

■原 著■ 2017 年度神奈川大学総合理学研究所共同研究助成論文

記述式設問に対する自動採点モデルの試案

永松礼夫^{1,2} 後藤智範¹

Proposal of Model for Automated Marking Techniques for Description Type of Examination

Leo Nagamatsu^{1,2} and Tomonori Gotoh¹

¹ Department of Information Science, Faculty of Science, Kanagawa University, Hiratsuka City, Kanagawa 259-1293, Japan.

² To whom correspondence should be addressed. E-mail: lnag@kanagawa-u.ac.jp

Abstract: After AY2020, the new "national center test for university admissions" will start. Description-type questions that are not included in the current national center test are newly introduced. The marking process for description type problems -- evaluating scores for each answer script -- is modeled as the semantic similarity of answer sentences to the correct answer sentences. It can be handled within a framework of natural language analysis. In this paper, based on natural language analysis, we proposed a basic model and a prototype method for automated marking. We also discussed some problems regarding the actual implementation of this method.

Keywords: automatic exam marking, description type of examination, Japanese morpheme analysis, thesauri, knowledge base, dependency analysis, domain knowledge

序論

現在の形式での大学入試センター試験は 2019 年度 (2020 年 1 月) の実施を最後に廃止され、2020 年度からは新形式の「大学入学共通テスト」の実施に移行され、これに伴って記述式問題が導入される予定であり¹⁾、2017 年秋には試行調査が行われた²⁾。一方、神奈川県では県立高校の入学試験で、2016 年度から記述式問題が導入されている³⁾。

記述式問題の採点は、マークシート方式とは異なり機械的採点は困難と認識され、上記の試験では、人手による採点を実施されている。人手による採点は、マークシート方式・コンピュータによる採点と比較し以下に挙げる問題がある。

(1) 採点時間 (2) コスト

(1) については、センター試験、予備校の試験などは、特定期間内に受験者の得点を公表する必要があり、受験者数によっては多数の採点者を必要とされる。また、採点者の負担が大きくなる問題も発生しうる。結果的に人件費・コストが増加する(2)。したがって記述式設問に対してもコンピュータによる採点が期待される。

本研究は、正解文の長が 40 ~ 100 字に限定した

記述式問題について自動採点を意図したモデルおよび手法を提案し、その問題点を検討する。

設問の種類とその特徴

設問形式の種類

試験問題において設問・問いの形式としては、以下の代表的な 3 種類が挙げられる。

(1) 選択式、(2) 記入式 (穴埋め式)、(3) 記述式

(1) 選択式

問題文に続いて与えられた選択肢の中から解答を選ぶ方式である。大学入試センター試験の大半の設問では選択肢の数とした五つ程度をこの方式が採用されている。受験者の解答パターンは出題者の想定したものに限られ、パターン数は有限個なので採点の自動化は容易である。

(2) 記入式 (穴埋め式)

文章の一部が空欄になった設問について、それを埋める語 (または値) を答える方式である。正解とされた単語の一致を見るだけなので (表記のゆれを除けば) 容易に自動採点できる。類似形式として複数の単語からなる句や数式を答えさせるものもある。

(3) 記述式

文あるいは文章を答えさせるものである。想定する文字数が少ないもの（文の数が1または2程度で、100文字程度までのもの）と長いもの（文の数が多く、想定文字数も多い）がある。採点を容易にするため文字数制限を課すことが多い。また、設問で何を問うかの方向性として以下の3項目が挙げられる。

(a) 知識、(b) 思考力、(c) 表現力

(a) 知識を問う：複数の単語（教科書にあるようなキーワード）の関係性を記述した文を書かせて、事実関係の連関を理解しているかを問うような使い方。例えば「黒船の来航を契機に幕府は開港を決めた」のような文を想定する。採点の自動化の立場からは必要なキーワードが含まれていて適切に関係づけられているかを判定すれば可能と思われる。

