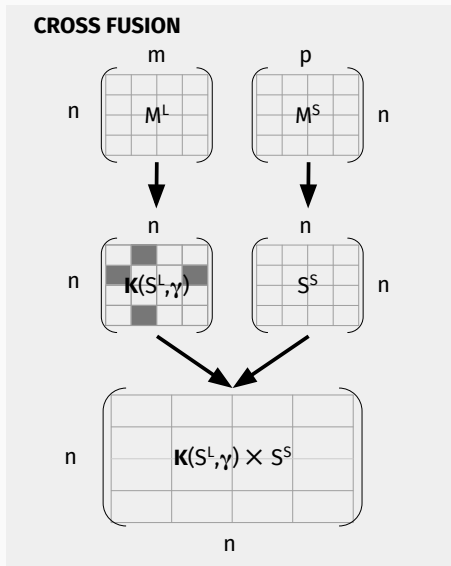


LATE FUSION: SIMILARITY FUSION

Matrices S^L and S^S have the same size



Cross Fusion



- We distinguish three levels of fusion operators
 - **First Degree**
 - $E(M^L, M^S), L(S^S, M^L)$
 - Cross Feature Fusion: $X_F(S^S, M^L)$
 - Cross Similarity Fusion: $X_S(S^S, S^L)$
 - **Second Degree**
 - Cross Feature Early Fusion: $X_F(S^T, E(M^L, M^S))$
 - Late Cross Feature Fusion: $L(M^T, X_F(S^T, M^T))$
 - **Higher Degree**
 - Triple Early Double Late Cross Feature Fusion:
 $E(M_L, E(E(M_T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$

Hybrid Fusion

$$E(M_L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$$

Hybrid Fusion

The diagram illustrates a nested function expression for hybrid fusion, represented as $E(M^L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$. The expression is enclosed in a blue rectangular box. Inside this box, the first argument is $E(M^L, \dots)$. The second argument is a nested expression $E(E(M^T, \dots), \dots)$. The innermost expression is $L(M^T, X_F(S^T, M^T))$, which is highlighted with a red background. This is followed by a comma and another nested expression $L(M^L, X_F(S^S, M^L))$, which is highlighted with a yellow background. The entire nested structure is enclosed in a pink rectangular box.

$$E(M^L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$$

$$L(M^L, X_F(S^S, M^L))$$

$$\begin{aligned} \begin{matrix} w_1 & w_2 & w_3 \\ \begin{pmatrix} S^S \end{pmatrix} \end{matrix} \times \begin{matrix} f_{L1} & f_{L2} & f_{L3} \\ \begin{pmatrix} M^L \end{pmatrix} \end{matrix} &= \begin{matrix} f_{L1} & f_{L2} & f_{L3} \\ \begin{pmatrix} X_F(S^S, M^L) \end{pmatrix} \end{matrix} \\ \begin{matrix} f_{L1} & f_{L2} & f_{L3} \\ \begin{pmatrix} M^L \end{pmatrix} \end{matrix} + \begin{matrix} f_{L1} & f_{L2} & f_{L3} \\ \begin{pmatrix} X_F(S^S, M^L) \end{pmatrix} \end{matrix} &= \begin{matrix} f_{L1} & f_{L2} & f_{L3} \\ \begin{pmatrix} L(M^L, X_F(S^S, M^L)) \end{pmatrix} \end{matrix} \end{aligned}$$

$$E(M_L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$$

$$L(M^T, X_F(S^T, M^T))$$

$$\begin{aligned} \begin{matrix} w_1 & w_2 & w_3 \\ \begin{pmatrix} S^T \end{pmatrix} \end{matrix} \times \begin{matrix} f_{T1} & f_{T2} & f_{T3} \\ \begin{pmatrix} M^T \end{pmatrix} \end{matrix} &= \begin{matrix} f_{T1} & f_{T2} & f_{T3} \\ \begin{pmatrix} X_F(S^T, M^T) \end{pmatrix} \end{matrix} \\ \begin{matrix} f_{T1} & f_{T2} & f_{T3} \\ \begin{pmatrix} M^T \end{pmatrix} \end{matrix} + \begin{matrix} f_{T1} & f_{T2} & f_{T3} \\ \begin{pmatrix} X_F(S^T, M^T) \end{pmatrix} \end{matrix} &= \begin{matrix} f_{T1} & f_{T2} & f_{T3} \\ \begin{pmatrix} L(M^T, X_F(S^T, M^T)) \end{pmatrix} \end{matrix} \end{aligned}$$

$$E(M_L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$$

$$E(M^T, L(M^T, X_F(S^T, M^T)))$$

$$\begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{T1} & f_{T2} & f_{T3} \\ M^T \end{pmatrix} \parallel \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{T1} & f_{T2} & f_{T3} \\ L(M^T, X_F(S^T, M^T)) \end{pmatrix} = \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{T1} & f_{T2} & f_{T3} \\ L(M^T, X_F(S^T, M^T)) \end{pmatrix}$$

$$E(M_L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$$

$$E(M_L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$$

$$\begin{matrix} f_{L1} & f_{L2} & f_{L3} \\ \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \end{matrix} \left(\begin{matrix} M^T \end{matrix} \right) \parallel \begin{matrix} f_{L1} & f_{L2} & f_{L3} \\ \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \end{matrix} \left(E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L)))) \right) =$$

$$\begin{matrix} f_{L1} & f_{L2} & f_{L3} \\ \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \end{matrix} \left(E(M_L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L)))) \right)$$

Contributions in Detail

Finding Communities in the Network

- **Language networks tend to be scale-free**

- **Language networks tend to be scale-free**
 - There are certain nodes (hubs) that are very well connected forming communities within the network

- **Language networks tend to be scale-free**
 - There are certain nodes (hubs) that are very well connected forming communities within the network
- **Seminal approaches**

