

試験問題の自動生成を可能とする知識ベース自動構築手法の 提案と評価

古舘昌伸^{†1} 高木正則^{†1} 高木輝彦^{†2 †3}

近年、TOEIC や情報処理技術者試験をはじめとし、様々な検定試験が実施されている。これらの検定試験では、問題の作成に多大な負担がかかっている。この作問における負担を軽減すべく、試験問題を自動生成する研究が行われている。しかし、これらの研究の多くは、問題の生成元となる知識ベースなどを作問者、あるいは開発者が事前に用意しなければならない。そこで、本研究では問題を自動生成するための知識ベース構築のための負担軽減を目的とし、過去問題を用いて知識ベースを自動構築する手法を提案する。具体的には、検定試験の過去問題を既存知識とみなし、過去問題で問われている知識（対象知識）を抽出する。この対象知識を分類基準とし、既存知識の体系化を行う。本稿では、ご当地検定の過去問題を手動で分析・構築した知識ベースを正解データとし、提案手法により自動構築された知識ベースの再現率を求めることで提案手法の評価を行う。

A Proposal and Evaluation on a Method of Automatic Construction of Knowledge Base for Automatic Generation of Exam Questions

MASANOBU FURUDATE^{†1} MASANORI TAKAGI^{†1}
TERUHIKO TAKAGI^{†2 †3}

In recent years, diverse examination such as TOEIC and Information Technology Engineer Examination are carried out. In these examinations, there is a growing problem a large burden creating of the exam questions. To this problem, there is a research that tries decrease burden creating of the exam questions by automatic generating a question to become the contents. However, can't decrease burden because exam preparer have to construct knowledge based for resource automatic generation. Therefore, this research, we aimed to decrease burden on construct knowledge base for automatic generating a questions and we proposal on a method of automatic construction of knowledge base use of existing exam questions. In specific, regard existing exam questions as existing knowledge and an extract from being questioned (subject knowledge) in existing exam questions. And systematize the existing knowledge by this subject knowledge is a classification criterion. This paper, we evaluate a proposal method to find a recall of knowledge base was constructed automatic that exist exam questions of Gotouchi Test analyzed and construct knowledge base is correct data.

1. はじめに

近年、大学入学や就職などにおいて資格の有無が個人の能力を示す1つの指標となっている[1]。資格の需要が高まるにつれ試験の種類も増加していき、日本国内だけでも1000を超える種類の資格・検定試験が存在している[2]。試験規模も情報処理技術者試験[3]や TOEIC[4]などの大規模な試験から、中小規模な検定試験まで多岐にわたる。しかし、多くの試験では作問における手間が大きいことが問題となっており、問題数（コンテンツ）が不足する事態を招いている。特に、小規模な検定試験ではこの問題が顕著に現れており、予算の都合からボランティアに近い状態で作

問をしているものもある。その結果、検定試験の中止を余儀なくされる試験もある。

これらの問題点に対し、作問負担軽減や学習コンテンツ不足の解消を目的とした、問題の自動生成に関する先行研究がいくつか行われている。まず、数学や物理の力学を対象とした問題の自動生成に関する研究が挙げられる[5][6][7]。これらの研究では、問題文や正答などの問題に含まれている情報（以下、問題情報）や試験の対象領域（以下、対象領域）の知識、問題の解法プロセスなどを知識ベース化し、それをリソースとして問題文や選択肢、解説などの自動生成を行っている。知識ベースから様々なパターンの問題を作ることが可能だけでなく、誤選択肢の解説文の生成を可能としたものもある[8]。しかし、これらの研究はある程度が構造化されたものや計算問題のような解法プロセスが定まっている領域の科目にしか対応していない。また、自動生成のリソースとなる知識ベースは、提案されているシステムごとに決められた構造で作問者や開発者が予め用意しておかなければならない。知識ベースの作成は対

^{†1} 岩手県立大学大学院
Iwate Prefectural University Graduate School, Takizawa, Iwate 020-0193,
Japan

^{†2} 電気通信大学大学院情報システム学研究所
Graduate School of Information Systems, The University of
Electro-Communications, Chofu, Tokyo 182-8585, Japan

^{†3} 日本学術振興会特別研究員 PD
Research Fellow of Japan Society for the Promotion of Science, Chiyoda,
Tokyo 102-0083, Japan

象領域に対して十分な知識を持つ者が行う必要があり、作問者に新たな負担がかかってしまう。作問者に知識ベースなどを用意させることなく、既存の辞書（分類語彙表[9]）と問題情報のみを用いて問題の自動生成を行っている研究もある[10]。しかし、対象科目が限定されているほか、十分な知識を持つユーザでないと妥当な問題の生成が行えないという課題が残っている。

一方、数学や物理以外を対象とし、問題の自動生成を行っている研究もある[11][12]。これらは語彙情報や既存の教材を用いて多肢選択形式や穴埋め問題の自動生成を行っている。また、学習者の理解状況を動的に評価し、その時点での理解状況に合わせた問題の自動生成を可能としている研究もある[13][14]。しかし、これらの研究においても自動生成のリソースである知識ベースを作問者が予め用意しておく必要や、教材にタグを埋め込むなどの人手を加える必要がある。また、多肢選択形式問題の誤答選択肢自動生成システムに関する研究も行われている[15]。これは、(1) 一問一答形式問題を解答した際の誤回答情報、(2) 類似問題の選択肢情報、(3) Web上の情報の3つのリソースを用いて知識ベースを自動生成し、作問者に負担をかけることなく誤答選択肢の自動生成を行っている。しかし、問題文や正解選択肢の自動生成は行っていない。

更に、Wikipediaの記事やオンラインニュース記事、企業のホームページなどのWeb上の文章を基に穴埋め問題などを生成するサービスもある[16][17]。しかし、これらはWeb上の知識を基に生成しているほか、既に体系化されている領域である英語の文法が対象となっていたりする。また、生成できる問題形式に制限がある。

以上の既存研究により、問題の自動生成においては多くの有益な知見が示されてきた。しかし、自動生成のリソースとなる知識ベースなどの情報は、基本的に作問者がそれぞれのシステムに対応した知識ベースを用意しなければならない。また、問題の自動生成が可能な対象科目が制限されている。更に、検定試験ではWeb上にない知識を対象とすることが多く、Web上の知識のみを利用した知識ベースでは対応しきれない。