(b) 思考力を問う：受験者に推論の過程を示したような解答を要求するもの。採点のためには複数の文の接続として論理が整合して述べられているかの判定が必要となる。

(c) 表現力の評価：文科省／大学入試センターによる「表現力」の位置づけは、「知識・技能」と対比する「思考力・判断力・表現力」の一部とされ、単体での定義は明確ではない。「思考力・判断力・表現力を問う条件付記述式問題について」⁴⁾の「表現力」の項目では「目的に応じた文章の構成や展開を工夫し、論拠に基づいて自分の考えを文章にまとめる」とされている。試験問題の立場からは、解答文を作りあげる力であり国語科の採点項目であるが、他教科では「内容が適切に表現されている」ことを採点で考慮するのみでよいと考えられる。

記述式設問において、字数制限だけでは正解と同じ、もしくは類似の意味の文が多数発生し、正解／不正解の判別を困難にする。これを軽減する手段として、解答文に含まれるべき単語や句を指定、または設問文中に記載し、それらの並べ替えで文章を作成させる設問形態もある。「情報を SNS に発信」という例では、制限がなければ、「発信」の言い換えとして、「投稿」、「送信」、「アップロード」などとした答案があり得るが、使える単語群に「投稿」のみを設定すれば解答例を少なくできる。

設問形式の特徴

設問形式の種類で挙げた3種類の設問形式毎の比較表を表1に示し、比較項目について以下に説明する。解答パターン数：正答・誤答を含めて提出される答案が何通りあるかを示す。

正解の数：正答とされる答案が何通りあるかを示す。

正解判定：正解か／否かを判定する際に行う処理内

容である。

採点の公平性：人間の複数の採点者の間で評点に差が出る可能性があるかを示す。

自動化：機械的採点が可能かを示す。

部分点：採点結果が正解・不正解の二値でなく、数段階の部分点に分けた点数付与にできるかを示す。

表 1. 設問形式の特徴

	選択	記入	記述
解答パターン数	1	少数	多
正解の数	1	1+ α	多い*
正解／不正解の判定	一意	一意+言換え	複数条件で判定
採点の公平性の担保	可	可	ゆらぎ大
自動化	可	可	難
部分店の可能性	no	no	あり

* 正解と意味的に等価な言換えが多数発生する。

記述式設問では、人間（採点者）の関与が必要となる。例えば、試行調査の国語・記述式の「正答の条件」⁵⁾では、『○○○○』ということが書かれている」といった項目が複数設定され、それらを満たしているかで判定する方式である。客観的にこの判断を行うのは困難なので、公正性を保つには一つの答案に複数の採点者がつくことが望ましい。

方法
類似度に基づく評価

設問形式の特徴で述べたように、一般的には記述式設問に対する解答は、1つ以上の文で表現されるが、問題を単純化するために正解文および解答文について以下の2つの条件を満たす記述式設問を想定する。

(a) 文字数：40～100字、(b) 文数：1

文は構成単語、およびこれらの構文構造で規定される特定の意味を有する。一方、特定の意味を表現する文は唯一ではない。ある特定の文を構成する個々の単語について、同義語が存在しない文を想定することは現実的でないからである。このため他の形式の設問と異なり、記述式設問において、解答の評価を正解と完全に一致するか否かとするのは不適切である。記述式設問の解答の評価は以下のように定義することが妥当と考えられる。

$$\text{正解文に対する解答文の類似度} = [0, 1] \cdot \cdot (1)$$

解答文が正解文と意味的に完全に一致していれば1、そうでなければ類似の程度に応じて0までの値となる。左記の定義に従えば、ある記述式設問の配点を p とし、解答文の実際の評点を p' とすると、(1)式からは次式のように定義することができる。