- **Language networks tend to be scale-free**
 - There are certain nodes (hubs) that are very well connected forming communities within the network
- **Seminal approaches**
 - Hyperlex [VÓ4]

- **Language networks tend to be scale-free**
 - There are certain nodes (hubs) that are very well connected forming communities within the network
- **Seminal approaches**
 - Hyperlex [VÓ4]
 - University of York (UoY) [KM07]

- **Language networks tend to be scale-free**
 - There are certain nodes (hubs) that are very well connected forming communities within the network
- **Seminal approaches**
 - Hyperlex [VÓ4]
 - University of York (UoY) [KM07]
- **Limitations of existing approaches**

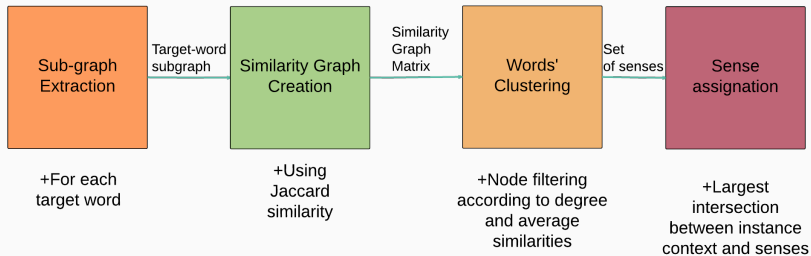
- **Language networks tend to be scale-free**
 - There are certain nodes (hubs) that are very well connected forming communities within the network
- **Seminal approaches**
 - Hyperlex [VÓ4]
 - University of York (UoY) [KM07]
- **Limitations of existing approaches**
 - Single typed networks

- **Language networks tend to be scale-free**
 - There are certain nodes (hubs) that are very well connected forming communities within the network
- **Seminal approaches**
 - Hyperlex [VÓ4]
 - University of York (UoY) [KM07]
- **Limitations of existing approaches**
 - Single typed networks
 - Large number of parameters

- **Language networks tend to be scale-free**
 - There are certain nodes (hubs) that are very well connected forming communities within the network
- **Seminal approaches**
 - Hyperlex [VÓ4]
 - University of York (UoY) [KMo7]
- **Limitations of existing approaches**
 - Single typed networks
 - Large number of parameters
- **Proposition**

- **Language networks tend to be scale-free**
 - There are certain nodes (hubs) that are very well connected forming communities within the network
- **Seminal approaches**
 - Hyperlex [VÓ4]
 - University of York (UoY) [KM07]
- **Limitations of existing approaches**
 - Single typed networks
 - Large number of parameters
- **Proposition**
 - Be able to exploit different types of linguistic information (lexical or syntactic co-occurrence)

- **Language networks tend to be scale-free**
 - There are certain nodes (hubs) that are very well connected forming communities within the network
- **Seminal approaches**
 - Hyperlex [VÓ4]
 - University of York (UoY) [KMo7]
- **Limitations of existing approaches**
 - Single typed networks
 - Large number of parameters
- **Proposition**
 - Be able to exploit different types of linguistic information (lexical or syntactic co-occurrence)
 - Keep the number of parameters low and allow for their automatic adjusting according to the network's nature



Applications to NLP

Hypergraph Model Instantiation

Hypergraph Model Instantiation

- **Apply our proposed linguistic model to a real world corpus**

Hypergraph Model Instantiation

- **Apply our proposed linguistic model to a real world corpus**
 - Use the English Wikipedia as input and generate a textual structure following the proposed network model

Hypergraph Model Instantiation

- **Apply our proposed linguistic model to a real world corpus**
 - Use the English Wikipedia as input and generate a textual structure following the proposed network model
- **We provide two resources**

Hypergraph Model Instantiation

- **Apply our proposed linguistic model to a real world corpus**
 - Use the English Wikipedia as input and generate a textual structure following the proposed network model
- **We provide two resources**
 - A syntactically annotated English Wikipedia corpus (SAEWD)

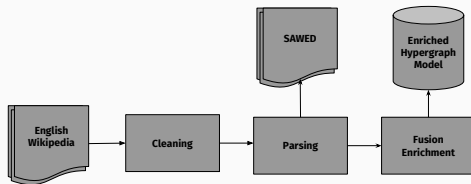
- **Apply our proposed linguistic model to a real world corpus**
 - Use the English Wikipedia as input and generate a textual structure following the proposed network model
- **We provide two resources**
 - A syntactically annotated English Wikipedia corpus (SAEWD)
 - An Wikipedia-based enriched hypergraph linguistic model

Hypergraph Model Instantiation

- **Apply our proposed linguistic model to a real world corpus**
 - Use the English Wikipedia as input and generate a textual structure following the proposed network model
- **We provide two resources**
 - A syntactically annotated English Wikipedia corpus (SAEWD)
 - An Wikipedia-based enriched hypergraph linguistic model
- **Steps performed**

Hypergraph Model Instantiation

- **Apply our proposed linguistic model to a real world corpus**
 - Use the English Wikipedia as input and generate a textual structure following the proposed network model
- **We provide two resources**
 - A syntactically annotated English Wikipedia corpus (SAEWD)
 - An Wikipedia-based enriched hypergraph linguistic model
- **Steps performed**



Hypergraph Model Instantiation

SAEWD: Parsed sample

FILENAME wiki.oo.parsed					
token	lemma	POS	constituency	head	dependency
%%#PAGE Anarchism					
⋮	⋮	⋮	⋮	⋮	⋮
%%#SEN 25 9					
A	a	DT	NP_22,S_97	3	det
great	great	JJ	NP_22,S_97	3	amod
brigand	brigand	NN	NP_22,S_97	4	nsubj
becomes	become	VBZ	VP_44,S_97	0	root
a	a	DT	NP_18,NP_20,VP_44,S_97	6	det
ruler	ruler	NN	NP_18,NP_20,VP_44,S_97	4	xcomp
of	of	IN	PP_57,NP_20,VP_44,S_97	9	case
a	a	DT	NP_18,PP_57,NP_20,VP_44,S_97	9	det
Nation	nation	NN	NP_18,PP_57,NP_20,VP_44,S_97	6	nmod

