

# 全地方議会会議録の横断検索に向けた データ収集とデータ構造の検討

井原大将<sup>1,a)</sup> 内田ゆず<sup>2,b)</sup> 高丸圭一<sup>3,c)</sup> 木村泰知<sup>4,d)</sup> 江崎浩<sup>1,e)</sup>

概要：全国には都道府県・市・特別区・町・村を合わせて、1788の地方自治体が存在しており、このうち約86%がウェブ上で地方議会会議録を公開している。しかしながら、ウェブ上での会議録の公開方法やデータ形式は自治体により異なっており、横断検索や集計などにおいてそれらを統一的に扱うのは難しい。そこで本稿では、収集や整理がしづらい会議録を対象として、それらのデータ収集方法における問題点とその対処法について述べ、横断検索や集計が可能となるデータ構造を提案する。

HIROMASA IHARA<sup>1,a)</sup> YUZU UCHIDA<sup>2,b)</sup> KEIICHI TAKAMARU<sup>3,c)</sup> YASUTOMO KIMURA<sup>4,d)</sup>  
HIROSHI ESAKI<sup>1,e)</sup>

## 1. はじめに

近年、地方自治体による公共データのウェブ公開、いわゆるオープンデータが進められている。地方政治における議論の過程がすべて記録された地方議会会議録も活用が期待される公共データの一つである。地方議会会議録の活用を推進する取り組みでは、複数の自治体の会議録を横断的に取り扱うことが要求される。たとえば「地方議会議事録横断検索 chiholog」では、798の自治体から収集した地方議会会議録の全文検索や単語出現頻度の時系列分析機能等をウェブサービスとして提供している。この取り組みはオープンデータの活用例であるとして、VLED 利活用普及委員会勝手表彰において表彰を受けている\*1。また、木村らは地方議会会議録コーパスを用いた学際的応用研究を進めており、情報工学、社会言語学、計量経済学などの分野で成果をあげつつある [1][2][3][4][5]。このような地方議会

会議録を対象としたウェブサービスやコーパスを構築するためには、地方自治体がウェブサイトにおいて個別に公開している会議録を収集し、共通のデータ形式に整理した上で、データベースに登録する必要がある。しかしながら、会議録のウェブ公開方法は統一されていないため、公開方法に応じた収集・整理の手法を用意する必要がある。地方議会会議録のウェブ公開は、全文検索システムによる公開と、PDF等のファイルによる公開に大別される。上述のウェブサービスやコーパスは現段階では、前者のみを対象としている。本研究の最終的な目標は、多様な形式で公開される地方議会会議録をすべて収集し、横断的な分析を可能にすることである。このため、PDF等で公開される会議録も収集・整理の対象とする必要がある。そこで本稿では、まず地方議会会議録の公開状況と公開形式について述べる。次に、さまざまな形式で公開される会議録を統一的に扱うためのデータ構造の提案と、実際のデータ収集において生じた問題への対処について述べる。最後に、提案したデータ構造を用いて行ったデータ収集・整理の結果について述べる。

## 2. 地方議会会議録の Web 公開の状況と公開形式

地方議会会議録の公開には大きく分けて、会議録システムとファイルの2つの形式がある。

<sup>1</sup> 東京大学  
The University of Tokyo  
<sup>2</sup> 北海学園大学  
Hokkai Gakuen University  
<sup>3</sup> 宇都宮共和大学  
Utsunomiya Kyowa University  
<sup>4</sup> 小樽商科大学  
Otaru University of Commerce  
a) taisyo@hongo.wide.ad.jp  
b) yuzu@eli.hokkai-s-u.ac.jp  
c) takamaru@kyowa-u.ac.jp  
d) kimura@res.otaru-uc.ac.jp  
e) hiroshi@wide.ad.jp  
\*1 <https://bitlet.co/>