$$0 \leq p' \leq p \cdot \cdot \cdot \cdot \cdot (2)$$

以下では、個々の得点 p' を何に基づいて、どのような方法で決定するかについて述べる。

採点モデル

考え方

一般に、記述式設問の正解文には当該問題において「重要な概念を表す語・句」を複数含んでいる。この「重要な概念を表す語・句」は、問題作成者が解答文に含まれてほしいと期待する語・句であり、ここでは、「評価対象表現」とよぶことにする。一般に、記述式設問では正解文中に含まれる、言い換えれば、正解文を構成する評価対象表現の数は、指定された字数の長さに応じて増加すると言える。例えば、指定された文字数が20字前後の場合は、正解文中の語数が数語となるので、評価対象表現は単語となる場合が多く、また50字以上では、動詞を含む句となる場合もありうる。

記述式設問の構成要素を以下のように定義する。

Sr : 正解文、 Sa : 解答文

e : 評価対象表現、 e' : 対応表現、 w : 語

k : 正解文中の評価対象表現の数

h : 解答文中の対応表現の数

正解文は w, e の並び、解答文は w, e' の並びとして以下のように表せる。

$Sr: w e_1 w \cdots w e_2 w \cdots w e_i w \cdots w e_k w$

$Sa: w \cdots e'_1 w \cdots e'_2 w \cdots e'_i w \cdots e'_h w$

$1 \leq i \leq k, 1 \leq j \leq h, h \leq k$

上式において、 i, j はそれぞれ出現順序を表すもので、例えば e_2 と e'_2 が対応しているわけではない。というのは、日本語文では、ヲ格、ニ格、副詞句などは、語順は任意であるからである。

設問形式の特徴で述べたように、記述式設問に対する解答に対しては部分点を与えることが一般的である。上記の考え方に従えば、部分点付与の基準、具体的には解答文中のどの表現に対して部分点をどのくらい付与するかは、概ね次の観点からとなる。

- (a) 正解文中の特定の評価対象表現に対応する解答文に含まれる表現の個数
- (b) 特定の評価対象表現に対応する解答文中の表現の意味的な類似の度合い

採点計算モデル

通常記述式設問では正解文中の個々の評価対象表現はそれぞれ部分点が設定される。これと解答文中の対応する表現の部分点を以下のように表す。

se : 評価対象表現(e)の部分点

se' : 対応表現(e')の部分点

正解文に対する解答文の類似度は、解答文中の対応表現の部分点の総和(s_a)として次式で定義することができる。

$$\sum_{j=1}^h se'_j = s_a \leq p \quad \cdots \cdots \cdots (3)$$

採点手法

解答文の評点は、(3)に示されるように部分点(se'_j)の総和となる。したがって、部分点(se'_j)をどのような方法で決定するかが問題となる。以下の段階を経て決定する。

- (1) 正解文中の評価対象表現の設定
- (2) 解答文に対する形態素一構文解析
- (3) 評価対象表現と対応表現の探索・照合

以下の説明文を読み、設問に答えなさい(10点)。

説明文

2017年5月に、改正個人情報保護法が施行された。この改正では個人情報とされるものの範囲が大きく拡大されている。第二条では、氏名・生年月日等の項目に加え、識別符号等の項目が追加された。

1) その情報に含まれる氏名、生年月日その他の記述などによって、特定の個人を識別できるもの(他の情報と容易に照合でき、それにより特定の個人を識別できるものを含む)

2) 個人識別符号が含まれるもの(特定の個人の身体的特徴を変換した文字、番号、記号などや、カードや書類で個人に割り当てられた文字、番号、記号などで、特定の個人を識別できるものを含む)

二番目の項目で、新たに「個人識別符号」が定義されている。具体的には、指紋データや顔認証データ、虹彩、声紋、DNAのような個人の身体的特徴を変換した文字、番号、記号などや、パスポート番号や運転免許証番号、住民票コード、基礎年金番号、保険証番号のような個人に割り当てられた文字、番号、記号なども個人情報に該当することとなった。