Hypergraph Model Instantiation

Hypergraph Incidence Matrix

		CONSTITUENT			DEPENDENCY	SENTENCE
		NP ₁ DT:NN	NP ₂ NP:PP:PP	NP ₃ NNS	nsubj contains <i>dobj</i> <i>contains</i>	S ₁
NN	report	1			1	1
	copies		1	1	1	1
	minutes		1			1
	meetings		1			1
VB	contains					1

	Lexical Features (5.49%) M^L	Syntactic Features (4.97%) M^S	Early Fusion (5.23%) $E(M^L, M^S)$	X_F Fusion (16.75%) $X_F(S^S, M^L)$	X_F Fusion (13.45%) $X_F(S^L, M^S)$
priest	priests	monk	sailor	vassal	sailor
	nun	regent	regent	regent	fluent
	canton	aedile	nuclei	nun	dean
	sailor	seer	nun	sailor	nuclei
	burial	meek	relic	monk	chorus

Applications to NLP

Solving Named Entity Recognition

- **NER Objective**

- **NER Objective**

- The goal is to automatically discover mentions that belong to a well-defined semantic category.

- **NER Objective**

- The goal is to automatically discover mentions that belong to a well-defined semantic category.

- **Classic entities types**

- **NER Objective**

- The goal is to automatically discover mentions that belong to a well-defined semantic category.

- **Classic entities types**

- Location (LOC)

- **NER Objective**

- The goal is to automatically discover mentions that belong to a well-defined semantic category.

- **Classic entities types**

- Location (LOC)
- Organization (ORG)

- **NER Objective**

- The goal is to automatically discover mentions that belong to a well-defined semantic category.

- **Classic entities types**

- Location (LOC)
- Organization (ORG)
- Person (PER)

- **NER Objective**

- The goal is to automatically discover mentions that belong to a well-defined semantic category.

- **Classic entities types**

- Location (LOC)
- Organization (ORG)
- Person (PER)
- Miscellaneous (MISC)

- **NER Objective**

- The goal is to automatically discover mentions that belong to a well-defined semantic category.

- **Classic entities types**

- Location (LOC)
- Organization (ORG)
- Person (PER)
- Miscellaneous (MISC)
- **None (O)**

- **NER Objective**

- The goal is to automatically discover mentions that belong to a well-defined semantic category.

- **Classic entities types**

- Location (LOC)
- Organization (ORG)
- Person (PER)
- Miscellaneous (MISC)
- None (O)

- **Our goal**

- **NER Objective**

- The goal is to automatically discover mentions that belong to a well-defined semantic category.

- **Classic entities types**

- Location (LOC)
- Organization (ORG)
- Person (PER)
- Miscellaneous (MISC)
- None (O)

- **Our goal**

- We assess the effectiveness of the classic fusion methods and propose new hybrid combinations

Representation Spaces

Lexical Space (L)

Word	Features
Australian	word:Australian, word+1:scientist, word+2:discovers
scientist	word-1:Australian, word:scientist, word+1:discovers, word+2:star
discovers	word-2:Australian, word-1:scientist, . . . , word+2:telescope
star	word-2:scientist, word-1:discovers, word:star, . . . , word+2:telescope
with	word-2:discovers, word-1:star, word:with, word+1:telescope
telescope	word-2:star, word-1:with, word:telescope

Syntactic Space (S)

Word	Contexts
Australian	scientist/NN/amod_inv
scientist	Australian/JJ/amod, discovers/VBZ/nsubj_inv
discovers	scientist/NN/nsubj, star/NN/dobj, telescope/NN/nmod:with
star	discovers/VBZ/dobj_inv
telescope	discovers/VBZ/nmod:with_inv

Standard Features Space (T)

- Each word
- Whether it is capitalized
- Prefix and suffix (of each word their surroundings)
- Part of Speech tag

Solving Named Entity Recognition

Experimental Protocol

- **Preprocessing**

Solving Named Entity Recognition

Experimental Protocol

- **Preprocessing**
 - Normalize numbers

Experimental Protocol

- **Preprocessing**
 - Normalize numbers
- **Test Corpora**

Solving Named Entity Recognition

Experimental Protocol

- **Preprocessing**
 - Normalize numbers
- **Test Corpora**
 - CoNLL-2003 (CONLL) [SM03]: Train: 219,554 lines. Test: 50,350

Solving Named Entity Recognition

Experimental Protocol

- **Preprocessing**
 - Normalize numbers
- **Test Corpora**
 - CoNLL-2003 (CONLL) [SM03]: Train: 219,554 lines. Test: 50,350
 - Wikiner (WNER) [NMC09]: No Train/Test split. 3.5 million words.
Evaluated in a 5-fold CV

Solving Named Entity Recognition

Experimental Protocol

- **Preprocessing**

- Normalize numbers

- **Test Corpora**

- CoNLL-2003 (CONLL) [SM03]: Train: 219,554 lines. Test: 50,350
- Wikiner (WNER) [NMCo9]: No Train/Test split. 3.5 million words.
Evaluated in a 5-fold CV
- Wikigold (WGLD) [Bal+09]: No Train/Test split. 41,011 words.
Evaluated in a 5-fold CV

Solving Named Entity Recognition

Experimental Protocol

- **Preprocessing**

- Normalize numbers

- **Test Corpora**

- CoNLL-2003 (CONLL) [SM03]: Train: 219,554 lines. Test: 50,350
- Wikiner (WNER) [NMC09]: No Train/Test split. 3.5 million words.
Evaluated in a 5-fold CV
- Wikigold (WGLD) [Bal+09]: No Train/Test split. 41,011 words.
Evaluated in a 5-fold CV

