

試験問題画像における XML 構造と属性の推定 Estimation of Attributes and Structured from Exam Questions

松本 涼[†] 遠藤 聡志[‡]
Ryo Matsumoto Satoshi Endo

1 はじめに

人工知能技術で大学入試問題を解くことに挑戦する「ロボットは東大に入れるか」(以下、「東ロボ」)プロジェクト [1] は、大学入試センター試験(以下、「センター試験」)の模擬試験において、2016 年時点で全教科合計の偏差値 57.1 を達成した。しかしながら、その入力人手で電子化、アノテーションされた XML データであり [2]、「XML 化にコストがかかる」、「人が XML データを作成しないと解答動作を行えない」などの課題がある。従って、センター試験の画像から、End-to-End で XML データを生成するシステムが必要となるが、そのためには、抽象的な複数の情報の抽出、推定および統合を行う必要がある。

磯崎ら [4] は、センター試験英語の画像データに文字認識(OCR)を行い、得られた文字列にルールベースで XML タグを付与する手法を提案しており、いくつかの情報について抽出・推定、および XML 化を実現した。しかしながら、全ての情報の抽出、推定は実現できていない。また、XML データの精度評価も行われていない。

そこで本研究では、XML の仕様を参考に、XML データを要素、属性、構造、内容の 4 種類の情報に分けて抽出、推定し、精度評価を行うことを提案する。著者ら [5] はこれまでに、要素の抽出と精度評価を行なった。よって本稿では、属性および構造の推定と精度評価を行う。また、End-to-End での XML データ生成について、考察を行う。

2 研究背景

2.1 東ロボくん プロジェクト

東ロボプロジェクト [1] は、大学入試センター試験および二次試験問題を解く人工知能の開発を目的としている。2015 年時点で、センター試験模試の合計点が学生の平均点を超えており、2019 年センター試験の英語筆記本試験において、185 点(偏差値 64.1)を達成した*1。しかしながら、その目的は「試験問題に正解すること」であり、問題を(解く前の)読み取る動作において、人の手作業を許容している。すなわち、入力データは画像形式ではなく、計算機にとって扱いやすい XML 形式を前提としており [1][2][3]、XML データの作成は、人が行う必要がある [4]。

2.2 センター試験 XML データ

センター試験 XML データは、1990~2017 年度の大学入試センター試験問題を、XML 形式で電子化したデー

[†] 琉球大学大学院理工学研究科情報工学専攻, Graduate School of Engineering and Science, University of the Ryukyus

[‡] 琉球大学工学部工学科知能情報コース, Computer Science and Intelligent Systems, University of the Ryukyus

*1 <https://www.nii.ac.jp/news/release/2019/1118.html>

タセットである。東ロボくんプロジェクトが公開*2しており、プロジェクトが独自に策定した XML ベースのアノテーション仕様(以下、「東ロボ XML 仕様」)に従って記述されている。

3 End-to-End 自動 XML 化システム

XML(eXtensible Markup Language)*3 は、あらゆる文書を電子化する目的で策定されたマークアップ言語である。本研究は、XML を要素、属性、構造、内容に分解し、それぞれについて情報の抽出・推定を行い、得られた情報を XML データとして統合する End-to-End モデル(以下、「E2E XML 生成システム」)を提案することを目指している。

3.1 XML タグ、内容について

XML において、デザインや構造、意味などのメタ情報は、図 1 のように、文や図表などを XML タグでマークアップすることで表現される。

```
<XMLタグ1>
  <XMLタグ2>
    単語、文、文章、図表、地図、...
  </XMLタグ2>
</XMLタグ1>
```

図 1 XML タグと内容の位置関係

XML タグには、start-tag と end-tag がある。一般的に、XML タグは start-tag と end-tag で囲まれた情報(以下、「内容情報」)にメタ情報を付与する。

内容情報には、文章情報や画像情報、音声、動画情報、XML データなどがある。センター試験 XML データには、問題文や図、表の各セルの値、図表のキャプション文、数式の記号などがある。

3.2 要素、属性について

XML タグは、要素(element)と属性(attribute)情報を持つ。センター試験 XML データの例を、図 2 に示す。図 2 は、2011 年の地理 B 科目、第 1 問の間 2 (一部省略)であり、左側が写真、右側が XML データである。XML データに、要素を赤四角で、属性を青四角でそれぞれ示す。

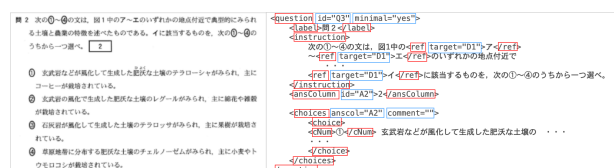


図 2 要素、属性の例

*2 <https://21robot.org/dataset.html>

*3 Extensible Markup Language (XML) 1.1 (Second Edition) W3C Recommendation <https://www.w3.org/TR/xml11/>