## 2.1 会議録システムによる公開

会議録システムとは、議会の会議録を Web 上で閲覧/検索できるようにした Web システムであり、このシステムを通して Web ページとして会議録が提供される(図1に千葉県議会を例を示す)。ユーザは開催日や議会の種別、発言者、発言内容などを指定して検索し、閲覧する。会議録システムは、基本的には地方議会ごとに異なる、ソフトウェアベンダーによって開発されたシステムを導入している場合が多い。主要な会議録システムを表1に挙げる。

表1 主要な会議録システム

会議録システム	ソフトウェアベンダー
DiscussNet	NTT アドバンステクノロジー
DB-Search	大和速記情報センター
Sophia	神戸総合速記
VOICES	フューチャーイン



図2 ファイルによる公開の画面例(茂原市議会)



図1 会議録検索システムによる公開の画面例(千葉県議会)

## 2.2 ファイルによる公開

ファイルによる公開では、地方議会や地方自治体のウェブサイト上で、ファイルとして会議録が提供される(図2に茂原市議会の例を示す)。ユーザはウェブサイトアクセスし、ファイルをダウンロードして閲覧する。多くの場合、ファイルの形式はPDFファイルやMicrosoft Wordファイルである。

## 3. Web 公開の会議録のデータ収集における問題

Web で公開されている会議録を収集するには、解決しなければならない2つの問題がある。これらの問題はどちらの公開形式にも共通である。

### 3.1 会議録のリソースをいかに見つけるか

Web 公開において、会議録のデータは地方議会のウェブサイト内のファイル(PDFファイルやMicrosoft Wordファイル)や会議録システム上のWeb ページとして提供される。しかしながら、それらはユーザがWeb ブラウザを通してアクセスすることが前提となっており、ユーザは視覚的なボタン配置やリンクテキストの情報に基づいて会議録のリソースへ辿りつく。そのため、コンピュータでクロールする場合は、クローラに会議録へのリソースへたどり着くための情報、例えば、リンクの辿り方や書き換え方などを与える必要がある。

### 3.2 会議録のデータをいかにデータ構造に落とし込むか

収集した会議録データは統一的なデータ構造で管理する必要がある。Web 公開において、会議録のデータのファイル形式、フォーマットは地方議会によって異なる。そのため、統一的に会議録を扱うためには、地方議会ごとの会議録のデータのファイル形式やフォーマットなどデータの構成を考慮し、データ構造を定義する必要がある。また、コンピュータでクロールする場合は、ファイル形式やフォー

マットの差異を吸収し、同じデータ構造にするための情報、例えば、”HTMLのあるタグでマークアップされたテキストは会議名に変換する”等の変換ルールを与える必要がある。

#### 4. 地方議会議録横断検索の要件とデータ構造の定義

ここで、地方議会議録横断検索の要件を下記に挙げる。

- (1) 対応する地方議会すべてに対して同時に横断的に検索できる。
- (2) 会議録の本文や会議名などを対象としたキーワード検索ができる。
- (3) 検索条件として開催日や自治体名を使用できる。

次に、この要件のためのデータ構造を定義し、その定義を表2に示す。斎藤ら[6]はさらに細かい粒度のデータ構造を定義している。例えば、会議録本文を発言者、発言内容のフィールドに分割し、会議名の他に回や号などのフィールドを定義している。本研究では、すべての地方議会に対応するためのクローラの作成における実装上の難しさや、すべての地方議会に対応できるデータ構造の抽象化の難しさを考慮し、簡素なデータ構造を採用する。

表2 地方議会議録横断検索のデータ構造定義

項目	型	備考
地方議会 ID	数値型	地方議会を表す番号。 地方公共団体コードを用いる。
会議録本文	文字列型	
会議録会議名	文字列型	
会議開催日	日付型	

#### 5. データ収集における問題への対処

会議録の収集では、地方議会の会議録の Web 公開ページにクローラでアクセスし、会議録データの探索と保存を行う。地方議会ごとに異なる URL や HTML マークアップを用いているため、クローラは地方議会ごとに作らなければならない。しかし一方で、同じソフトウェアベンダーによる会議録システムや同じファイル形式を使っている自治体間では、URL や HTML マークアップに共通部分が多いこともある。そこで、今回は会議録システムごとのクローラとファイルによる公開に分けてクローラを作成する。本章では、クローラの具体的な動作を会議録の公開形式別に述べる。