- **Annotation Scheme**

Solving Named Entity Recognition

Experimental Protocol

- **Preprocessing**

- Normalize numbers

- **Test Corpora**

- CoNLL-2003 (CONLL) [SM03]: Train: 219,554 lines. Test: 50,350
- Wikiner (WNER) [NMCo9]: No Train/Test split. 3.5 million words.
Evaluated in a 5-fold CV
- Wikigold (WGLD) [Bal+09]: No Train/Test split. 41,011 words.
Evaluated in a 5-fold CV

- **Annotation Scheme**

- **Beginning, Inside, Outside**

Solving Named Entity Recognition

Experimental Protocol

- **Preprocessing**

- Normalize numbers

- **Test Corpora**

- CoNLL-2003 (CONLL) [SM03]: Train: 219,554 lines. Test: 50,350
- Wikiner (WNER) [NMCo9]: No Train/Test split. 3.5 million words.
Evaluated in a 5-fold CV
- Wikigold (WGLD) [Bal+09]: No Train/Test split. 41,011 words.
Evaluated in a 5-fold CV

- **Annotation Scheme**

- **B**eginning, **I**nside, **O**utside

- **Learning Algorithm**

Solving Named Entity Recognition

Experimental Protocol

- **Preprocessing**

- Normalize numbers

- **Test Corpora**

- CoNLL-2003 (CONLL) [SM03]: Train: 219,554 lines. Test: 50,350
- Wikiner (WNER) [NMCo9]: No Train/Test split. 3.5 million words.
Evaluated in a 5-fold CV
- Wikigold (WGLD) [Bal+09]: No Train/Test split. 41,011 words.
Evaluated in a 5-fold CV

- **Annotation Scheme**

- **B**eginning, **I**nside, **O**utside

- **Learning Algorithm**

- Structured Perceptron [Colo2]

Solving Named Entity Recognition

Experimental Protocol

- **Preprocessing**

- Normalize numbers

- **Test Corpora**

- CoNLL-2003 (CONLL) [SM03]: Train: 219,554 lines. Test: 50,350
- Wikiner (WNER) [NMC09]: No Train/Test split. 3.5 million words.
Evaluated in a 5-fold CV
- Wikigold (WGLD) [Bal+09]: No Train/Test split. 41,011 words.
Evaluated in a 5-fold CV

- **Annotation Scheme**

- **B**eginning, **I**nside, **O**utside

- **Learning Algorithm**

- Structured Perceptron [Colo2]

- **Evaluation Metrics**

Solving Named Entity Recognition

Experimental Protocol

- **Preprocessing**

- Normalize numbers

- **Test Corpora**

- CoNLL-2003 (CONLL) [SM03]: Train: 219,554 lines. Test: 50,350
- Wikiner (WNER) [NMCo9]: No Train/Test split. 3.5 million words.
Evaluated in a 5-fold CV
- Wikigold (WGLD) [Bal+09]: No Train/Test split. 41,011 words.
Evaluated in a 5-fold CV

- **Annotation Scheme**

- **B**eginning, **I**nside, **O**utside

- **Learning Algorithm**

- Structured Perceptron [Colo2]

- **Evaluation Metrics**

- Precision, Recall, F-measure

Solving Named Entity Recognition

Evaluation

A	B	Early Fusion (EF)		
		CONLL	WNER	WGLD
M^L	M^S	72.01	70.59	59.38
M^L	M^T	78.13	79.78	61.96
M^S	M^T	77.70	78.10	60.93
M^L	$E(M^S, M^T)$	78.90	80.04	63.20
		Late Fusion (LF)		
		CONLL	WNER	WGLD
S^L	S^S	61.65	58.79	44.29
S^L	S^T	55.64	67.70	48.00
S^S	S^T	50.21	58.41	49.81

Solving Named Entity Recognition

Evaluation

A	B	Early Fusion (EF)		
		CONLL	WNER	WGLD
M^L	M^S	72.01	70.59	59.38
M^L	M^T	78.13	79.78	61.96
M^S	M^T	77.70	78.10	60.93
M^L	$E(M^S, M^T)$	78.90	80.04	63.20
		Late Fusion (LF)		
		CONLL	WNER	WGLD
S^L	S^S	61.65	58.79	44.29
S^L	S^T	55.64	67.70	48.00
S^S	S^T	50.21	58.41	49.81

Cross Feature Fusion (X_{FF})				
		CONLL	WNER	WGLD
S^L	M^T	49.90	70.27	62.69
S^S	M^T	47.27	51.38	48.53
S^T	$b_{X_{FF}}^*$	52.89	62.21	50.15
Cross Similarity Fusion (X_{SF})				
		CONLL	WNER	WGLD
S^L	S^T	27.75	59.12	38.35
S^S	$b_{X_{SF}}^*$	36.87	40.92	39.62
S^T	$b_{X_{SF}}^*$	41.89	52.03	39.92

$$b_{X_{FF}}^* \in \{M^L, M^T\}$$

$$b_{X_{SF}}^* \in \{S^L, S^S\}$$

Solving Named Entity Recognition

Evaluation

$$E(M_L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$$

Solving Named Entity Recognition

Evaluation



		Triple Early Double Late Cross Feature Fusion (EEELX _F LX _F)		
		CONLL	WNER	WGLD
M^L	$\hat{b}_{EEELX_F LX_F}$	65.01	78.02	62.34
$M^L_{\alpha=0.95}$	$\hat{b}_{EEELX_F LX_F}$	79.67	81.79	67.05
EF Baseline		78.90	80.04	63.20

Analyzing the Best Fusion Operator

- **Understand how the evolution towards and enriched space helps the model take the correct decision**