これを踏まえ、本研究では資格・検定試験において、作問者に新たな負担をかけることなく、問題の自動生成を可能とすることを目的とし、問題の自動生成に活用可能な知識ベースの自動構築手法の提案を行う。本手法では、検定試験の過去問題に含まれる知識を体系化し、知識ベースとすることで作問者に新たな手間をかけずに、かつWeb上にない知識にも対応可能とする。これにより、過去問題さえあれば対象領域（試験）を固定せず、幅広い検定試験で知識ベースの自動生成が可能となる。なお、本研究では岩手県盛岡市のご当地検定試験である「盛岡もの知り検定試験（以下、もりけん）[18]」の問題を対象とする。もりけんでは盛岡に関連する人物や事柄についての知識を問う問題

が出題される。試験は1級から3級があり、問題形式は多肢選択形式がほとんどだが、1級では一問一答形式の問題も出題される。

本稿では、まず2章で知識の体系化（知識ベースの構築）の自動化を試みている関連技術や関連研究について述べる。3章では知識ベースの要件定義をし、更に自動構築を行う際の課題を明確にした上で検討する。検討結果を基に4章で自動構築手法を提案する。更に5章では提案手法で構築した知識ベースの有効性を検証する。6章で今後の課題について述べ、7章では構築した知識ベースを用いた試験問題の自動生成についての構想を述べる。最後に8章で全体をまとめる。

2. 知識の体系化に関する先行研究

知識ベースの自動構築を試みている研究は多々存在するが、近年は特に、日本語版 Wikipedia の情報をリソースとして、汎用的な大規模オントロジーの構築を試みる研究が盛んに行われている。その手法は様々なものがあり、既存の日本語オントロジーを用いたものや[19]、日本語版 Wikipedia のみをリソースとして構築をしているものもある[20][21]。また、Wikipedia の Infobox と呼ばれる半構造情報[22]を RDF に変換することによって、大規模なデータベースを構築している DBpedia がある[23]。DBpedia は英語版 Wikipedia をリソースとしているが、日本語版 Wikipedia を対象とし、独自でマッピング作業を行っている DBpedia Japanese も存在する[24]。これらの研究では、これまでより低コストで非常に大規模なオントロジーを（半）自動的に構築することに成功している。また、これらのオントロジーはユーザ参加型という特徴を持つ Wikipedia をリソースとしているため、他にはないほど情報の網羅性が高く、更に即時更新制に優れている。しかし、これらのオントロジーにおける実際の値となるインスタンスやトリプルのプロパティ及びその値は、Infobox からの抽出や、Wikipedia の一覧記事からのスクレイピング、記事のリスト構造からの抽出などによって生成されており、記事の構造化されていない文章は対象としてない。Infobox などの情報は、記事の基礎情報や箇条書きでまとめたものを記載するものであり、検定試験の問題（特に難易度が高いもの）に利用可能な詳細な情報は記載されていない。そのため、Wikipedia の記事に問題となり得る情報があったとしても、オントロジーにはその情報が入っていない。例えば、新渡戸稲造の Wikipedia の記事には、「生誕の地である盛岡市と、客死したビクトリア市は、新渡戸が縁となって現在姉妹都市となっている」という文章がある（2015年1月現在）。もりけんでも、「新渡戸稲造が客死したことが縁で盛岡の姉妹都市になったのはどこですか」という問題が出題されているが、Wikipedia オントロジーではその記述は要素として含まれていない。また、クラス-インスタンスの関係に誤りがある

ものもあり、問題作成においてノイズになってしまう情報もある。以上のことから、構築されたオントロジーをそのまま用いた問題の自動生成は困難であると考えられる。また、Wikipedia に過去問題の情報を基に記事を追加したとしても、問題情報がそのままインスタンスやトリプルとして表現されるわけではないため、問題の生成には使用できない。

一方、オントロジーには特定の領域の知識を体系化したドメインオントロジーもあり、我々が構想している知識ベースも特定の領域（今回は盛岡の領域）に沿った知識ベースである。ドメインオントロジーの構築については、森田らにより提案・開発された DODDLE-OWL がある[25]。これは、Wikipedia やフリーテキストなどの様々な情報資源からドメインオントロジーを半自動的に構築可能なツールである。しかし、入力された概念間の関係はユーザが入力する必要があるなど、依然ユーザの構築コストは高い。また、既存の観光情報サイトのロコミ情報からドメインオントロジーを自動構築し、観光情報推薦を行っている研究もある[26]。しかし、クラス階層は事前に手動で構築しているほか、ロコミ情報はある程度体系化されて投稿されているため、この手法を問題情報に適用することはできない。

更に、シラバスにおいて未分類の教科を自動分類する手法も提案されている[27]。しかし、シラバスの文書集合と検定試験の問題情報では文書の構成や出現する語の性質が異なるため、この手法を問題情報に適用することは難しい。

以上のように、Wikipedia やシラバスなどの情報を（半）自動的に体系化する研究は多数存在している。しかし、本研究が目指すような検定試験等の問題情報を自動的に体系化する研究は、著者らの知る限り存在しない。

3. 知識ベースの要件定義と検討課題

3.1 要件定義

1章, 2章より、本研究における知識ベース構築の要件を以下のようにまとめた。

- 資格・検定試験において作問者に新たな負担をかけることなく知識ベースを構築可能
- 計算問題以外の問題形式に対応可能
- Web 上の知識に依存せず、試験の既存知識（過去問題）を用いて知識ベースを構築可能
- 問題の自動生成（問題文及び選択肢の自動生成）に活用可能

本研究ではこの要件定義に従い知識ベースを構築し、試験問題の体系化を行っていく。

3.2 分類基準の明確化

試験問題の体系化を行う上で、分類基準を明確にし、その基準となる情報を問題情報から抽出する必要がある。もりけんの作問現場において、作問者が試験問題を作成する際には、対象領域（もりけんの場合は盛岡市）の中から分

表 1 知識ベースの構成要素

Table 1 Knowledge base component.