##### 5.1 会議録システムによる公開の場合

###### 5.1.1 会議録データの探索

会議録システムは、会議録の検索と閲覧の機能を提供する。そして、会議録の閲覧や検索結果の表示のための URL はシステムによって動的に生成される。そのため、会議録

の検索結果を表すページや会議録データを表すページのリンクには規則性がある。例えば、図3の例では、議事録データページの URL は Template=DocOneFrame というパラメータを持っている。これを利用し、通常ユーザーが、検索機能によって目的の会議録一覧を表示し、そこから目的の会議録を閲覧する場合と同じように、クローラも同様の手順でアクセスを行う。具体的には、すべての会議録に対して検索条件を指定せずに検索を行い、すべての会議録の検索結果を表示し、個々の会議録データページへのリンクを取得する。ページが分かれている場合は、再帰的に次のページの会議録一覧を表示し、会議録リンクを取得する。会議録システムや地方議会によっては、会議録リンクが、会議録全文表示に対するリンクでなく、発言単位表示のリンクや HTML フレームへのリンクとなっている。その場合は、URL のパラメータを置換し、会議録本文表示への会議録リンクへと修正する。例えば、発言単位表示のテンプレートを指定するパラメータを会議録全文表示へ置換する。

```

1 <TR>
2 <TD align="right" nowrap><SPAN class="BasicRed">1</SPAN></TD>
3 <TD><SPAN class="Basic">2017.04.17</SPAN></TD>
4 <TD nowrap>
5 <SPAN class="Basic">
6 <A href="/chiba-c/dsweb.cgi/documentframe!!guest01
!!4744!!1!!1,-1,1!4629!312126!1,-1,1!4629!312126!2,0,2!4629!312126!319064!!1?
Template=DocOneFrame" target="_self">平成
29年新庁舎整備調査特別委員会 名簿</A>
7 </SPAN>
8 <SPAN class="Basic">
9 <A href="?Template=DocOneFrame" target="_self"></A>
10 <BR>
11 </TD>
12 </TR>
13 <TR>
14 <TD align="right" nowrap><SPAN class="BasicRed">2</SPAN></TD>
15 <TD><SPAN class="Basic">2017.04.17</SPAN></TD>
16 <TD nowrap>
17 <SPAN class="Basic">
18 <A href="/chiba-c/dsweb.cgi/documentframe!!guest01
!!4744!!1!!1,-1,1!4629!312126!1,-1,1!4629!312126!2,0,2!4629!312126!319066!!1?2
Template=DocOneFrame" target="_self">平成
29年新庁舎整備調査特別委員会 本文</A>
19 </SPAN>
20 <SPAN class="Basic">
21 <A href="?Template=DocOneFrame" target="_self"></A>
22 <BR>
23 </TD>
24 </TR>
25 <TR>
26 <TD align="right" nowrap><SPAN class="BasicRed">3</SPAN></TD>
27 <TD><SPAN class="Basic">2017.03.15</SPAN></TD>
28 <TD nowrap>
29 <SPAN class="Basic">
30 <A href="/chiba-c/dsweb.cgi/documentframe!!guest01
!!4744!!1!!1,-1,1!4629!312126!1,-1,1!4629!312126!2,0,2!4629!312126!326714!!1?3
Template=DocOneFrame" target="_self">平成
29年第1回定例会(第9日目) 議事日程・名簿</A>
31 </SPAN>
32 <SPAN class="Basic">
33 <A href="?Template=DocOneFrame" target="_self"></A>
34 <BR>
35 </TD>
36 </TR>

```

図3 会議録検索システムの検索結果における URL の記述例(千葉市議会)

###### 5.1.2 会議名抽出

検索結果のページにある会議録データページへのリンクテキスト (<a>タグのテキスト) や議事録データページのページタイトル (<title>タグのテキスト) を会議名として抽出する。