Analyzing the Best Fusion Operator

- **Understand how the evolution towards and enriched space helps the model take the correct decision**
 - Decompose the large fusion operator into 4 separate representations

Analyzing the Best Fusion Operator

- **Understand how the evolution towards and enriched space helps the model take the correct decision**
 - Decompose the large fusion operator into 4 separate representations
 - Train a model with each individual operator (4 models: M_1 , M_2 , M_3 , M_4)

Analyzing the Best Fusion Operator

- **Understand how the evolution towards and enriched space helps the model take the correct decision**
 - Decompose the large fusion operator into 4 separate representations
 - Train a model with each individual operator (4 models: M_1 , M_2 , M_3 , M_4)
 - Investigate how the features added at each step help the model predict the correct class

Analyzing the Best Fusion Operator

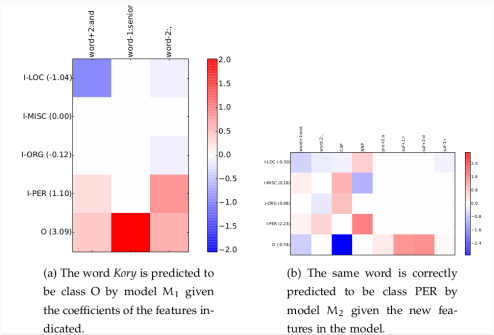
- **Understand how the evolution towards and enriched space helps the model take the correct decision**
 - Decompose the large fusion operator into 4 separate representations
 - Train a model with each individual operator (4 models: M_1, M_2, M_3, M_4)
 - Investigate how the features added at each step help the model predict the correct class

$$E_{\alpha=0.95}(\underbrace{M^L}_{\textcircled{1}}, \underbrace{M^T}_{\textcircled{2}}, \underbrace{L(M^T, X_F(S^S, M^T))}_{\textcircled{3}}, \underbrace{L(M^L, X_F(S^S, M^L))}_{\textcircled{4}})$$

Solving Named Entity Recognition

Analyzing the Best Fusion Operator

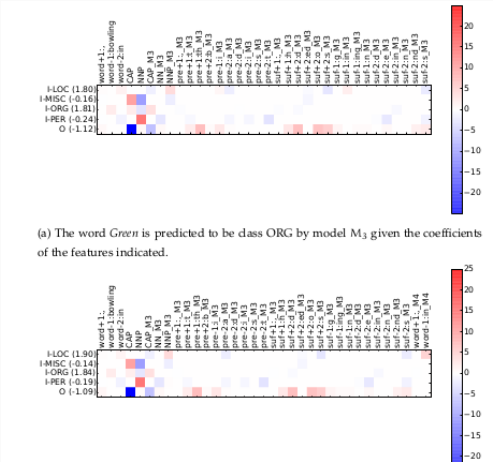
We focus on the word *Kory*, and its performance from model M_1 to M_2



Solving Named Entity Recognition

Analyzing the Best Fusion Operator

We focus on the word *Green*, and its performance from model M_3 to M_4



Applications to NLP

Solving Word Sense Induction and Disambiguation

Introduction

- **WSI/WSD Objective**

- **WSI/WSD Objective**

- The goal is to determine a set of possible senses to a given word according to its possible contexts (WSI). Then, assigning a correct sense to a particular instance of said word

- **WSI/WSD Objective**

- The goal is to determine a set of possible senses to a given word according to its possible contexts (WSI). Then, assigning a correct sense to a particular instance of said word

- **Our goal**

- **WSI/WSD Objective**

- The goal is to determine a set of possible senses to a given word according to its possible contexts (WSI). Then, assigning a correct sense to a particular instance of said word

- **Our goal**

- Again, to assess the effectiveness of the fusion enriched spaces and to evaluate the pertinence of our community discovering algorithm

- **Preprocessing**
 - Remove very frequent and very infrequent words
- **Test Corpora**
 - Semeval 2007 [SM03]: Train: 219,554 lines. Test: 50,350
- **Clustering Algorithm**
 - Spectral Clustering [SM00]
 - Proposed Community Algorithm
- **Evaluation Metrics**
 - Supervised Recall
 - Unsupervised F-measure
 - Proposed: H-measure

$$\text{H-measure} = \frac{1}{2} \left(2 * \frac{\text{SR} * \text{UF}}{\text{SR} + \text{UF}} + \frac{\delta}{\delta + |\#\text{cl} - \delta|} \right)$$

δ is the average true number of senses of the words in a test corpus

Solving Word Sense Induction and Disambiguation

Spectral Clustering Evaluation

Cross Feature Cross Similarity Fusion (X _F X _S F)					
X _F (X _S (S ¹ , S ⁵), M ¹)	78.40	80.40	76.10	3.11	
X _F (X _S (S ¹ , S ⁵), M ⁵)	78.90	81.80	75.60	3.16	
Early Cross Feature Fusion (EX _F F)					
E(M ¹ , X _F (S ¹ , M ¹))	79.20	82.40	75.70	3.57	2F
E(M ⁵ , X _F (S ¹ , M ¹))	78.30	80.50	75.80	1.95	
Late Cross Feature Fusion (LX _F F)					
L(M ⁵ , X _F (S ¹ , M ⁵))	78.60	81.10	75.80	4.22	
L(M ¹ , X _F (S ¹ , M ¹))	79.50	82.80	75.70	3.96	
Early Late Cross Feature Fusion (ELX _F F)					
E(M ¹ , L(M ⁵ , X _F (S ¹ , M ⁵)))	78.50	81.40	75.40	4.26	HF
E(M ¹ , L(M ¹ , X _F (S ¹ , M ¹)))	79.50	82.70	75.90	3.99	
Baseline MFS	78.70	80.90	76.20	1.00	