他にも、マイナンバー、各種カード番号、端末ID、アカウントIDなども該当すると考えられ、これらは慎重に扱うべき情報とされ、本人の許諾なく公開することなどには問題がある。

設問

「財布を拾ったので、持ち主を探すため、財布の中に入っていたクレジットカードの番号と拾った状況をSNSで公開して情報提供を求めた」ことが問題となる理由を60文字以内で書きなさい。

図1. 記述式設問の実例。

(4) 得点計算

上記の過程について、図 3.1 に示す高校教科「情報」分野の記述式設問を具体例として以下で上記の処理手続きを説明する。

評価対象表現の設定

正解文を構成する個々の語句に対し、評価対象表現を決定し、それぞれについて部分点を設定する。図 1 に示す記述式設問の正解例を以下に挙げる。

正解文 [44 字]:

クレジットカード番号は個人情報に相当し、
所有者の許可なくそれを公開しているの
問題である。

下線が引かれた語句を評価対象表現とする。この正解文には 4 か所あり、 $p=10$ とし、それぞれの部分点 (se_i) を以下に示す。

- e_1 クレジットカード番号 3 (se_1)
- e_2 個人情報 2 (se_2)
- e_3 所有者の許可なく 3 (se_3)
- e_4 公開している 2 (se_4)

形態素－構文解析

正解文の評価対象表現のうち、 e_1 、 e_2 においては例えば以下に挙げる同義表現がある。

- 同義表現 (e_1): クレジットカードの番号
- 同義表現 (e_2): 個人の情報、個人に関わる情報
プライバシー情報

上の例は、いずれも単語、複合語およびこれらから構成される名詞句である。一方、 e_3 、 e_4 では、前者は動詞（「許可」、動詞としてのサ変名詞）を含み、後者は動詞そのものであるため、同義の様々な活用形態の表現がありうる。したがって、次の (3) 評価対象表現と対応表現の照合の精度を高めるために、正解文、解答文の両者の構成単語と構文構造を明確にする必要があり、以下の解析を行う。

(1) 形態素解析 → (2) 係り受け -- 表層格解析

形態素解析

図 2 に正解文に対する形態素解析の結果（分かち書きされた文）および係り受け解析－表層格解析の結果を示す。

形態素解析では名詞性要素に対しチャンキングを行い、さらに、元の単語とチャンキングによって得られた複合語の対応情報を保持する必要がある。例えば、上記の例では以下となる
「クレジットカード」「番号」→「クレジットカード番号」「個人」「情報」-「個人情報」

この操作により、例えば解答文に「クレジットカ

ードの番号」という句が含まれていても、 e_1 に対応する表現と同定することができる。

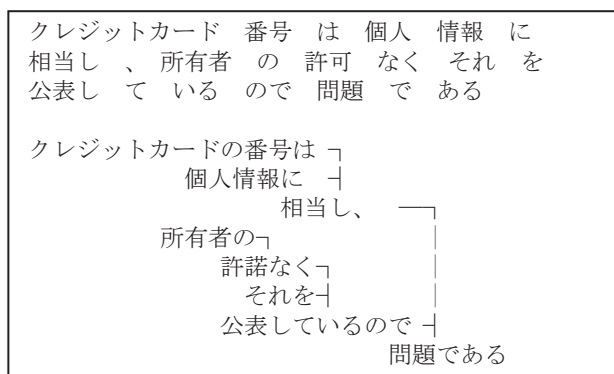


図 2. 正解文の解析結果の実例。

係り受け -- 表層格解析

図 2 の係り受け - 表各解析結果において、図中の「-」の前方の語が係り元、下方の語が係り先、イタリック体は述部、ゴチック体は連用修飾する助詞、アンダーラインは名詞性要素を示している（以下の解析結果の実例も同様の表記）。