###### 5.1.3 会議録本文抽出

会議録データページの HTML テキストから HTML タグを取り除き、本文を抽出する。ただし、一部の HTML タグが議事録の書式の一部となっている場合は必要に応じてその書式へ置換する。例えば、改行を表すタグ (<br>タグ) が該当する。

### 5.1.4 開催日抽出

会議名や会議録本文から抽出するが、開催日の表記には様々なパターンがあるため、それを考慮した抽出のルール(正規表現)を作成する。考慮すべき揺れや書式の例を次に示す。

- 半角数字や全角数字の混在
- 空白文字の混在
- 西暦表記, 元号表記
- / (スラッシュ), - (ハイフン), , (カンマ) などの年月日区切り文字

### 5.1.5 JavaScript で動的に生成されるページコンテンツ

一部の会議録システムでは、会議録データを見つけるのに必要な情報が、JavaScript によって生成されることがある。例えば、図4の例では、議事録データページへのURLがJavaScriptのdocument.write()によって生成されている。通常クローラはHTTP[S]クライアントとして動作し、サーバーに対してHTMLドキュメントを要求し、取得するだけであるため、Webブラウザ上で生成されるコンテンツには対応していない。JavaScriptで動的にコンテンツを生成する会議録システムや会議録システムのあるページに対しては、クローラ標準のHTTP[S]クライアントではなく、Webブラウザを外部からプログラムで操作してJavaScriptを実行し、コンテンツを得る必要がある。

```

1 <SCRIPT type="text/javascript">
2 <!--
3   if( "1" == "1" ){
4     document.write('<SPAN class="BoldRed">1</SPAN>');
5   }
6   else{
7     document.write('<A href="/kisarazu-c/dsweb/cgi/query!!1!guest02
8       !!2200!!1!011,-1,!!2986!209134!1,-1,!!2986!209134!2,0,2!2986!209134!!1!7Template
9       =List&QuerySelect=No&List=now&DocumentType=" target="_self">1</A>');
10  }
11 </SCRIPT>
12 <SCRIPT type="text/javascript">
13 <!--
14   if( "2" == "1" ){
15     document.write('<SPAN class="BoldRed">2</SPAN>');
16   }
17   else{
18     document.write('<A href="/kisarazu-c/dsweb/cgi/query!!1!guest02
19       !!2200!!1!011,-1,!!2986!209134!1,-1,!!2986!209134!2,0,2!2986!209134!!1!3?
20       Template=List&QuerySelect=No&List=now&DocumentType=" target="_self">2</A>');
21  }
22 </SCRIPT>
23 <SCRIPT type="text/javascript">
24 <!--
25   if( "3" == "1" ){
26     document.write('<SPAN class="BoldRed">3</SPAN>');
27   }
28   else{
29     document.write('<A href="/kisarazu-c/dsweb/cgi/query!!1!guest02
30       !!2200!!1!011,-1,!!2986!209134!1,-1,!!2986!209134!2,0,2!2986!209134!!1!25?
31       Template=List&QuerySelect=No&List=now&DocumentType=" target="_self">3</A>');
32  }
33 </SCRIPT>

```

図4 会議録検索システムのJavaScriptによって生成されたコンテンツの記述例(木更津市議会)

## 5.2 ファイルによる公開の場合

### 5.2.1 会議録データの探索

ファイルによる公開では、地方議会の会議録公開Webページに一覧として会議録データへのリンクがあるか、階層的になっていて例えば、"平成29年会議録" → "定例会" → "第1回定例会第4号会議録本文.pdf"のようにリンクが階層化されていることが多い(図5に一宮町議会の例を示す)。そこで、クローラには地方議会ごとにリンク深さ、

リンクURLフォーマット、リンクテキストフォーマットを与え、指定されたリンク深さのリンクをすべて走査し、そのうちURLのフォーマットとリンクテキストに合致するリンクを会議録データへのリンクとして取得する。URLのフォーマット(正規表現)の例としては".\*\.\pdf", リンクテキストには"会議録"や"本文"などを指定する。



図5 ファイルによる公開のリンクを辿る画面例(一宮町議会)