Figure 1: Supervised Recall

Solving Word Sense Induction and Disambiguation

Spectral Clustering Evaluation

Cross Feature Cross Similarity Fusion (X _F X _S F)				
X _F (X _S (S ^L , S ^S), M ^L)	78.40	80.40	76.10	3.11
X _F (X _S (S ^L , S ^S), M ^S)	78.90	81.80	75.60	3.16
Early Cross Feature Fusion (EX _F F)				
E(M ^L , X _F (S ^L , M ^L))	79.20	82.40	75.70	3.57
E(M ^S , X _F (S ^L , M ^L))	78.30	80.50	75.80	1.95
Late Cross Feature Fusion (LX _F F)				
L(M ^S , X _F (S ^L , M ^S))	78.60	81.10	75.80	4.22
L(M ^L , X _F (S ^L , M ^L))	79.50	82.80	75.70	3.96
Early Late Cross Feature Fusion (ELX _F F)				
E(M ^L , L(M ^S , X _F (S ^L , M ^S)))	78.50	81.40	75.40	4.26
E(M ^L , L(M ^L , X _F (S ^L , M ^L)))	79.50	82.70	75.90	3.99
Baseline MFS	78.70	80.90	76.20	1.00

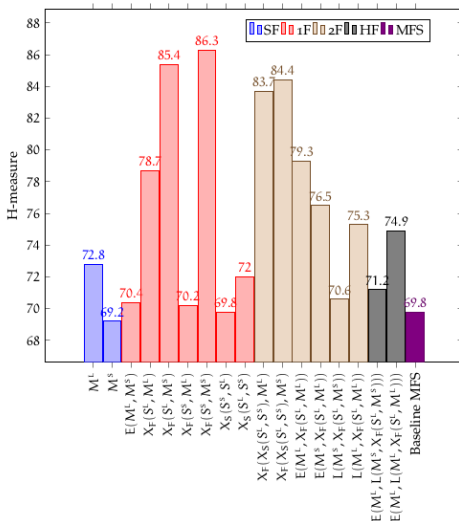
Figure 1: Supervised Recall

Early Fusion (EF)				
E(M ^L , M ^S)	74.00	76.66	71.11	4.46
Cross Feature Fusion (X _F F)				
X _F (S ^L , M ^L)	76.20	79.60	72.50	3.63
X _F (S ^L , M ^S)	74.60	75.10	73.90	3.08
X _F (S ^S , M ^L)	78.90	80.70	76.90	1.08
X _F (S ^S , M ^S)	73.70	77.70	70.00	2.72
Cross Similarity Fusion (X _S F)				
X _S (S ^S , S ^L)	78.90	80.80	76.80	1.01
X _S (S ^L , S ^S)	78.70	80.50	76.80	1.33

Figure 2: Unsupervised F-measure

Solving Word Sense Induction and Disambiguation

Spectral Clustering Evaluation



Solving Word Sense Induction and Disambiguation

Proposed Algorithm Evaluation

	Early Fusion (EF)			
$E(M^L, M^S)$	78.80	81.00	76.40	2.43
	Cross Feature Fusion (X_F F)			
$X_F(S^L, M^L)$	78.70	80.90	76.20	3.11
$X_F(S^L, M^S)$	78.50	81.10	75.60	1.92
$X_F(S^S, M^L)$	79.10	81.60	76.40	1.73
$X_F(S^S, M^S)$	78.60	80.90	76.00	1.81
	Cross Similarity Fusion (X_S F)			
$X_S(S^S, S^L)$	78.60	80.80	76.20	1.44
$X_S(S^L, S^S)$	78.70	80.90	76.20	1.10

Figure 4: Supervised Recall

Solving Word Sense Induction and Disambiguation

Proposed Algorithm Evaluation

Early Fusion (EF)				
$E(M^L, M^S)$	78.80	81.00	76.40	2.43
Cross Feature Fusion ($X_F F$)				
$X_F(S^L, M^L)$	78.70	80.90	76.20	3.11
$X_F(S^L, M^S)$	78.50	81.10	75.60	1.92
$X_F(S^S, M^L)$	79.10	81.60	76.40	1.73
$X_F(S^S, M^S)$	78.60	80.90	76.00	1.81
Cross Similarity Fusion ($X_S F$)				
$X_S(S^S, S^L)$	78.60	80.80	76.20	1.44
$X_S(S^L, S^S)$	78.70	80.90	76.20	1.10

Figure 4: Supervised Recall

Early Fusion (EF)				
$E(M^L, M^S)$	76.90	80.20	73.10	2.43
Cross Feature Fusion ($X_F F$)				
$X_F(S^L, M^L)$	71.00	68.10	74.20	3.11
$X_F(S^L, M^S)$	77.70	79.60	75.50	1.92
$X_F(S^S, M^L)$	75.20	75.50	74.90	1.73
$X_F(S^S, M^S)$	77.60	80.50	74.30	1.81
Cross Similarity Fusion ($X_S F$)				
$X_S(S^S, S^L)$	74.10	72.10	76.50	1.44
$X_S(S^L, S^S)$	78.30	79.70	76.80	1.10

Figure 5: Unsupervised F-measure

Solving Word Sense Induction and Disambiguation

Proposed Algorithm Evaluation

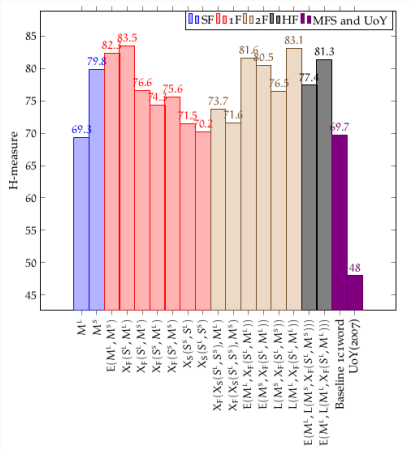


Figure 6: Proposed H-measure

Conclusions

Insights From our Contributions

- **Hypergraph Linguistic Model**

Insights From our Contributions

- **Hypergraph Linguistic Model**
 - Considering heterogeneous features to link words together at once using a hypergraph structure