| 構成要素 | 説明 |
|------------------|---|
| 対象知識 | 問題で問われている知識や解決の中心となる知識。ほとんどが専門用語。 |
| カテゴリ (Category) | 問題のカテゴリ。対象知識の上位概念。 |
| プロパティ (Property) | 対象知識や対象知識に関連する事柄についての説明。Object と他要素の繋がりとの関係性。 |
| オブジェクト (Object) | Property が示す実際の値となる要素。 |

野毎に試験の受験者に知っていてほしい、あるいは理解してほしい知識を定め、それを問題として出題している。そのため、問題情報には必ず問われている知識・問題の中心となる知識があり、その知識が問題情報の中で最も重要であると我々は考えた。これはもりけんだけでなく、一般的な資格・検定試験にも言えることであると考えられる。高木らの研究ではこのような知識を「問題で問われている知識や解決の中心となる知識（以下対象知識）」と定義している[28]。本研究でもこの対象知識を分類基準として用いることとした。なお、対象知識はほとんどが専門用語となる。更に、問題の管理がしやすいように対象知識を基に「類似問題」と「関連問題」を定義し、似ている問題や関連のある問題同士のリンク付けを行うこととした。これにより、問題の自動生成を行いやすくすることが狙いである。4.1節で後述するが、対象知識は問題によっては複数になる場合もある。そこで、本研究では類似問題、関連問題を以下のように定義した。

- **類似問題**：対象知識が同じ問題
- **関連問題**：対象知識が問題情報の各要素（3.3 節で後述）に含まれている問題

3.3 知識ベースの構成要素

問題情報には、問題で問われている知識である対象知識と、対象知識について、あるいは対象知識に関連する事柄についての説明が必ず記述されている。これを踏まえ、我々は DBpedia[23]でも用いていた RDF のトリプル形式[29]を参考にし、本研究で構築する知識ベースの構成要素を表 1 のように決定した。表 1 に示すように、知識ベースは対象知識を中心として構成した。まず、分類基準を対象知識としたことから、問題情報のカテゴリは対象知識の用語を基に決定していく。そのため、対象知識はカテゴリの要素の子要素となり、カテゴリは対象知識の上位概念となる。次に、その対象知識について具体的に問われている内容や説

明などをプロパティで示す。問われている内容や説明の答え（つまりプロパティが示す実際の値）を表すのがオブジェクトである。プロパティ・オブジェクトはそれぞれほとんどが問題文・正答の要素となることが考えられるが、問題の出題パターンによってはプロパティに正答の要素が入ったり、オブジェクトに対象知識が入ったりすることもある（4章で後述）。

知識ベースの自動構築を行う際には、一問一問の問題情報から表 1 にある要素を表すものをそれぞれ自動で抽出・決定していく。

4. 知識ベースの各要素の自動抽出方法の検討

本章では、表 1 に示した知識ベースの各要素について、問題情報から自動で抽出・決定する方法について検討する。分析・検討を行う上で、我々はまず手動で問題情報から各要素の抽出を行った。現在もりけんの過去問題は全部で 2000 問存在するが、今回は 2013 年度の 1 級、2 級、及び 2011 年度から 2013 年度の 3 級問題の全 450 問（1 級：50 問、2,3 級：100 問）を対象とし分析を行った。以下、この分析結果を基に検討を行っていく。なお、ご当地検定の場合もそうだが、試験によっては過去問題に解説がないものもあるため、本稿では解説情報を取り扱わないこととした。そのため問題情報において本稿で取り扱うものは、問題文、正答、誤答、問題の出題年度と出題級の情報となる。

4.1 対象知識の自動抽出

4.1.1 高木らの研究における対象知識の自動抽出

本研究では、高木らの研究[28]をベースとし対象知識の抽出を行っているため、ここで高木らの研究における対象知識の自動抽出手法を紹介する。

高木らの研究では、「コンピュータネットワーク論 I」という講義において、学生が作成した問題を対象とし、対象知識の自動抽出を試みている。アプローチとしては、作成された問題にある問題文、正答、誤答、キーワード（解説は除外）の情報から、対象知識の出現箇所を絞り込み特定するという手法を採用している。具体的には、対象知識である専門用語を日本語・英語別、及び単名詞（それ以上分割できない名詞）・複合名詞別に 5 種類の単位の語に分ける。その 5 種類の語と問題の出題パターンを基に、形態素解析等を用いて対象知識の出現箇所の自動決定を行っている。高木らの研究における対象知識の出現箇所自動決定手順を図 1 に示す[28]。出現箇所を問題ごとに特定後、決定された出現箇所に対して形態素解析を行い、先程の 5 種類の単位の語を抽出する。その語が対象知識となる。なお、対象知識は 1 つの問題に複数存在する場合もある。

4.1.2 本研究における対象知識の自動抽出

高木らの研究で扱っている問題と我々が対象としている問題は、対象領域や文書構成、出題パターンが異なるため、

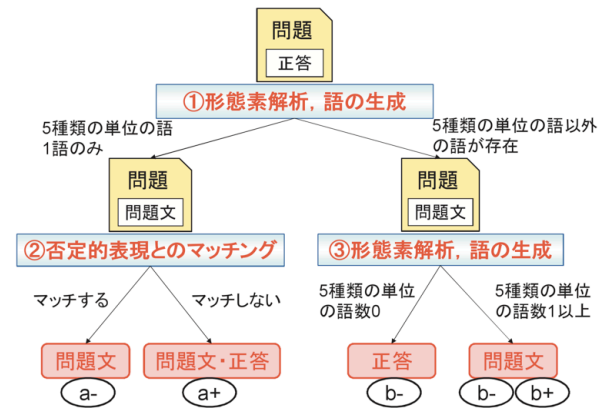


図 1 高木らの研究における対象知識の出現箇所自動決定手順

Fig 1 Procedure of automatically identifying the occurrence part in research of Takagi.