3.2.1 要素の抽出

要素はXMLタグの先頭(左端)に1つだけ指定でき、XMLタグの名前を表す。なお、著者らはこれまでに、センター試験画像からの要素情報の抽出に物体検出モデルを適用する実験を行い、全ての要素が、平均適合率 mAP=0.8 以上の精度で抽出可能であることを示した [5]。要素情報の抽出に物体検出を適用する手法は、論文 PDF[6] や、古典・近代書籍 [7][8]、などにも適用され、高精度で抽出できることが確かめられている。

3.2.2 属性の推定

属性は、要素の補足説明であり、要素の右側に必要なだけ記述できる。また、属性は name と value からなり、name=value の形で表記される。例えば、文書内の写真に、要素名が data である XML タグを付与し、その type が“写真”であることを表現したい場合、XML データは図3のようになる。

```
<data type="写真" size="640x320" ... >
  写真データ
</data>
```

図3 属性の例

なお、本研究では属性を @name と表記する。例えば図3は、「(要素 data) の @type は“写真”である」となる。表1に、東ロボ XML 仕様における属性(一部抜粋)を示す。

表1 東ロボ XML 仕様における属性(一部抜粋)

要素	属性	意味
exam	@year	年度
	@subject	科目名
	@range_of_options	選択問題の範囲
	@num_of_options	選択する選択問題の数
question	@minimal	それ以上小問を含まない?
	@answer_style	回答方式(選択, 記述, ...)
	@answer_type	答えのタイプ
	@knowledge_type	必要となる知識
	@anscol	対応する解答欄の @id の値
data	@type	text, image, table, complex
cell	@rowspan	(表内に、横方向のセルの連結がある場合)結合されているセルの数
	@colspan	(表内に、縦方向のセルの連結がある場合)結合されているセルの数
ref	@target	参照記号の参照先の要素の @id の値
blank	@digits	空欄に指定された桁数

表1の他に、@id (refの @target などの参照のための、XMLタグの識別番号)や @comment (コメント(任意))、examにおける @source (問題の出題機関(大学入試センター))、@srcTxtURL (問題のURL(任意))、誤字脱字や欠落に対応するための @correction (文字化け箇所の正しい文字)、@string (テキスト欠落部分の正しい文字)がある。

3.3 構造について

要素が複数ある場合、要素の順序は木構造(特に内包、並列関係)を表現する。例えば、図3の写真が、センター試験の選択問題における選択肢であるとき、XMLタグは図4のように表現される。

```
<choices>
  <choice> <data type="写真">
    写真データ1
  </data> </choice>
  <choice> <data type="写真">
    写真データ2
  ...
</data> </choice>
</choices>
```

図4 構造の例

図4において、choices と choice は親子関係にあるため、順序を入れ替えられない。choice と data は、start-tag と end-tag の位置関係が同じため、入れ替えることは可能である。ただし、センター試験においては、choice と data の間に cNum (選択肢番号。①, ②, ...) が入ることがあり、その場合は、choice は cNum と data の親要素である必要がある。このように、XMLタグの順序は重要な意味(以下、「構造情報」)を持つ。

3.4 E2E XML 生成システム

XMLデータは、図5のような配置で、要素(赤)、属性(青)、構造(タグの順序)、内容(緑)の4種類の情報が記述されたデータである。

```
<要素1 @属性1 @属性2 ... >
  <要素2 @属性3 @属性4 ... >
    内容
  </要素2 >
</要素1 >
```

図5 要素、属性、構造、内容の記述方法

センター試験の画像から End-to-End で XML を自動生成するシステムの実現には、これら4つの情報を抽出、推定し、XMLに統合する必要がある。本稿では、そのうちの属性と構造の推定実験を行う。また、E2E XML 生成システムの実装の課題について、考察を行う。

4 実験

センター試験 XML データの全情報の抽出、推定という目標のために、属性および構造情報の推定実験を行った。本実験の目的は、以下の2つである。

- 属性と構造の推定において利用可能な情報を、要素、属性、構造、内容の視点で整理する。
- 一部の属性と構造の情報について、機械学習アルゴリズムで推定実験を行い、推定精度のベースラインを示す。

4.1 属性、構造の推定に利用可能な情報

本実験で利用可能な情報は、「推定対象の情報が含まれるページの画像データ(以下、「ページ画像」)」と「XMLデータ」である。

「XMLデータ」のうち、属性、構造の推定どちらにも利用可能な情報を、「ページ画像」、「要素(名前、領域)」、「内容(文章)」とする。また、属性の推定には、「構造」および「推定対象ではない属性」の情報を、構

造の推定には、「属性」の情報を、それぞれ利用できるものとする。すなわち、属性または構造を推定する際、その他の情報は完全に取得できた(理想的な)状態を前提として、実験を行う。

4.2 属性推定

4.2.1 推定対象の情報

本実験では、属性 @answer_type と @knowledge_type の属性値の推定を行った。表2に、それぞれの属性値の例を示す。@answer_type は、問題の答えのタイプを14クラス(複数指定可能)で分類し、@knowledge_type は、問題の解答に必要なとなる知識を、10以上のクラス(複数指定可能)で分類している。("DM"のみ、科目ごとに異なる子クラスを持つ。英語科目には、辞書的知識 (PRN, DIC.O,...) や文章読解 (R_QA, R_SUM,...)、グラフ・図表の読み取り (IC_G, IC_P, ...) など、計21クラスがある)