### 5.2.2 会議名抽出

先程述べた指定されたリンクの深さまで走査するとき、各ページにおいてページタイトル(<title>タグのテキスト)とリンクテキスト(<a>タグのテキスト)と直上のHeadingテキスト(<h6>, <h5>, <h4>, <h3>, <h2>, <h1>タグのテキスト)を収集し、それらから、"定例会"や"回", "号"を含むものを抽出する。例えば、図5の一宮町議会では各ページにおいて、ページタイトルが"議会会議録 | 一宮町役場" → "平成28年 会議録 | 一宮町役場", リンクテキストが"平成28年 会議録" → "平成28年第4回定例会会議録", Headingテキストが"議会会議録" → "平成28年 会議録"になり、会議名として"平成28年第4回定例会会議録"が抽出される。

### 5.2.3 会議録本文抽出

ファイル形式に応じて、変換ツールを用いて本文抽出を

行う。今回のクローラでは、PDF形式のファイルのためにPDFMiner<sup>\*2</sup>を用いた。

#### 5.2.4 開催日抽出

会議録システムの場合と同様の開催日抽出を行う。

#### 5.2.5 ファイルのメタ情報を用いた情報抽出

ファイル形式によっては、メタ情報に会議名や開催日に結びつく情報が格納されていることがある。具体的には、PDFファイルの場合は、タイトルやアウトラインに会議録の情報があることが多い。そのため、5.2.2で述べた会議名抽出が失敗したときに代わりとして利用できる。

### 6. データ収集結果

先に述べたデータ構造とデータ収集に基づいて、千葉県内の市町村議会の会議録を収集した結果を表3に示す。千葉県には地方議会が全部で54議会(37市16町1村)あり、会議録公開に関する内訳としては、そのうち49議会(同県内全議会91%)が会議録のWeb公開を行っている。公開形式でみると、会議録システムが35議会(同県内全議会65%)、ファイルによる公開が15議会(同県内全議会28%)である。ファイルによる公開をしている議会のファイル形式は、すべてPDFファイルである。作成したクローラによって、会議録システムによる公開をしている地方議会に関しては、市川市議会を除く34議会、ファイルによる公開をしている地方議会に関しては、勝浦市議会を除く14議会の会議録を収集できた。今回作成したクローラは大手4会議録システムに対応したのみである。市川市は独自の会議録システムを採用しているため、会議録データの収集ができなかった(図6に市川市議会の会議録公開画面を示す)。勝浦市は5.2.2で述べた会議名抽出に失敗したため、会議録データの収集ができなかった。具体的には、表組みされたセルの中に、会議録へのリンクや会議名に関する情報が格納されており、ページタイトル、リンクテキスト、Headingテキスト内に会議名に関する情報が含まれていなかった(図7に勝浦市の会議録公開画面を示す)。この形式に対応するためには、表組みを考慮した会議名抽出を行う必要がある。旭市議会と九十九里町議会、長柄町議会、大滝町議会は5.2.5で述べたファイルのメタ情報としてPDFのアウトラインに会議名が入っていたが、5.2.2の会議名抽出で会議名を抽出できなかったため利用しなかった。同様に茂原市議会はPDFのタイトルに会議名が入っていたが、5.2.2の会議名抽出で会議名を抽出できなかったため利用しなかった。

### 7. まとめ

本稿ではまず、地方議会会議録のWeb公開の2つの形式、会議録システムによる公開とファイルによる公開の形



図6 市川市の会議録公開画面



図7 勝浦市の会議録公開画面

について説明し、それぞれの公開形式に共通する収集における2つの問題を示した。そして、それらを解決するためのデータ構造の定義とデータ収集法について述べた。提案手法を千葉県内の地方議会に対して適用した結果、会議録システムによる公開をしている地方議会に関しては、35議会中34議会、ファイルによる公開をしている地方議会に関しては、15議会中14議会の議事録を収集することができた。今後の課題としては、提案手法の他の地方議会への適用が挙げられる。