表 2 問題例と対象知識

Table 2 Example quiz and targeted knowledge.

| 設問番号 | 問題文 | 正答 | 対象知識 |
|------|---------------------------------------|------|---------------|
| 1 | 盛岡市にある施設で、平成 25 年に 100 周年を迎えたのはどこですか。 | 盛岡劇場 | 盛岡劇場 |
| 2 | 言語学者の金田一京助が生涯をかけて研究した言語は何語ですか。 | アイヌ語 | 金田一京助 アイヌ語 |

図 1 の手法をそのまま適用出来るとは考えにくい。そこで、我々はまずもりけんの問題について対象知識を手動で抽出し、出題パターンを分類した。

分析した 450 問において、対象知識は 514 個抽出された。ほとんどが人名や建築物名、地域名などの対象領域である盛岡に関連する専門用語となり、それらは高木らの研究と同じように単名詞か複合名詞となった。このことから、問題情報に含まれている専門用語が対象知識となり得ることが分かる。しかし、隣接している市町村の数や駅の数などを答える問題もあり、対象知識が専門用語とならない問題もいくつか存在した。また、対象知識を複数とした問題もあった。表 2 に、対象知識が 1 つとなる問題と 2 つとなる問題の例を示す。表中の設問 1 は盛岡に存在する施設である「盛岡劇場」についての説明から、「盛岡劇場」を答える問題である。この問題の場合は正答である「盛岡劇場」についての知識があれば答えられる問題であるため、設問 1 の対象知識を盛岡劇場とした。表中設問 2 は、「アイヌ語」を答える問題だが、「アイヌ語」は対象領域独自の専門用語ではないほか、この問題を解くには「アイヌ語」の知識だ

表 3 問題の出題パターンと各級ごとの割合 (%)

Table 3 Quiz types and percentage of each class.

| パターン ID | 出題パターン | 1 級 | 2 級 | 2013 年 3 級 | 2012 年 3 級 | 2011 年 3 級 | 計 |
|---------|---|-----|-----|---------------|---------------|---------------|------|
| Pa+ | ある専門用語について正しい例や説明などを選択する問題 例：盛岡弁で「せっこき」とはどのような意味ですか。 | 2 | 7 | 4 | 2 | 5 | 4.2 |
| Pa- | ある専門用語について誤った例や説明などを選択する問題 例：「国分通り」の名称の由来となった「国分謙吉」について、正しくないのはどれですか。 | 0 | 0 | 1 | 0 | 0 | 0.2 |
| Pb- | ある専門用語についてその種類・属性と異なる専門用語を選択する問題 例：次のうち、実在しない学校はどれですか。 | 0 | 0 | 0 | 0 | 1 | 0.2 |
| Pc+ | ある専門用語に関連する事柄について正しい例や説明などを選択する問題 例：志波城が造られて約 10 年後に徳丹城へ移転した理由は何ですか。 | 4 | 7 | 1 | 1 | 0 | 2.4 |
| Pd+ | ある専門用語に関連する事柄について関係する専門用語を選択する問題 例：旧盛岡藩主・南部家の家紋（表紋）に描かれた鳥の種類は何ですか。 | 14 | 11 | 24 | 20 | 21 | 18.4 |
| Pd- | ある専門用語に関連する事柄について関係しない専門用語を選択する問題 例：次のうち、実業家「三田義正」の残した事業でないものはどれですか。 | 4 | 1 | 2 | 1 | 0 | 1.3 |
| Pf | ある専門用語についての例や説明からその専門用語を選択する問題 例：中津川の中の橋と下の橋の間に架かる橋は何ですか。 | 70 | 63 | 55 | 62 | 63 | 61.8 |
| Pg | ある専門用語について正しい値（数値や日付）を選択する問題 例：盛岡市に隣接している市町村の数はいくつですか。 | 2 | 7 | 9 | 6 | 5 | 6.2 |
| Ph | ある専門用語について正しい読み方を選択する問題 例：江戸時代、盛岡城下の小人町は「御得道具丁」とも呼ばれました。この読み方を平仮名で書きなさい。 | 2 | 0 | 0 | 1 | 0 | 0.4 |
| Pi+ | 各選択肢に記述された専門用語同士の組み合わせのうち正しい組み合わせを選択する問題 例：次の伝説と伝えられている寺名が正しいものはどれですか。 | 0 | 1 | 0 | 0 | 0 | 0.2 |
| | その他（穴埋め、画像、地図、並び替え） | 2 | 3 | 4 | 7 | 5 | 4.4 |

けではなく、「金田一京助」の知識も必要となる。そのため、対象知識を「金田一京助」、「アイヌ語」とした。このように、対象知識となり得る専門用語が対象領域独自の専門用語ではない、かつその用語の知識だけでは解けない問題の場合は、対象知識を複数とした。この分析結果から、もりけんにおける対象知識となり得る専門用語を以下のように定義した。

- (1) 対象領域である盛岡に存在する施設名や地域名など、盛岡独自の専門用語
- (2) 上記以外の専門用語

(3) 人名や氏族の名前

(2)の専門用語については、(1)の用語と組み合わせて複数の対象知識とする。また、もりけんの問題には人物名を答える問題が多数あり、そのような問題の場合は人物名を対象知識とすることとした。

出題パターンについては高木らの研究を参考に、対象知識の問われ方によって分類した。高木らの研究と比較すると、もりけんには出題されなかったパターンもあるが、どのパターンにも当てはまらない問題もあったため、新しくパターンを作ったものもある。もりけんに出題されたパタ

表 4 出題パターンごとの対象知識出現箇所の割合

Table 4 Percentage of targeted knowledge occurs in each quiz content.

| 3 級/パターン ID | Pa+ | Pa- | Pb- | Pc+ | Pd+ | Pd- | Pf | Pg | Ph | Pi+ |
|-------------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 問題文 | 100 | 100 | 100 | 100 | 50 | 100 | 0 | 100 | 100 | 0 |
| 正答 | 0 | 0 | 0 | 0 | 50 | 0 | 100 | 0 | 0 | 100 |
| 誤答 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

ーンとその例, 各級ごとの割合を表 3 に示す. 今回分析した問題に出題されていなかったパターンについては記載していない. 「その他」には穴埋め問題や画像・地図を用いた問題, 並び替え問題などが含まれる. これらの問題は, 対象知識をはじめ知識ベースの各要素を自動で抽出・決定することが困難なため, 本研究では対象外とした. もりけんにおいては, どの級でも出題パターンの傾向が似ており, パターン Pf, Pd+の問題が多く, 分析した全問題の約 80% がこのどちらかのパターンとなっていた. よって, まずはこの Pf, Pd+のパターンを持つ問題に対応することが重要であるといえる. 更に, 出題パターンごとに対象知識の出現箇所を分析した (表 4). 表中に示している割合は, 各出題パターン, 要素に出現した対象知識の数を, 全問題から抽出された対象知識の数 (514) で割ることで求めている. 表 4 から, 対象知識は問題文に出現している問題が多いが, Pf, Pi+では正答, Pd+では問題文と正答に出現していることが分かる. また, 対象知識が複数となる問題は出題パターン Pd+の問題のみであることが分かる.

