


วิธีการกำจัดความกำกวมของคำหลายความหมายโดยใช้ความหมายเชิงการกระจายและความหมายแฝง



นางสาวสุนีย์ พงษ์พินิจภิญโญ

ศูนย์วิทยพัทยากร

จุฬาลงกรณ์มหาวิทยาลัย

วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิศวกรรมศาสตรดุษฎีบัณฑิต

สาขาวิชาวิศวกรรมคอมพิวเตอร์ ภาควิชาวิศวกรรมคอมพิวเตอร์

คณะวิศวกรรมศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย


ปีการศึกษา 2547

ISBN 974-53-1612-1

ลิขสิทธิ์ของจุฬาลงกรณ์มหาวิทยาลัย

A DISTRIBUTIONAL SEMANTICS AND LATENT SEMANTICS APPROACH  
FOR WORD SENSE DISAMBIGUATION

Miss Sunee Pongpinigpinyo



ศูนย์วิทยทรัพยากร  
จุฬาลงกรณ์มหาวิทยาลัย

A Dissertation Submitted in Partial Fulfillment of the Requirements  
for the Degree of Doctor of Philosophy in Computer Engineering

Department of Computer Engineering

Faculty of Engineering

Chulalongkorn University


Academic year 2004

ISBN 974-53-1612-1

Thesis Title                    A DISTRIBUTIONAL SEMANTICS AND LATENT  
SEMANTICS APPROACH FOR WORD SENSE  
DISAMBIGUATION  
By                                    Miss Sunee Pongpinigpinyo  
Field of Study                    Computer Engineering  
Thesis Advisor                   Associate Professor Dr.Wanchai Rivepiboon

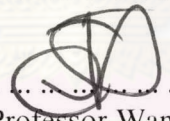
---

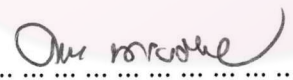
Accepted by the Faculty of Engineering, Chulalongkorn University in  
Partial Fulfillment of the Requirements for the Doctor's Degree

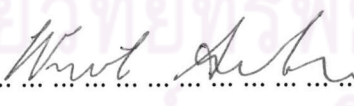
  
..... Dean of the Faculty of Engineering  
(Professor Direk Lavansiri, Ph.D.)


THESIS COMMITTEE

  
..... Chairman  
(Associate Professor Boonserm Kijisirikul, D.Eng.)

  
..... Thesis Advisor  
(Associate Professor Wanchai Rivepiboon, Ph.D.)

  
..... Member  
(Yunyong Teng-amnuay, Ph.D.)

  
..... Member  
(Assistant Professor Wirote Aroonmanakun, Ph.D.)

  
..... Member  
(Assistant Professor Arnon Rungsawang, Ph.D.)

ศุณีย์ พงษ์พินิจภิญโญ : วิธีการกำจัดความกำกวมของคำหลายความหมายโดยใช้ความหมายเชิงการกระจายและความหมายแฝง. (A DISTRIBUTIONAL SEMANTICS AND LATENT SEMANTICS APPROACH FOR WORD SENSE DISAMBIGUATION) อ. ที่ปรึกษา : รศ. ดร. วันชัย รั้วไพบุลย์ จำนวนหน้า 86 หน้า. ISBN 974-53-1612-1.

วิทยานิพนธ์นี้มีวัตถุประสงค์เพื่อพัฒนาวิธีการแก้ปัญหาความกำกวมของคำหลายความหมายในภาษาไทยโดยใช้วิธีการเรียนรู้แบบไม่ต้องใช้ผู้สอน (unsupervised learning technique) ซึ่งไม่จำเป็นต้องใช้แหล่งความรู้ต่างๆมากมาย เช่น ข้อความที่กำกับความหมาย หรือ พจนานุกรม วิธีการที่นำเสนอในวิทยานิพนธ์นี้เป็นวิธีการใหม่เพราะยังไม่เคยมีใครนำวิธีการนี้มาใช้กับภาษาไทยที่ยังไม่ได้รับความสนใจมากนักในงานวิจัยด้านการประมวลผลภาษาธรรมชาติ วิธีการของผู้วิจัยที่นำเสนอคือการใช้ความหมายเชิงกระจายที่มีพื้นฐานจากสมมติฐานของการกระจายที่ว่า คำที่มีความหมายคล้ายกันจะปรากฏอยู่ในบริบทที่คล้ายกัน การนำเสนอวิธีนี้เกี่ยวข้องกับการใช้ฐานข้อมูลของภาษา (corpora) เพื่อที่ตรวจสอบบริบทของแต่ละคำที่ปรากฏอยู่แล้วจึงคำนวณหาความคล้ายกันระหว่างบริบทที่กระจาย กลุ่มความหมายของคำเป็นกลุ่มคำที่เกิดขึ้นด้วยกันเสมอกับคำหลายความหมาย ซึ่งแต่ละกลุ่มความหมายประกอบไปด้วยกลุ่มคำที่เกิดขึ้นด้วยกันเสมอที่ให้ความหมายเหมือนกัน