Insights From our Contributions

- **Hypergraph Linguistic Model**
 - Considering heterogeneous features to link words together at once using a hypergraph structure
 - Yields a multi-layered representation of text

Insights From our Contributions

- **Hypergraph Linguistic Model**
 - Considering heterogeneous features to link words together at once using a hypergraph structure
 - Yields a multi-layered representation of text
- **Combining Features and Dealing with Sparsity**

Insights From our Contributions

- **Hypergraph Linguistic Model**
 - Considering heterogeneous features to link words together at once using a hypergraph structure
 - Yields a multi-layered representation of text
- **Combining Features and Dealing with Sparsity**
 - Using fusion operators

Insights From our Contributions

- **Hypergraph Linguistic Model**
 - Considering heterogeneous features to link words together at once using a hypergraph structure
 - Yields a multi-layered representation of text
- **Combining Features and Dealing with Sparsity**
 - Using fusion operators
 - Intuitive way to leverage the different points of view of each heterogeneous feature while increasing the density of the representation

Insights From our Contributions

- **Hypergraph Linguistic Model**
 - Considering heterogeneous features to link words together at once using a hypergraph structure
 - Yields a multi-layered representation of text
- **Combining Features and Dealing with Sparsity**
 - Using fusion operators
 - Intuitive way to leverage the different points of view of each heterogeneous feature while increasing the density of the representation
- **Applications to NLP**

Insights From our Contributions

- **Hypergraph Linguistic Model**
 - Considering heterogeneous features to link words together at once using a hypergraph structure
 - Yields a multi-layered representation of text
- **Combining Features and Dealing with Sparsity**
 - Using fusion operators
 - Intuitive way to leverage the different points of view of each heterogeneous feature while increasing the density of the representation
- **Applications to NLP**
 - Solving NER and WSI/WSD with fusion enriched representations and our community-driven algorithm

- **Hypergraph Linguistic Model**
 - Considering heterogeneous features to link words together at once using a hypergraph structure
 - Yields a multi-layered representation of text
- **Combining Features and Dealing with Sparsity**
 - Using fusion operators
 - Intuitive way to leverage the different points of view of each heterogeneous feature while increasing the density of the representation
- **Applications to NLP**
 - Solving NER and WSI/WSD with fusion enriched representations and our community-driven algorithm
 - A high degree combination of fusion operators are the ones that yield the improvements

- **Hypergraph Linguistic Model**
 - Considering heterogeneous features to link words together at once using a hypergraph structure
 - Yields a multi-layered representation of text
- **Combining Features and Dealing with Sparsity**
 - Using fusion operators
 - Intuitive way to leverage the different points of view of each heterogeneous feature while increasing the density of the representation
- **Applications to NLP**
 - Solving NER and WSI/WSD with fusion enriched representations and our community-driven algorithm
 - A high degree combination of fusion operators are the ones that yield the improvements
 - The community finding algorithm improves over similar algorithms while being simpler and allows for heterogeneous features

- **Hypergraph Linguistic Model**
 - Considering heterogeneous features to link words together at once using a hypergraph structure
 - Yields a multi-layered representation of text
- **Combining Features and Dealing with Sparsity**
 - Using fusion operators
 - Intuitive way to leverage the different points of view of each heterogeneous feature while increasing the density of the representation
- **Applications to NLP**
 - Solving NER and WSI/WSD with fusion enriched representations and our community-driven algorithm
 - A high degree combination of fusion operators are the ones that yield the improvements
 - The community finding algorithm improves over similar algorithms while being simpler and allows for heterogeneous features
 - The Wikipedia-based instantiation serves as a NLP system starting point

- **Hypergraph Linguistic Model**

- **Hypergraph Linguistic Model**
 - Implementing a dataframe-like structure allowing for queries and exploration of large corpora using the proposed model

- **Hypergraph Linguistic Model**
 - Implementing a dataframe-like structure allowing for queries and exploration of large corpora using the proposed model
- **Combining Features and Dealing with Sparsity**

- **Hypergraph Linguistic Model**
 - Implementing a dataframe-like structure allowing for queries and exploration of large corpora using the proposed model
- **Combining Features and Dealing with Sparsity**
 - Finding a more principled way to determine what type of context with what type of fusion operation according to the task at hand

- **Hypergraph Linguistic Model**
 - Implementing a dataframe-like structure allowing for queries and exploration of large corpora using the proposed model
- **Combining Features and Dealing with Sparsity**
 - Finding a more principled way to determine what type of context with what type of fusion operation according to the task at hand
 - Exploring with other modal features

- **Hypergraph Linguistic Model**
 - Implementing a dataframe-like structure allowing for queries and exploration of large corpora using the proposed model
- **Combining Features and Dealing with Sparsity**
 - Finding a more principled way to determine what type of context with what type of fusion operation according to the task at hand
 - Exploring with other modal features
- **Applications to NLP**

- **Hypergraph Linguistic Model**
 - Implementing a dataframe-like structure allowing for queries and exploration of large corpora using the proposed model
- **Combining Features and Dealing with Sparsity**
 - Finding a more principled way to determine what type of context with what type of fusion operation according to the task at hand
 - Exploring with other modal features
- **Applications to NLP**
 - Using the large Wikipedia-based network as a background corpus to further enrich domain-specific corpora

- **Hypergraph Linguistic Model**
 - Implementing a dataframe-like structure allowing for queries and exploration of large corpora using the proposed model
- **Combining Features and Dealing with Sparsity**
 - Finding a more principled way to determine what type of context with what type of fusion operation according to the task at hand
 - Exploring with other modal features
- **Applications to NLP**
 - Using the large Wikipedia-based network as a background corpus to further enrich domain-specific corpora
 - Test more feature weighting schemes, validate findings on more datasets