以上の結果を踏まえ, 本研究における対象知識の出現箇所の特定及び出題パターンの特定手法を図 2 に示す. まず, 問題文に対し「～の読み方を」などの読み方を答える問題表現とマッチするかを調査する (図 2 中①). マッチした場合は出題パターンを Ph とみなし, 対象知識の出現箇所を問題文とする. マッチしない場合は正答に対し形態素解析を行い, 単名詞か複合名詞のみが存在するかを調査する (図 2 中②). 名詞以外 (助詞など) が出現した場合は, 誤答にも形態素解析を行い, 正答, 誤答の全てに数値が存在するかを調査する (図 2 中③). 存在する場合は出題パターンを Pg, 存在しない場合は出題パターンを Pa+, Pa-, Pc+とみなし, 対象知識の出現箇所を問題文とする. 一方, 正答が名詞のみの場合は, 問題文に対し「正しくないのはどれか」などの否定的表現とのマッチング調査を行う (図 2 中④). マッチした場合は出題パターンを Pb-, Pd-とし, 対象知識の出現箇所を問題文とする. 否定的表現とマッチしなかった場合は, 正答, 誤答に対し形態素解析を行い, 名詞がそれぞれに 2 つずつ存在しているかを調査する. 存在する場合は出題パターンを Pi+とし, 対象知識の出現箇所を正答とする. 存在しない場合は正答が人名かを調査し, 人名の場合は出題パターンを Pf とし, 対象知識の出現箇所を正答とする. 正答が人名でない場合は, 正答が対象領域

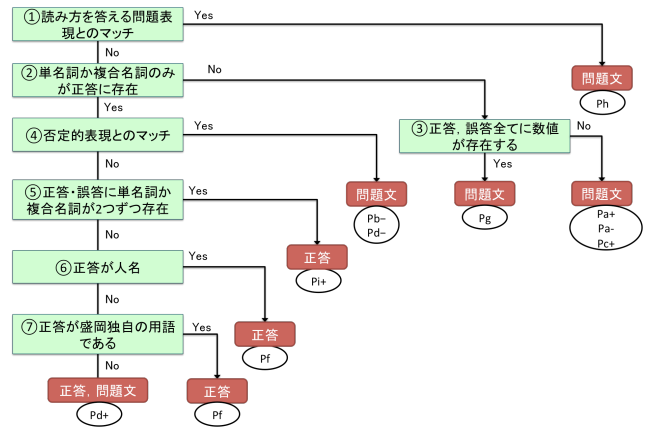


図 2 対象知識の出現箇所及び出題パターンの特定手順
 Fig 2 Procedure of automatically identifying the occurrence part and quiz types.

(盛岡) 独自の用語であるかを調査し, 独自の用語である場合は出題パターンを Pf とし, 正答を対象知識とする. 独自の用語でない場合は出題パターンを Pd+とし, 対象知識の出現箇所を正答と問題文とする. なお, 盛岡独自の用語としては, 用語に「盛岡」や「岩手」の単語があるもの, 教育施設や寺院, 橋などの建築物の名前, 町の名前, 山の名前とした. しかし, 盛岡市にある「石割桜」など, この条件に当てはまらないものもいくつかある. よってこの条件の他に, 対象の専門用語の Wikipedia 記事がある場合は, その記事を参照し, 記事内に「盛岡市にある」という単語がある場合は盛岡独自の用語とすることとした.

4.2 カテゴリ (Category) の自動決定

3.3 節で述べたように, カテゴリは対象知識の上位概念となる. そのため, ここでは対象知識である用語の上位下位関係を自動決定する方法について検討する. なお, 問題を分類する手法として, 問題情報にキーワードやタグを付与し, その情報を基に分類する手法も考えられる. しかし, これらの情報に基づいた場合, 分類はユーザの主観となってしまう, 作問者が何人もいた場合では統一性がとれなくなる可能性がある. また, キーワード等が付与されていない既存の問題に対し, これらの情報を付与していくことは非常に手間がかかる.

ある日本語の単語 (名詞) について意味別に分類を行うために, 日本語語彙大系 [30] や分類語彙表 [9], 日本語 WordNet [31], EDR 電子化辞書 [32] を用いることが多い. こ

れらは単語の意味（概念）を階層的に分類，体系化した日本語シソーラスであり，全て電子化されている．これらは莫大な数の日本語単語の上位・下位関係を収録しているが，専門用語や固有名詞の網羅率は低い．また，これらは人手により構築されているため，新たな語を追加する際にもその意味をよく検討する必要がある，容易に追加することはできない．ご当地検定試験を含め，多くの資格・検定試験においては専門用語や新語を用いることが多いため，これらのシソーラスをそのまま用いるのは難しいと考えられる．そこで，2.3節でも紹介した日本語 Wikipedia オントロジー（JWO: Japanese Wikipedia Ontology）を用いることとする．JWOは専門用語や固有名詞などの語彙網羅性，及び即時更新性に優れているため[33]，他と比べ時事用語などの新単語にも対応可能である．しかし，JWOには以下の2点の問題点がある．

- ① 上位概念，中間概念が不足している
- ② 全ての用語に対応していない

①については，JWOは日本語 WordNet などと比べると単語の上位概念，中間概念が不足している．例えば，JWOにおいて「人物」の上位概念は「事物」であるが，日本語 WordNet においては「生物」，「生き物」などがある．また，「人物」から「ギタリスト」の間にある中間概念においては，JWOは「音楽家」，「演奏家」の2つとなっているが，日本語 WordNet では「エンターテイナー」，「公演者」，「ミュージシャン」の3つとなっている[20]．