งานวิทยานิพนธ์นี้ให้ประโยชน์ไม่เพียงสำหรับการพัฒนาขั้นตอนต่อไปของวิธีการแก้ปัญหาความกำกวมของคำหลายความหมายในภาษาไทยแต่ยังสามารถนำไปปรับใช้กับลักษณะของปัญหา (domain) ที่แตกต่างกัน และขอบเขตตัวอย่างข้อมูลขนาดใหญ่ขึ้นได้

ศูนย์วิทยทรัพยากร  
จุฬาลงกรณ์มหาวิทยาลัย

ภาควิชา วิศวกรรมคอมพิวเตอร์  
สาขาวิชา วิศวกรรมคอมพิวเตอร์  
ปีการศึกษา 2547

ลายมือชื่อนิสิต..... กนก นอนโหวน  
ลายมือชื่ออาจารย์ที่ปรึกษา.....

## 4371818121 : MAJOR COMPUTER ENGINEERING

KEY WORD: WORD SENSE DISAMBIGUATION / NATURAL LANGUAGE PROCESSING / DISTRIBUTIONAL SEMANTICS / LATENT SEMANTICS

SUNEE PONGPINIGPINYO : (A DISTRIBUTIONAL SEMANTICS AND LATENT SEMANTICS APPROACH FOR WORD SENSE DISAMBIGUATION. THESIS ADVISOR: ASSOC. PROF. WANCHAI RIVEPIBOON, Ph.D., 86 pp. ISBN 974-53-1612-1.

This thesis aims at developing a methodology of word sense disambiguation in Thai using unsupervised learning techniques that do not rely on any knowledge intensive resources like sense-tagged text or dictionaries. The thesis approach is novel since the work focuses on Thai language which has not received much attention in the Natural Language Processing literature. Our proposed method is the distributional semantics which is based on the *distributional hypothesis* that *similar words appear in similar contexts*. This approach involves using corpora to examine the contexts which each word appears in and then calculating the similarity between context distributions. The sense clusters are occurrences of words that have been grouped into clusters based from raw text where each cluster consists of occurrences having the same meaning.

This thesis work also gives benefit not only to the further development of word sense disambiguation for Thai language but also to be portable to different domains, and can be scaled up easily for larger samples of text.

ศูนย์วิทยทรัพยากร  
จุฬาลงกรณ์มหาวิทยาลัย

Department    Computer Engineering  
Field of study    Computer Engineering  
Academic year    2004

Student's signature... *Sunee Pongpinigpinyo*  
Advisor's signature..... *[Signature]*



## ACKNOWLEDGEMENTS

I sincerely feel that this thesis would not have been possible, without the tremendous help, support, guidance and lastly, but not least, the true friendship and love, given by various people who I have both the honour and the luck to get acquainted with.

Firstly, I would like to thank my supervisor, Associate Professor Dr. Wanchai Rivepiboon, for the guidance and advice. Next, I would like to thank my thesis committee: Associate Professor Dr. Boonserm Kijsirikul, Dr. Yunyong Teng-amnuay, Assistant Professor Dr. Arnon Rungsawang, and Assistant Professor Dr. Wirote Aroonmanakun who provided valuable thesis comments. Especially, a very grateful thank you is given to Assistant Professor Dr. Arnon Rungsawang, without him I would probably never completed this thesis. He gives his time for thesis guidance and discussion on various aspects of the research.

I am also specially thankful to Associate Professor Dr. Prabhas Chongstitvatana, Yodthong Rodkaew and Chaiwat Jessadapakorn who lift me up when I felt down with stress. They give me the strength and encouragement thru all my Ph.D student life. They cheer me up and get me ready for another day of research.

Also I would like to express my gratitude towards my good friends Dr. Judith Young for taking the trouble to help by proof reading drafts of this thesis.

My sincere thanks are given to all fellows of the Software Engineering Laboratory for their helps with programming, technical computer problems and being my recreation. My special thanks are given to Miss Jutapuck Pugsee for her invaluable and kindness helps whenever I ask for.

I am indebted to the Commission on Higher Education Ministry of Education for funding source for my study.

Last but not least, I would like to express my greatest gratitude to the persons to whom I owe the most, my beloved parents who have always given me pure love, devotion and moral support throughout the various stages of my life, understanding and great patience waiting for me during these years to finish my study. Thank you for being who you are; the best parents for me.

Finally, I want to offer my wholehearted praises and thanksgiving to our Lord Buddha for giving enlightenment Dharma to make my body and soul be in peace when I feel suffer from attachment and stress.