表層格解析の結果として、述語とその格要素、および副詞句が決定され、各格要素内の語（句）- 名詞性要素が決定する。具体的には、以下のとおりである。

- 動詞: 相当する
- 格要素 八: クレジットカード番号
- 二: 個人情報

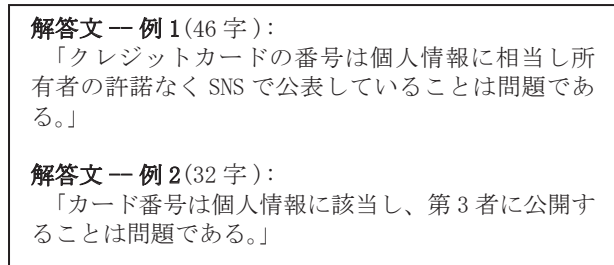


図 3. 解答文の実例。

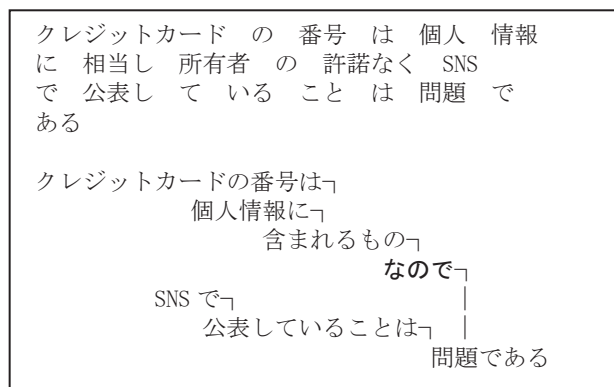


図 4. 解答文 -- 例 1 の解析結果の実例。

解答文の例として、図 3.3 に解答文 - 例 1、解答文 - 例 2 を挙げる。上記解析は、解答文に対しても同様の処理を実施する。

図 4 に、解答文 - 例 1 の形態素解析、係り受け - 表層格解析の結果を示す。

評価対象表現と対応表現の照合

この処理では、正解文および解答文に対する自然言語解析の結果に基づき、それらを照合することで、前者の評価対象表現 (e) に対し、対応表現 (e') 探索し照合する過程である。

解答文には「クレジットカードの番号」という表現があるが、(1) における処理単語 - 複合語の対応から、正解文の「クレジットカード番号」の同義表現であると認定される。また、解答文中には語として、言い換えれば終端記号として「許可」、「公開」は無いが、それぞれの同義語である「許諾」、「公表」が用いられている。以上から、正解文および解答文 - 例 1 について、評価対象表現 (e_i) とその対応表現 (e') は以下のようになる。

評価対象表現 (e)	対応表現 (e')[解答文-例 1]
1 クレジットカード番号	クレジットカードの番号
2 個人情報	個人情報
3 所有者の許可なく	所有者の許諾なく
4 公開している	公表している

図 2 の解答文 - 例 2 では、評価対象表現 (e_i) とその対応表現 (e') は以下のようになる。

評価対象表現 (e)	対応表現 (e')[解答文-例 2]
1 クレジットカード番号	カード番号
2 個人情報	個人情報
3 所有者の許可なく	
4 公開している	公開する

得点計算

この処理では、前項で説明した (3) の結果に基づき、次の 2 つの段階を経て得点計算をする。

(a) 部分点の算出：個々対応表現 (se' j) の評価

(b) 部分点の総和

解答文 - 例 1 では、4 つの評価対象表現に対し同数の対応表現があり、いずれも同義表現であるため、上記の (a) 段階では、全ての対応表現について事前に設定された s_i 部分点が付与される。