②については，1章や2章で述べたように，検定試験によっては Web 上の知識にないものを取り扱っているものもあり，また JWO にインスタンスとして登録されない知識がある．そのため，対象知識として抽出した用語全てに対応可能なわけではない．

上記を踏まえ，我々はカテゴリの決定方法として以下の3つの方法を提案する．

JWO と既存の日本語シソーラスの組み合わせ
問題文からカテゴリを表す用語の抽出
専門用語の一般名詞化

(1)の手法については，対象知識の用語が JWO にインスタンスとしてある場合，この手法を用いる．本稿では既存の日本語シソーラスとして，唯一オープン化されている日本語 WordNet を用いることとした．まず JWO で対象知識の用語の上位概念であるクラスを取得する．不足している上位概念については，取得したクラス情報を日本語 WordNet にある用語と結びつけることで補う．例えば，JWO で「盛岡劇場」を参照する．JWO では「盛岡劇場」のクラスは「劇場」となっている．その「劇場」を今度は日本語 WordNet で参照する．そうすることで，日本語 WordNet に登録されている「劇場」の上位語である「建物」，「建築物」などが取得できる．(2)，(3)は JWO に対象知識の用語が存在しない（つまり Wikipedia の記事がない）場合に用いる．

もりけんの問題においては，例えば橋の名前を答える問題の場合は，「～の橋は何ですか」や「～の橋の名前は何か」という問い方をしている．この文章構成を利用し，問題文の語尾の直前に出現する名詞を対象知識のカテゴリとする((2)の手法)．問題文の語尾の直前に名詞がない場合は，(3)の手法として対象知識の用語に対し形態素解析を行い，一般名詞の抽出を試みる．例えば「三ツ石神社」の場合は，後方文字列処理により「神社」という名詞を抽出し，これをカテゴリとする．

また，いくつか例外処理を行った問題もある．方言（盛岡弁）の意味を問う問題や，方言の意味から方言を答える問題については，カテゴリを「方言」とするようにした．この判断は，問題文に「盛岡弁」，「方言」の単語が出現しているか否かとした．また，対象知識が「盛岡市」としかなり得ない問題がある．例えば，「平成 25 年現在，盛岡市の人口はおおよそ何人ですか」という問題や，「盛岡市に隣接している市町村の数はいくつですか」というような問題である．このような問題は，「基本情報」というカテゴリとすることとした．基本情報となる問題については，対象知識が「盛岡市」になる問題だけでなく，盛岡市の Wikipedia 記事の Infobox を参照し，問題文に「人口」や「市のシンボル」など特定の単語が含まれた場合，その問題のカテゴリを基本情報とすることとした．

4.3 プロパティ（Property）の自動抽出

プロパティの自動抽出方法の検討にあたり，我々はまず対象知識と同じように問題の出題パターンごとにプロパティとなる情報が問題情報のどこに出現しているのかを分析した．その結果，Pa-以外の全てのパターンにおいて，プロパティは問題文に出現しており，Pa-は誤答に出現していた．よって，図 2 の手法を利用して出題パターンを特定し，Pa-以外のパターンの場合は問題文から，Pa-パターンの場合は誤答から抽出することとした．ただし，もりけんの場合，問題文には複数の情報が含まれている場合がある．例えば，「平成 24 年暮れに盛岡駅前商店街に登場した，盛岡三大麺の一つをアレンジした新ご当地料理は何ですか」という問題文には，「平成 24 年暮れに盛岡駅前商店街に登場した」と，「盛岡三大麺の一つをアレンジした新ご当地料理」の2つの情報がある．このように複数の情報を記述することで，解答者にヒントを与えることができ，または解答を絞り込むことが可能となる．このように複数の情報がある場合は，分割して抽出することとした．よって，上記の問題の場合はプロパティが2つとなる．なお，分析の結果，問題文等に複数の情報が含まれている場合は，ほとんどが句読点で区切られていることが分かったので，本研究ではプロパティにおいては，1つの情報を句読点で区切ることにした．また，プロパティの抽出の際，不要語の削除と置き換えを行っている．不要語としては，問題文の語尾である「何ですか」や「誰ですか」などが該当する．また，もりけんの

問題には時事問題が多く出題されるため、「今年」や「昨年」などの単語が出現する。これらの単語は問題の出題年度の情報を用いて「2012年」などに置き換える。

4.4 オブジェクト (Object) の自動抽出

オブジェクトについても、同じように出題パターンごとのオブジェクトの出現箇所を分析した。オブジェクトは、説明を選択する出題パターンである Pa-, Pd-以外では、正答に出現しており、Pa-では問題文、Pd-では誤答に出現していた。ただし、Pa-の場合は対象知識と同じ用語がオブジェクトとなった。また、Pd-では誤答である 3 つの用語がオブジェクトとなり得た。よって、Pa-では対象知識を、Pd-は誤答の 3 つの用語をオブジェクトとして扱う。それ以外の出題パターンの場合は正答をオブジェクトとした。

5. 知識ベース自動構築手法の提案

本研究で提案する知識ベース自動構築手法の全体構成を図 4 に示す。4 章の検討結果を基に、本提案手法では 2 つのモジュールを作成した。「問題変換モジュール」では、問題情報から表 1 の要素の抽出を行う。「クラス作成モジュール」では、4.2 節の手法で問題情報のカテゴリの自動決定を行う。また、本手法により構築された知識ベースを用いて、「作問モジュール」で問題の自動生成を行う。作問モジュールについては、現段階で構想している内容を 8 章で述べる。なお、本手法では形態素解析を用いているが、形態素解析には MeCab[34]を利用している。

6. 実験・評価

構築した知識ベースについての評価は、手動で構築した知識ベースを正解データとし、比較結果によって有用性を検証していく方法を取る。まず、知識ベースのそれぞれの要素について、問題情報から自動抽出されたものが合っているかどうかについて評価し、更に問題のカテゴリについて比較していく。

7. 今後の課題

今後は本研究で提案した手法を、他の資格・検定試験においても適用し、検証を行っていく必要がある。また、本研究により自動構築された知識ベースを用いて、実際に試験問題を自動生成する必要がある。その手法については構想中ではあるが、次章に現段階で検討している内容を述べる。