# CONTENTS

	Page
ABSTRACT (THAI).....	iv
ABSTRACT (ENGLISH).....	v
ACKNOWLEDGEMENTS.....	vi
CONTENTS.....	vii
LIST OF TABLES.....	x
LIST OF FIGURES.....	xi
CHAPTER 1 INTRODUCTION.....	1
1.1 Motivation.....	1
1.2 Objectives.....	3
1.3 Scopes and Assumption.....	3
2.1.1 Scopes.....	4
2.1.2 Assumption.....	4
1.4 Contribution.....	5
1.5 Organization of the Thesis.....	5
CHAPTER 2 BACKGROUND.....	6
2.1 Theoretical Background.....	6
2.1.1 Vector Space Based Model.....	6
2.1.2 The Advantages and Disadvantage of Vector Space Based Model.....	7
2.2 Latent Semantic Indexing (LSI).....	8
2.3 Data Clustering Methods.....	11
2.3.1 Feature Selection.....	12
2.3.2 Object Representation.....	13
2.3.3 Similarity Measurement and Clustering Method.....	15
2.3.3.1 Similarity Measurement.....	15
2.3.3.2 Clustering Methods.....	16

## CONTENTS (continued)

viii  
Page

2.4 Evaluation.....	20
2.4.1 Precision.....	20
2.4.2 Recall.....	21
2.4.3 F-measure.....	21
CHAPTER 3 LITERATURE REVIEW.....	22
3.1 What is a Word Sense?.....	22
3.2 Word Sense Disambiguation Approaches in Literatures.....	23
3.2.1 Knowledge-Based Approaches.....	24
3.2.2 Corpus-Based Approaches.....	26
3.2.2.1 Supervised Learning Approach.....	28
3.2.2.2 Unsupervised Learning Approach.....	32
3.2.2.3 Semi-Supervised Learning Approach.....	36
3.3 Baseline Evaluation.....	37
3.3.1 An Upper-bound Performance.....	37
3.3.2 A Lower-bound Performance.....	37
3.3.3 Baseline .....	37
3.4 Our Work Approach.....	38
CHAPTER 4 DISTRIBUTIONAL SEMANTICS.....	39
4.1 The Distributional Semantics.....	39
4.2 Types of Features.....	40
4.2.1 Co-occurrences.....	41
4.2.1.1 Co-occurrences within a large window.....	41
4.2.1.2 Co-occurrences within a small window.....	41
4.3 Contextual Representation.....	42
4.3.1 Word Vector.....	42
4.3.2 Dimension Reduction.....	44
4.3.3 Context Vector.....	44
4.4 Sense Clusters.....	45
4.5 Methodology.....	46
4.6 Comparison between Schütze and Thesis Approach.....	49



## CONTENTS (continued)

ix  
Page

CHAPTER 5 EXPERIMENTAL SETTING AND RESULTS.....	50
5.1 Dataset.....	50
5.1.1 Thai Dataset.....	50
5.1.1.1 Source of the Data.....	50
5.1.1.2 Scope of the Data.....	51
5.1.2 English Dataset.....	52
5.2 Experimental Setting.....	52
5.3 Experimental Results.....	58
5.4 Discussion.....	60
5.5 Result Comparison with Other Thai Word Sense Disambiguation Methodology.....	60
5.6 Further Investigated Parameters.....	61
 CHAPTER 6 CONCLUSION AND FUTURE WORKS.....	 64
6.1 Conclusion.....	64
6.2 Future Works.....	65
 REFERENCES.....	 68
 APPENDICES.....	 78
APPENDIX A Sense Definitions.....	78
APPENDIX B Publications.....	85
 BIOGRAPHY.....	 86

## LIST OF TABLES

x

Table	Page
4.1 An example of co-occurrence matrix.....	41
4.2 Summary of the differences between Schütze (Schütze, 1998) and thesis approach.....	49
5.1 Sense Distribution of หัว /hua4/.....	53
5.2 Sense Distribution of เก็บ /kep1/.....	53
5.3 Sense Distribution of <i>interest</i> .....	54
5.4 Co-occurrence Matrix.....	56
5.5 Accuracy of disambiguation of หัว /hua4/.....	59
5.6 Accuracy of disambiguation of เก็บ /kep1/.....	59
5.7 Accuracy of disambiguation of <i>interest</i> .....	60
5.8 Comparison Results of Average Precision rate of disambiguation of หัว /hua4/ and เก็บ /kep1/ between Decision List and our work.....	61
A.1 Definitions, derived senses of หัว /hua4/ found in a corpus.....	78
A.2 Definitions, derived senses of เก็บ /kep1/ found in a corpus.....	82
A.3 Definitions, derived senses of <i>interest</i> found in a corpus.....	84



ศูนย์วิทยทรัพยากร  
จุฬาลงกรณ์มหาวิทยาลัย

## LIST OF FIGURES

xi

Figure	Page
2.1 Reducing a Matrix to K Dimensions with SVD.....	10
2.2 Vector Representations of Objects.....	14
2.3 Point Representations of Objects.....	14
2.4 Example of Dendrogram.....	17
2.5 Single Link Clustering.....	18
2.6 Complete Link Clustering.....	19
2.7 Average Link Clustering.....	19
4.1 An Example of Window Sliding .....	42
4.2 Show the Context of an Example Context of <i>suit</i> .....	45
4.3 The Steps of Methodology .....	46
4.4 The Derivation of Sense Vectors.....	48
5.1 Sense Vectors .....	57
5.2 Compare Context Vector of Test Data with Sense Vector.....	58



ศูนย์วิทยทรัพยากร  
จุฬาลงกรณ์มหาวิทยาลัย