บทคัดย่อ

การประมวลผลภาษาธรรมชาติสำหรับภาษาไทยนั้นเป็นเรื่องที่ยากเนื่องจากภาษาไทยเป็น
ภาษาที่ไม่เข้มงวดในโครงสร้างของไวยกรณ์ แต่เน้นเอาความหมายเป็นหลัก คำที่ถูกสร้างขึ้นมาใหม่ส่วน
ใหญ่เกิดขึ้นจากเป็นการนำเอาคำโดดมาผสมกันเพื่อให้เกิดเป็นคำใหม่ที่อาจจะมีความหมายคล้ายกับ
คำเดิมหรือแตกต่างไป ซึ่งบางครั้งหากต้องการเน้นความหมายก็อาจใช้วิธีเขียนคำเดิมซ้ำหรือเพิ่มคำ
ขยาย หรืออาจจะใช้วิธีการผันเสียงที่สูงหรือหัวนขึ้น สำหรับคำในภาษาไทยจะไม่มีการเปลี่ยนรูปคำ ตาม
กาลเวลา หรือความเป็นเอกพจน์และพนูพจน์ หรือการใช้อักษรตัวใหญ่ ตัวเล็ก เพื่อบ่งบอกคำเฉพาะ ซึ่ง
ทำให้ไม่สามารถวิเคราะห์หน้าที่ของคำจากโครงสร้าง ไวยกรณ์ได้ อีกทั้งคำที่เขียนในภาษาไทยจะ
เขียนในลักษณะที่ติดต่อกันโดยไม่มีช่องว่างระหว่างคำ ซึ่งจำเป็นต้องมีขบวนการการตัด ในส่วนของการ
วิเคราะห์

โครงงานวิจัยนี้จึงมีแนวคิดในการจัดทำเหมืองข้อมูลสำหรับเอกสารงานวิจัยและผลงานทาง
วิชาการของสำนักงานกองทุนสนับสนุนการวิจัย ทั้งที่มีอยู่ในปัจจุบันและในอดีตที่มีมากกว่า 5000 เรื่องที่
อาจอยู่ในไฟล์รูปแบบต่างๆหรืออาจจะยังอยู่ในรูปของกระดาษ โดยจะทำการวิเคราะห์เอกสารโดยการ
ประยุกต์ใช้หลักการประมวลผลภาษาธรรมชาติสำหรับการวิเคราะห์โครงสร้างของข้อความภาษาไทย
เพื่อคัดแยกโครงสร้าง คำสำคัญและจัดเอกสารออกเป็นหมวดหมู่ต่างๆ พร้อมทั้งจัดทำดัชนี สำหรับการ
สืบคัน พร้อมทั้งจัดทำเป็นโปรแกรมประยุกต์เพื่อให้เหมาะสมต่อการนำไปใช้งานและสำหรับผู้ใช้งานใน
ระดับต่างๆ

ลิขสิทธิ์มหาวิทยาลัยเชียงใหม่ Copyright[©] by Chiang Mai University All rights reserved

Abstract

Natural Language Processing (NLP) for Thai language is difficult. This is because Thai language gives less important in the grammartical structure when forming the sentences and rather more concerned about the meaning. Words usually are formed from several root words which provides similar or total new meaning to the original words. Specific meanings are sometimes presented by repeating words or adding more words to extending the meanings, and sometimes by changing the tones. Thai grammar doesn't have word inflection. The word form remain when presenting plurals or tenses. Thai language also has no capital letters for specific words or as beginning of sentences. Words are written altogether without a word break, therefore, word segmentation is needed during the analysis processes.

The objective of this research is to apply natural language processings to analyse the research papers of Thailand Research Funds (TRF) which have been collected for more than 5000 topics. The papers are in different formats and mostly are still in paper formats. The analysis processes consists of the grammatical structure analysis, words and keywords extraction. The documents are then classified into groups with pre-defined indexes as purposes of documents retrieval. Finally, the software application can be developed for an ease of uses.

ลิขสิทธิ์มหาวิทยาลัยเชียงใหม่ Copyright[©] by Chiang Mai University All rights reserved