

Abstract (บทคัดย่อ)

Project Code : MRG5580030

(รหัสโครงการ)

Project Title : Knowledge Generalization from Causality Knowledge Extracted from

Texts for Answering Why- Question Expressed by Text and Image

(ชื่อโครงการ) การหาความรู้โดยรวมจากความรู้เหตุและผลที่สกัดได้จากเอกสาร

ภาษาไทยสำหรับการตอบคำถามเหตุและผลโดยอาศัยคำถามจาก

ข้อความและภาพ

Investigator : Assoc. Prof. Chaveevan Pechsiri (Dhurakij Pundit University)

(ชื่อนักวิจัย) รองศาสตราจารย์ ดร. ฉวีวรรณ เพ็ชรศิริ (มหาวิทยาลัยธุรกิจบัณฑิตย์)

E-mail Address :itdpu@hotmail.com

Project Period : July 2012 – July 2014

(ระยะเวลาโครงการ) กรกฎาคม 2555 – กรกฎาคม 2557

Abstract

The research aims to extract and generalize the causality knowledge for supporting a Why Question Answering (QA) system integrated with image processing (called the Embedded-Image Why-QA system) for providing the knowledge used in the problem diagnosis, especially in plant diseases. The image expression is applied on the Why-question part for providing Why-question contents (i.e. plant symptoms) that are difficult to be explained by text. There are three main problems involved with this current research. The first problem of the causality knowledge extraction, especially the effect boundary determination problems, is confronted after applying the verb-pair (a causative verb and an effect verb) rules to identify the causality. Then, the research applies Maximum Entropy, Supported Vector Machine, and Naïve Bayes for the comparative study of the effect boundary determination, having the effect verb concepts from the effect EDUs as the features. The second problem is the knowledge generalization problems which come from the extracted causality knowledge containing the uncertainty nuance expression and the incompleteness knowledge. Thus, the research proposes applying the basic linguistic rules to solve the uncertainty problems and the Monte Carlo simulation technique to solve incompleteness problems by imputation of the effect unit. And then, we apply the fuzzy function right after the imputation to determine the generality value of each effect event expressed by the effect verb concept feature of the effect EDU from several documents having the same cause concept. The third problem is from the Embedded-Image Why-QA system which consists of how to determine the Why-question type from the text part of the question, how to determine the Why-question contents from the image part of the question, and how to determine the corresponding answer to the Why-question from the extracted causality. Therefore, the research applies a Why-question cue set to solve the Why-question type, a Bag-of-Visual-Words to solve the Why-question contents, and determining the corresponding answers by ranking the similarity scores between the question content and the extracted causality knowledge including the symptom generality value. Then, the results of this research have shown that the effect boundary determination based on ME has the highest correctness 92% on average and the extracted causality can support the embedded image Why-QA system by answering correctly at 78% correctness at the first rank.

Keywords: Generality value, effect boundary, Embedded-Image Why-QA system, visual word, Why-question cue

บทคัดย่อ

งานวิจัยนี้มีวัตถุประสงค์สกัดและหาความรู้เหตุและผลโดยรวม เพื่อสนับสนุนระบบการตอบคำถามเหตุและผล ที่ได้บูรณาการกับการประมวลผลภาพเข้าไว้ (เรียกระบบคำถามนี้ว่า ระบบเอมเบดเดด-อิมเมจวายคิวเอ Embedded-Image Why-QA system) ทั้งนี้เพื่อการจัดหาความรู้ ให้ชึ้นใช้ในการวินิจฉัยปัญหา โดยเฉพาะปัญหาโรคพีช การนำภาพมาประยุกต์ใช้ในระบบการตอบคำถามเหตุและผลนี้ก็เพื่อย่างต่อการอธิบายลักษณะอาการของโรค ซึ่งบางครั้งยากที่จะ อธิบายด้วยข้อความ ปัญหาสำหรับงานวิจัยนี้ประกอบด้วยสามปัญหาหลักดังนี้ ปัญหาแรกเป็น ปัญหาด้านการสกัดความรู้เหตุและผล โดยเฉพาะเรื่องการหาขอบเขตของผลหลังจากที่ได้ใช้กฎ คู่กริยา (Verb-Pair Rule, กริยาเหตุและกริยาผล) ระบุความรู้เหตุและผล ฉะนั้นงานวิจัยนี้จึงทำการประยุกต์ใช้แม็กซิมัมเอนโทรพีหรือเอ็มอี (Maximum Entropy, ME) :ชั้บพอร์ตเวคเตอร์แมชชีนหรือเอสวีเอ็ม (Support Vector Machine, SVM) และเนย์อีฟเบย์หรือเอนบี (Naïve Bayse, NB) เพื่อศึกษาการเปรียบเทียบการหาขอบเขตของผล โดยมีฟีเจอร์ (Feature) ที่ใช้คือกริยาผล ปัญหาที่สองเป็นปัญหาเกี่ยวกับการหาความรู้โดยรวม ซึ่งเกิดจากความรู้เหตุและผลที่สกัดได้นั้น พบปัญหาเกี่ยวกับ ความไม่แน่นอนของนูแอนซ์ (Nuance) ที่แสดงอาการโรค ปัญหาความไม่ สมบูรณ์ของข้อมูลที่สกัดได้ ดังนั้นงานวิจัยนี้ขอเสนอภูทางภาษาศาสตร์ แก้ปัญหาเกี่ยวกับ ความไม่แน่นอนของนูแอนซ์ และใช้เทคนิคการจำลอง มองติการ์โล (Monte Carlo Simulation Technique) ทำการเติมเต็มข้อมูลเพื่อแก้ปัญหาความไม่สมบูรณ์ของข้อมูล ก่อนที่จะทำการหา ความรู้โดยรวมด้วยฟังก์ชันฟัชชี (Fuzzy Function) สำหรับหาค่าเงินเนอรัลลิตี้ (Generality Value) ของแต่ละกริยาผลที่มาจากสาเหตุเดียวกัน ปัญหาที่สามเป็นปัญหาเกี่ยวกับระบบเอม เบดเดด-อิมเมจวายคิวเอ ที่ประกอบด้วยปัญหาการระบุประเภทคำถามเหตุและผลจากส่วนที่ เป็นข้อความของคำถาม ปัญหาการหาเนื้อหาคำถามจากส่วนที่เป็นภาพของคำถาม และปัญหา การหาคำตอบจากความรู้เหตุและผลที่สกัดได้พร้อมด้วยค่าเงินเนอรัลลิตี้ ดังนั้นงานวิจัยนี้จึง ประยุกต์ใช้คุชตประเภทคำถามเหตุและผล (Why-question cue set) ระบุประเภทคำถามเหตุ และผล ใช้ถุงของวิสചวอลเวอร์ด (Bag of Visual Word) หาเนื้อหาคำถาม และใช้การจัดลำดับ คะแนนความคล้าย ระหว่างเนื้อหาคำถามกับความรู้เหตุและผล ที่สกัดได้มาทำการหาคำตอบ จากการทดลองของงานวิจัยนี้ ได้แสดงให้เห็นว่า การหาขอบเขตของผลด้วยวิธีเอ็มอีให้ความ ถูกต้องเฉลี่ยสูงสุดคือ 92% และความรู้เหตุและผลที่สกัดได้สามารถใช้ตอบคำถามได้ถูกต้อง 78% ที่ลำดับที่ 1 (Rank1)

คำสำคัญ: Generality value, effect boundary, Embedded-Image Why-QA system, visual word, Why-question cue