8. 知識ベースを用いた試験問題自動生成の構想

本章では本研究により構築した知識ベースを用いて、資格・検定試験問題の自動生成を行う手法について述べる。

作問モジュールでは、出題テンプレートを用いて問題の自動生成を行うことを想定している。本モジュールの全体構成を図 4 に示す。出題テンプレートは表 3 で示した出題パターンに従い、数通り用意する。テンプレートは穴埋め

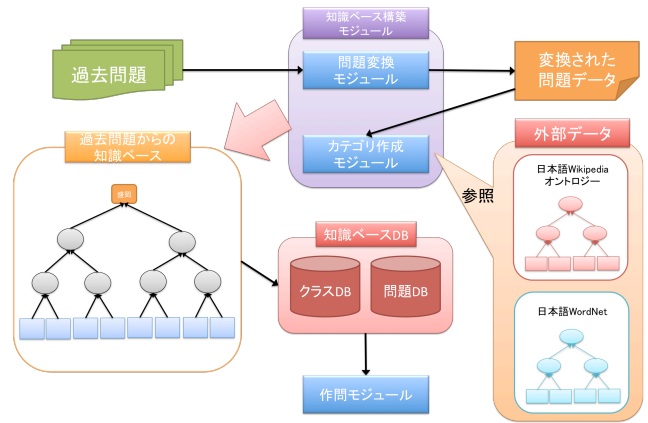


図 3 知識ベース自動構築手法の全体概要

Fig 3 Proposal model of Method of Automatic Construction of Knowledge Base.

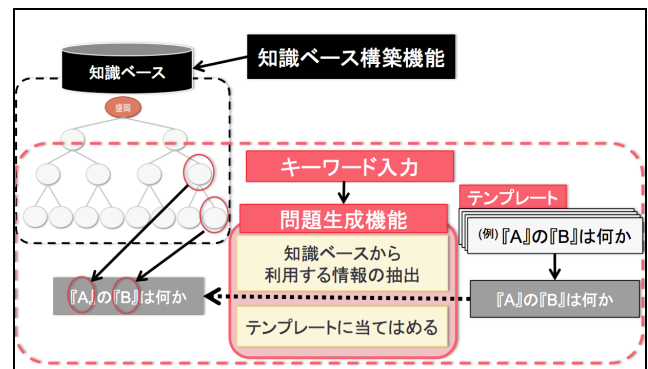


図 4 作問モジュール全体構成図

方式となっており、表 1 の知識ベースの各要素を適切な場所に入れることで、問題の生成を行う。本モジュールでは新問題の作成、及び類似問題の作成を行う。そのため、テンプレートへの要素の格納については、生成する問題によって適切な位置が変わる。新問題の自動生成例を図 5 に示す。現段階では、新問題を「対象知識が変化する場合」、あるいは「対象知識は同じだが、問う内容が変化した問題」と定義している。類似問題については、3.2 節で「対象知識が同じ問題」と定義している。そのため、類似問題の自動生成については、対象知識が同じ問題を生成することを前提としている。図 6 に類似問題の自動生成例を示す。この例では、知識ベースのプロパティの要素を組み合わせることにより、類似問題を生成している。このような問題の自動生成を行うことで、作問者の更なる負担軽減を狙う。本モジュールの想定している利用例を図 7 に示す。この例では、作問者が「金田一京助」についての問題を作成したいと考えている。その際に、作問したい問題の対象知識である「金田一京助」を入力することで、「金田一京助」に関する新問題、類似問題が本モジュールにより複数生成される。作問者は生成された問題に対し、編集等で手を加える

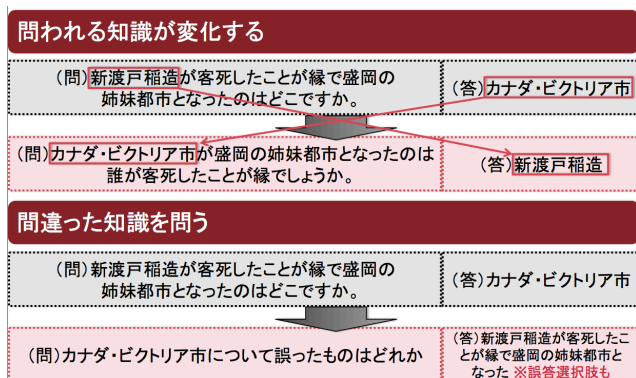


図 5 新問題の自動生成例



図 6 類似問題の自動生成例

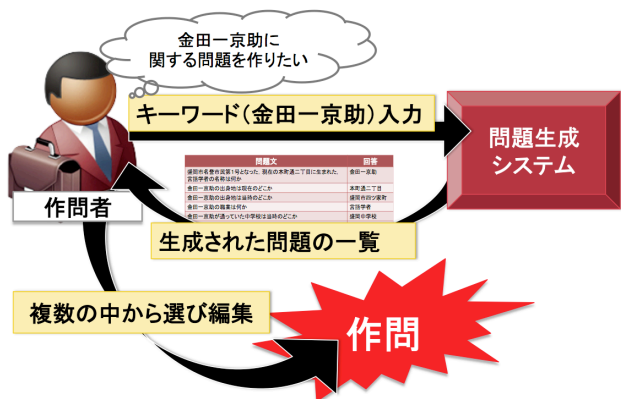


図 7 想定している作問モジュールの利用例

ことで、より良い問題とする。以上が、本モジュールにおいて現在検討している機能と、想定している利用例である。

9. まとめ

本研究では、作問者に新たな負担をかけることなく、問題の自動生成を可能とすることを目的とし、問題の自動生成に活用可能な知識ベースの自動構築手法の提案を行った。知識ベースの構築においては、問題情報から対象知識を中心として、定義した知識ベースの要素ごとに抽出し、カテゴリ分類を行った。提案手法に対しては、手動で構築した知識ベースとの比較により有用性を検証した。今後は、他試験への適用や、構築した知識ベースを用いた問題自動生

成方法の確立などを行っていく必要がある。

謝辞 本研究の一部は科学研究費(若手研究(B), No.24700904)の研究助成を受けたものである。また、本研究に協力して頂いた盛岡商工会議所と文化地層研究会の皆様方に感謝を申し上げます。