表2 属性値の例

属性名	属性値の例
answer_type	sentence, term(person, location, ...), formula, image(graph, map, ...), referenceSymbol, orthography, symbol-symbol-symbol, ... など
knowledge_type	KS(外部の知識源), RT(資料テキストからの判断), IC(グラフ、図表読み取り), GK(一般的知識), DM(ドメイン依存モデル), 'DIC.O,GK', 'R.QA,IC.O' ... など

@answer_type と @knowledge_type どちらも複数ラベルを指定可能であるため、マルチラベル分類問題であるが、本実験では、XML データセットに存在するクラスのみが多クラス分類問題として推定した。

4.2.2 データセット

1987 ~ 2017 年度の英語(筆記)科目のXMLデータからデータセットを作成した。@answer_type、@knowledge_type は、要素 question(小問)の属性であるため、入力データは question に関する情報が望ましいと考えられる。よって本実験では、「直前の大問の問題文(instruction)」(内容)と、「question タグで囲まれたXMLデータ」(要素(名前)、属性、構造、内容)を入力データとした。入力データの例を、図6に示す。

Unnamed: 0	answer_type	<instruction/>	contents
0	(symbol-sentence)*2	次の問い(問1・問2)において、 下線部(a)・(b)の単語のアクセント(強勢)の位置が正しい...	<label>問1</label> <ansColumn id="A1">1</ansC...
1	(symbol-sentence)*2	次の問い(問1・問2)において、 下線部(a)・(b)の単語のアクセント(強勢)の位置が正しい...	<label>問2</label> <ansColumn id="A2">2</ansC...
2	sentence	次の会話の下線部(1)~(4)について、 それぞれ以下の問い(問1~4)に示された○~◎のうち...	<label>問1</label> <ansColumn id="A3">3</ansC...

図6 属性推定における入力データの例

4.2.3 モデル

@answer_type、@knowledge_type 属性それぞれについて、LSTM 文章分類モデルで学習を行った。Tokenizer は、入力がXMLデータであることから、Sentence Piece[9]を使用した。Tokenize 前後のデータの例を、図7に示す。

```
<label>問1</label>
<data id="D20" type="text">
  <ansColumn id="A21">21</ansColumn>
</data>

['<', 'label', '>', '問', '1', '</', 'label', '>', '_type', '=', '1', 'text', '1', '>',
'_id', '=', 'D', '20', 'type', '=', '1', 'text', '1', '>',
'_id', '=', 'A', '21', '1', '>', '2', '1', '</', 'ansColumn', '>',
'</', 'data', '>', ]
```

図7 Sentence Piece による Tokenize 前後の入力データ

4.2.4 結果と考察

@answer_type、@knowledge_type の属性値の、LSTM 文章分類モデルによる推定精度(正解率)を表3に示す。

表3 属性値の推定精度

属性名	正解率
@answer_type	0.93
@knowledge_type	0.42

この結果から、@answer_type より @knowledge_type の方が推定が難しい情報であると言える。原因はいくつか考えられるが、まず、@answer_type は29クラスの分類問題であるのに対し、@knowledge_type は138クラスあることが挙げられる。また、@answer_type は「選択肢のデータ形式」から予測できるが、@knowledge_type は「問題を解くにはどこを検索すれば良いか」という情報であり、問題文をより詳細に読む必要があることも大きく影響していると考えられる。

なお、3つ目の原因として、入力データのXMLタグが問題文を詳細に読む際のノイズになる可能性もあるため、XMLタグを除いての実験も行なった。結果は、正解率が0.45であり、多少の精度向上は見られたが、本質的な改善には至らなかった。@knowledge_type の推定精度向上は今後の課題であるが、解決方針として、モデルの変更(マルチラベル分類モデルや、Bidirectional LSTM モデルなどの利用)を考えている。

4.3 構造推定

本実験では、構造情報(要素の順序)の推定を行った。E2E XML生成システムにおいて、構造推定までの処理の過程は、図8になる予定である。

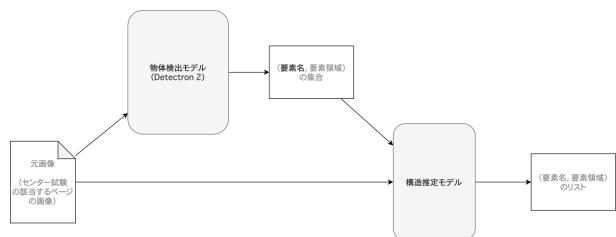


図8 構造推定までの処理の流れ

構造推定は、要素の抽出(物体検出)モデルから得られた「要素(名前, 領域)の集合」と「元画像」を入力として、正しい順序にソートされた「要素のリスト」を出