\*2 <https://euske.github.io/pdfminer/>

表 3 千葉県内の市町村議会における  
会議録の公開形式とデータ収集結果

議会	公開形式	備考	収集結果
千葉市	会議録システム	DBsearch	○
木更津市	会議録システム	DBsearch	○
南房総市	会議録システム	DBsearch	○
山武市	会議録システム	DBsearch	○
長生村	会議録システム	DBsearch	○
白井市	会議録システム	DBsearch	○
銚子市	会議録システム	Discuss	○
館山市	会議録システム	Discuss	○
松戸市	会議録システム	Discuss	○
野田市	会議録システム	Discuss	○
成田市	会議録システム	Discuss	○
東金市	会議録システム	Discuss	○
習志野市	会議録システム	Discuss	○
柏市	会議録システム	Discuss	○
市原市	会議録システム	Discuss	○
流山市	会議録システム	Discuss	○
八千代市	会議録システム	Discuss	○
我孫子市	会議録システム	Discuss	○
鎌ヶ谷市	会議録システム	Discuss	○
君津市	会議録システム	Discuss	○
浦安市	会議録システム	Discuss	○
四街道市	会議録システム	Discuss	○
袖ヶ浦市	会議録システム	Discuss	○
印西市	会議録システム	Discuss	○
富里市	会議録システム	Discuss	○
匝瑳市	会議録システム	Discuss	○
香取市	会議録システム	Discuss	○
いすみ市	会議録システム	Discuss	○
大網白里市	会議録システム	Discuss	○
栄町	会議録システム	Sophia	○
船橋市	会議録システム	Voices	○
佐倉市	会議録システム	Voices	○
鴨川市	会議録システム	Voices	○
富津市	会議録システム	Voices	○
市川市	会議録システム	独自	x
茂原市	ファイル	PDF	○
旭市	ファイル	PDF	○
勝浦市	ファイル	PDF	x
八街市	ファイル	PDF	○
東庄町	ファイル	PDF	○
九十九里町	ファイル	PDF	○
横芝光町	ファイル	PDF	○
一宮町	ファイル	PDF	○
睦沢町	ファイル	PDF	○
長柄町	ファイル	PDF	○
長南町	ファイル	PDF	○
大多喜町	ファイル	PDF	○
御宿町	ファイル	PDF	○
鋸南町	ファイル	PDF	○
酒々井町	ファイル	PDF	○
神崎町	なし		x
多古町	なし	議会だよりあり	x
芝山町	なし	議会だよりあり	x
白子町	なし	議会だよりあり	x

- [2] 葦原史敏, 木村泰知, 荒木健治: 地方議会会議録における節単位による議員の要望抽出, 電子情報通信学会論文誌, Vol.J98-D, No.11, pp.1390-1401, 2015.
- [3] 高丸圭一, 内田ゆず, 乙武北斗, 木村泰知: 地方議会会議録コーパスにおけるオノマトペ-出現傾向と語義の分析-, 人工知能学会論文誌, Vol.30, No.1, pp.306-318, 2015.
- [4] Yasutomo Kimura et al., Creating Japanese Political Corpus from Local Assembly Minutes of 47 Prefectures, Coling 2016, The 12th Workshop on Asian Language Resources, pp.78-85, 2016.
- [5] 木村泰知, 小林暁雄, 坂地泰紀, 内田ゆず, 高丸圭一, 乙武北斗, 吉田光男, 川浦昭彦, 地方政治コーパス構築における従来の成果と現在の課題-政治・経済分野の応用研究に向けたパネルデータの構築-, 言語処理学会第 23 回年次大会, C1-5, 2017.
- [6] 齋藤誠, 大城卓, 菅原晃平, 永井隆広, 渋谷英潔, 木村泰知, 森辰則: 地方議会会議録の収集とコーパスの構築, 言語処理学会第 17 回年次大会, P2-21, 2011.

参考文献

- [1] 木村泰知, 関根聡: 主辞に基づく政治問題抽出手法, 人工知能学会論文誌, Vol.28, No.4, pp.370-378, 2013.