参考文献

- 1) 資格世界の評論-資格試験の歴史【平成の資格王・中村一樹】、<http://www.shikakuou.com/02/02-2.html>
- 2) 日本の資格を全部掲載！1000以上の資格を解説 | 資格キング、<https://www.shikaku-king.com/>
- 3) 情報処理技術者試験、<http://www.jitec.ipa.go.jp/>
- 4) TOEIC | コミュニケーション英語能力を測る世界共通のテスト、<http://www.toeic.or.jp/>
- 5) 松岡秀朗, 金西計英, 光原弘幸, 松浦健二, 矢野米雄: 演習問題作成システムのための領域知識作成支援について, 電子情報通信学会技術研究報告, Vol.103, No.697, pp.179-184 (2004).
- 6) 高野敦子, 橋本淳: 知識ベースに基づいた学習者個別演習問題生成手法について, 情報処理学会研究報告, NL-160, pp.23-28 (2003).
- 7) 大川内祐介, 上野拓也, 平嶋宗: 派生問題の自動生成機能の開発とその実験的評価, 人工知能学会論文誌, No.27, Vol.6A, pp.391-400 (2012).
- 8) 舟生日出男, 穉山雅史, 平嶋宗: 問題解決プロセスを利用した選択問題の誤選択肢および解説の自動生成, 電子情報通信学会論文誌 D, Vol.93-D, No.3, pp.292-302 (2010).
- 9) 分類語彙表, <http://www.kokken.go.jp/>
- 10) 小島一晃, 三輪和久: 作問事例を用いて数学文章問題を生成するシステムの実現と評価, 人工知能学会誌, Vol.21, No.4, pp.361-370 (2006).
- 11) 北岡大輔, 松田憲幸, 平島宗, 滝寛和: 補助教材のための選択問題と誤答解説文の自動生成の構想, 電子情報通信学会技術研究報告, Vol.103, No.320, pp.55-58 (2003).
- 12) 宮地功: e ラーニングにおける客観テスト問題の自動生成方法の提案, 電子情報通信学会技術研究報告, Vol.108, No.247, pp.1-4 (2008).
- 13) 菅沼明, 峯恒憲, 正代隆義: 学生の理解度と問題の難易度を動的に評価する練習問題自動生成システム, 情報処理学会論文誌, Vol.46, No.7, pp.1810-1818 (2005).
- 14) 津森伸一, 海尻賢二: 理解状況に適応した選択問題生成方法の検討, 教育システム情報学会誌, Vol.26, No.3, pp.240-251 (2009).
- 15) 菅原達彦, 高木正則: 誤回答情報を用いた多肢選択形式作問支援システムにおける誤答選択肢の評価, 情報処理学会研究報告, Vol.2014-CE-123, No.13, pp.1-9 (2014).
- 16) 穴埋め択一問題の自動生成(日本語、英語) | 学びing 株式会社, <http://manabing.jp/services/autogenerate>
- 17) Ayako Hoshino, Hiroshi Nakagawa: Assisting cloze test making with a web application, SITE 2007--Society for Information Technology & Teacher Education International Conference, pp.2807-2814 (2007).
- 18) 地元学検定『もりけん』, <http://www.ccimorioka.or.jp/jinzai/moriken.html>
- 19) 柴木優美, 永田昌明, 山本和英: 日本語語彙大系を用いた Wikipedia からの汎用オントロジー構築, 情報処理学会研究報告, NL194-4 (2009).
- 20) 玉川奨, 桜井慎弥, 手島拓也, 森田武史, 和泉憲明, 山口高平: 日本語 Wikipedia からの大規模オントロジー学習, 人工知能学会論文誌, Vol.25, No.5, pp.623-636 (2010).
- 21) 玉川奨, 森田武史, 山口高平: 日本語 Wikipedia からプロバティを備えたオントロジーの構築, 人工知能学会論文誌, Vol.26,

No.4, pp.504-517 (2011).

22) Help:Infobox – Wikipedia,

<http://ja.wikipedia.org/wiki/Help:Infobox>

23) Soren Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, Zachary Ives: DBpedia: A Nucleus for a Web of Open Data, 6th International Semantic Web Conference, Vol.4825, pp.722-735 (2007).

24) DBpedia Japanese, <http://ja.dbpedia.org/>

25) 森田武史, 手島拓也, 和泉憲明, 山口高平: 既存情報資源を活用したオントロジー構築支援環境, 電子情報通信学会技術研究報告, Vol.106, No.617, pp.21-24 (2007).

26) 保科宗淳, 大河原渉, 山口高平: 領域オントロジーと Linked Data を利用した観光情報推薦, 第 26 回人工知能学会全国大会 (2012).

27) 太田晋, 美馬秀樹: 課題志向別シラバス自動分類システムの設計と実装: 言語処理学会論文誌, Vol.16, No.4, pp.91-106 (2009).

28) 高木輝彦, 高木正則, 勅使河原可海: 学生が作成した問題の類似度算出手法の提案と評価, 情報処理学会論文誌, Vol.50, No.10, pp.2426-2439 (2009).

29) RDF Primer,

<http://www.w3.org/TR/2004/REC-rdf-primer-20040210/>

30) 池原悟, 宮崎正弘, 白井諭, 横尾昭男, 中岩浩巳, 小倉健太郎, 大山芳史, 林良彦: 日本語語彙大系 CD-ROM 版, 岩波書店 (1999).

31) Francis Bond, Hitoshi Isahara, Sanae Fujita, Kiyotaka Uchimoto, Takayuki Kuribayashi and Kyoko Kanzaki: Enhancing the Japanese WordNet, 7th Workshop on Asian Language Resources, pp. 1-8 (2009).

32) 日本電子化辞書研究所: 電子化辞書使用説明書 第 2 版 (2001).

33) 玉川奨, 香川宏介, 森田武史, 山口高平: 日本語 Wikipedia オントロジーの構築と利用, 人工知能学会 第 29 回セマンティックウェブとオントロジー研究会 (2013).

34) MeCab: Yet Another Part-of-Speech and Morphological Analyzer, <http://mecab.googlecode.com/svn/trunk/mecab/doc/index.html>