Publications Produced by Our Research

- Edmundo-Pavel Soriano-Morales, Julien Ah-Pine, Sabine Loudcher: **Fusion Techniques for Named Entity Recognition and Word Sense Induction and Disambiguation**. DS 2017
- Edmundo-Pavel Soriano-Morales, Julien Ah-Pine, Sabine Loudcher: **Using a Heterogeneous Linguistic Network for Word Sense Induction and Disambiguation**. CICLING 2016
- Edmundo-Pavel Soriano-Morales, Julien Ah-Pine, Sabine Loudcher: **Hypergraph Modelization of a Syntactically Annotated English Wikipedia Dump**. LREC 2016
- Adrien Guille, Edmundo-Pavel Soriano-Morales, Ciprian-Octavian Truica: **Topic modeling and hypergraph mining to analyze the EGC conference history**. EGC 2016
- Adrien Guille, Edmundo-Pavel Soriano-Morales: **TOM: A library for topic modeling and browsing**. EGC 2016:

- Julien Ah-Pine, Edmundo-Pavel Soriano-Morales: **A Study of Synthetic Oversampling for Twitter Imbalanced Sentiment Analysis**. DMNLP@PKDD/ECML 2016
- Sabine Loudcher, Wararat Jakawat, Edmundo-Pavel Soriano-Morales, Ccile Favre: **Combining OLAP and information networks for bibliographic data analysis: a survey**. Scientometrics 103(2)

Thank you for your attention

References



Christopher D Manning, Hinrich Schütze, et al.
Foundations of statistical natural language processing.
Vol. 999. MIT Press, 1999.



Jianbo Shi and Jitendra Malik. “Normalized Cuts and Image Segmentation”. In: *IEEE Trans. Pattern Anal. Mach. Intell.* 22.8 (Aug. 2000), pp. 888–905. ISSN: 0162-8828. DOI: 10.1109/34.868688. URL: <http://dx.doi.org/10.1109/34.868688>.



Michael Collins. “Discriminative Training Methods for Hidden Markov Models: Theory and Experiments with Perceptron Algorithms”. In: *Proceedings of the ACL-02 Conference on Empirical Methods in Natural Language Processing - Volume 10*. EMNLP '02. Stroudsburg, PA, USA: Association for Computational Linguistics, 2002, pp. 1–8. DOI: 10.3115/1118693.1118694.



Erik F. Tjong Kim Sang and Fien De Meulder. “Introduction to the CoNLL-2003 Shared Task: Language-Independent Named Entity Recognition”. In: *CoNLL*. ACL, 2003, pp. 142–147.



Rada Mihalcea, Paul Tarau, and Elizabeth Figa. “PageRank on Semantic Networks, with Application to Word Sense Disambiguation”. In: *Proceedings of the 20th International Conference on Computational Linguistics*. COLING '04. Geneva, Switzerland: Association for Computational Linguistics, 2004. DOI: 10.3115/1220355.1220517.



Jean Véronis. “HyperLex: lexical cartography for information retrieval”. In: *Computer Speech & Language* 18.3 (2004), pp. 223 –252. ISSN: 0885-2308. DOI: 10.1016/j.cs1.2004.05.002.



Ioannis P. Klapaftis and Suresh Manandhar. “UOY: A Hypergraph Model for Word Sense Induction & Disambiguation”. In: *Proceedings of the 4th International Workshop on Semantic Evaluations*. SemEval '07. Prague, Czech Republic: Association for Computational Linguistics, 2007, pp. 414–417.



Ioannis P. Klapaftis and Suresh Manandhar. “Word Sense Induction Using Graphs of Collocations”. In: *Proceedings of the 2008 Conference on ECAI 2008: 18th European Conference on Artificial Intelligence*. Amsterdam, The Netherlands, The Netherlands: IOS Press, 2008, pp. 298–302. ISBN: 978-1-58603-891-5.



Dominic Balasuriya et al. “Named Entity Recognition in Wikipedia”. In: *Proceedings of the 2009 Workshop on The People’s Web Meets NLP: Collaboratively Constructed Semantic Resources*. People’s Web ’09. Suntec, Singapore: Association for Computational Linguistics, 2009, pp. 10–18. ISBN: 978-1-932432-55-8. URL: <http://dl.acm.org/citation.cfm?id=1699765.1699767>.



Monojit Choudhury and Animesh Mukherjee. “The Structure and Dynamics of Linguistic Networks”. English. In: *Dynamics On and Of Complex Networks*. Ed. by Niloy Ganguly, Andreas Deutsch, and Animesh Mukherjee. Modeling and Simulation in Science, Engineering and Technology. Birkhäuser Boston, 2009, pp. 145–166. ISBN: 978-0-8176-4750-6. DOI: 10.1007/978-0-8176-4751-3_9.



Joel Nothman, Tara Murphy, and James R. Curran. “Analysing Wikipedia and Gold-standard Corpora for NER Training”. In: *Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics*. EACL '09. Athens, Greece: Association for Computational Linguistics, 2009, pp. 612–620.



Yong-Yeol Ahn, James P Bagrow, and Sune Lehmann. “Link communities reveal multiscale complexity in networks”. In: *Nature* 466.7307 (2010), pp. 761–764.



Pradeep K. Atrey et al. “Multimodal fusion for multimedia analysis: a survey”. In: *Multimedia Syst.* 16.6 (2010), pp. 345–379.



Rada F. Mihalcea and Dragomir R. Radev. *Graph-based Natural Language Processing and Information Retrieval*. 1st. New York, NY, USA: Cambridge University Press, 2011. ISBN: 0521896134, 9780521896139.