$$s_i = se'_j \quad (1 \leq i, j \leq 4 = h)$$

したがって、(b) は以下となる。

$$\sum_{j=1}^4 se'_j = 3 + 2 + 3 + 2 = 10 = p$$

解答文 - 例 2 では、(a) において「カード番号」が「ク

レジットカード番号」と同義とはみなされず、se' 1 = 2 と仮定すれば、se' 3 = 0 であるので、(b) は以下となる。

$$\sum_{j=1}^4 se'_j = 2 + 2 + 0 + 2 = 6$$

採点手法の性能評価

自動採点手法の性能評価は、これを実装した自動採点プログラムが出力した得点と、人間の採点とがどの程度乖離しているか、乖離の度合いということになる。ある 1 つの記述式設問に対して、n 人の受験者 (解答数 = n) について、採点者の得点を sh、自動採点プログラムの得点を sc としたとき、全受験者間について両者間での相違の平均を d とすると、d は以下の式で与えられる。

$$\frac{1}{n} \sum_{i=1}^n |sh_i - sc_i| = d \cdot \dots \cdot (4)$$

上式において、d が 0 に近づくほど人間の採点との一致度合いが高くなる。試験という公正性を最重要課題とする状況では、そうする必要はある。

討論

採点計算モデルの拡張

設問形式の種類で述べたように記述式設問では、受験者の能力の評価項目として (a) 知識、(b) 思考力、(c) 表現力を挙げている。当該設問の属する科目、また作成者の意図によってはこれらの評価項目の比重は異なってくる。例えば図 1 の設問の解答文 - 例 3 として、以下を仮定する。

「クレジットカード番号はこじんじょうほうに相当し、しょゆうしゃのきよかなく公開しているので問題である。」

採点手法で述べた採点方法では、解答文 - 例 3 の得点 p = 10 となる。表現力を考慮した場合、p = 10 は不適当であろう。したがって、採点計算モデルとして (3) 式を拡張し以下の式 (5) を導入する。

$$\sum_{j=1}^h se'_j + C = s_a \leq p \quad \dots \dots (5)$$

上式で C は 評価対象表現以外の採点上考慮すべき項目の値で、試験個別に設定する値である。

同義表現の問題

採点手法の性能評価の (4) で d が大きくなる要因として、同義表現の問題が挙げられる。これは次の 2 項目が挙げられる。

(a) 個々の評価対象表現 (se_i) の同義表現の数

(b) 正解文に含まれる評価対象表現の数 (k)

(a) については、1つの評価対象表現 (se_i) について、語、句レベルで様々な同義表現が用いられる。同義表現の個数は、それぞれの評価対象表現によって異なる。同義表現の個数が多いほど、採点手法で述べた (3) の処理、すなわち解答文中の対応表現を同定することが困難となり、採点者が認識した対応表現の数 h よりも小さくなる。

(b) は指定された文字数と関連する問題である。文字数が長いほど構成単語数が多くなり、評価対象表現数 (k) が増加する。 k が増加すると、(a) の問題が累積し、(4) で d の増加をもたらす。上記問題の軽減手段としては、次の2つが考えられる。

(1) 文字数以外の制約、(2) 同義表現の収集

解答文に対する制約

長さ以外の制約条件として以下が挙げられる。

(a) 文の構成要素 (b) 文の構文

(a) としては、表 2.1 の記述式設問の説明で述べたように、解答文中に使用すべき語・句を、具体的に例示することである。(b) としては、以下の制約が挙げられる。

- ・ 文の構文の規定 (例: 「・・・の理由を述べよ」)
- ・ 文頭または、文末表現の規定

同義表現の収集

同義表現の収集のアプローチとして次の2つが挙げられる。

(1) 当該問題に限定 (2) 当該科目全般

(1) は、個別の記述式設問毎に評価対象表現 (se_i) の同義表現を収集するというものである。したがって、試験毎に毎回收集する必要がある。(2) は、記述式設問の科目について、学ぶべき概念を中心に、同義語を収集するというアプローチである。科目についてある種のシソーラスを構築すると考えてよい。シソーラスに含まれる用語が網羅的であれば、採点者の有する知識に近づき、採点計算モデルの拡張で述べた問題の減少をもたらす。また、記述式設問作成で (1) の作業の支援としても有用である。

同義語の収集方法

特定科目の教科書/参考書で説明される用語、および一般的な用語についての同義語を収集する方法としては以下が挙げられる。

(1) 既存の言語資源の利用

(2) 当該科目のコーパスに対する自然言語解析

(1) としては、電子媒体の辞書 (EDR 電子辞書等)、

類語辞典 (分類語彙表等)、シソーラス (JST 科学技術用語シソーラス等) が利用できる。

(1) では、語 (単語、複合語) 形態のいわゆる同義語は収集可能であるが、句レベル、言い換えれば<名詞 (句)、動詞> の対としての同義表現を得ることは困難である。(2) として、当該科目について電子媒体の教科書・参考書を対象に、形態素一構文解析で述べた手法を適用する必要がある。この作業としては、自然言語研究で使用されているオープンソースの形態素解析エンジン MeCab⁶⁾、日本語係り受け解析器 CaboCha⁷⁾ 等の解析ツールが利用できる。

自動採点システムの運用

はじめにで挙げた記述式設問では、採点に際し複数名による採点が行われている。採点手法の性能評価で説明した (4) 式において、 sh を採点者 A による得点 sA 、 sc を採点者 B による得点 sB に置き換えれば、 d は両採点者間の相違を示すことを意味する。採点者間の相違について、例えば $d < 0.1$ とならないと、当該問題に対して公正性の問題が生じ、採点が無効となる可能性が生ずる。記述式設問に対し自動採点システムを導入するためには、上記の条件を満たしていることを前提とすることが妥当と考えられる。

終わりに

今後の課題として以下が挙げられる。

(a) 正解文の別解・多義表現の研究

(b) 本モデルに基づく実証研究

(a) について、記述式設問の設問形式上の最大の問題は、別解・多義表現が多数出現し、採点者間に差異が生じることである。自然言語解析の観点からの実証的な研究が必要である。

(b) について、本稿はモデル提示に留まったが、具体的な問題・設問を作成し、数十名程度の学生に対して実施し当モデルによる実証実験が求められる。

謝辞

本研究は 2016 年度神奈川大学総合理学研究所共同研究助成「短い記述式解答の自動採点に向けた日本語文解析手法の検討」(RIIS201705) を受けて行った。記して感謝する。

文献

- 1) 大学入試センター (2017) 大学入学共通テスト実施に向けた検討状況. [http://www.dnc.ac.jp/daigaku-nyugakukibousyagakuryokuhyoka_test/progress.html].
- 2) 大学入試センター (2017) 大学入学共通テスト・平成 29 年度試行調査・問題、正解表、解答用紙等.

- [http://www.dnc.ac.jp/daigakunyugakukibousyagakuryokuhyoka_test/pre-test_h29_01.html].
- 3) 神奈川県教育委員会 (2016) 県立高等学校入学者選抜調査改善委員会中間とりまとめ検討資料. [<http://www.pref.kanagawa.jp/uploaded/attachment/826441.pdf>].
 - 4) 文部科学省 (2015) 思考力・判断力・表現力を問う条件付記述式問題について (たたき台). 高大接続システム改革会議 (第9回) 配付資料, 別紙 2-2, 2015年12月22日. [http://www.mext.go.jp/b_menu/shingi/chousa/shougai/033/shiryo/_icsFiles/afieldfile/2015/12/22/1365554_04_1.pdf].
 - 5) 大学入試センター (2017) 大学入学共通テスト・平成29年度試行調査: 国語・正答例および正答の条件. [http://www.dnc.ac.jp/albums/abm.php?f=abm00011249.pdf&n=01_%E5%9B%BD%E8%AA%9E.pdf].
 - 6) Kudo T (2006) MeCab: Yet Another Part-of-Speech and Morphological Analyzer. [<http://taku910.github.io/mecab/>].
 - 7) Kudo T (2005) CaboCha/ 南瓜: Yet Another Japanese Dependency Structure Analyzer. [<http://taku910.github.io/cabocha/>].