

自治体PDF文書に含まれる
図表解析における課題分析

2022073 小笠原 廉

2022222 高橋 孝輔

令和7年度提出

目次

第1章	はじめに	3
第2章	関連研究	5
2.1	レイアウト・構造情報を用いた文書理解と残る課題	5
2.2	OCR 依存手法と OCR 非依存手法の限界	5
2.3	Document VQA ベンチマークが示す構造理解の要求	6
2.4	日本語文書および公的文書における構造的課題	6
第3章	図表解析に向けたデータ分析	7
3.1	対象データの概要	7
3.2	予備調査による図表要素の観察	7
3.3	課題ラベル体系の設計	8
3.4	分析手法	9
第4章	課題ラベル体系の適用結果	10
4.1	図表ラベルの出現傾向	10
4.2	ラベルの複合出現パターン	10
第5章	課題ラベル体系の考察	12
5.1	ラベル体系の射程と限界	12
5.2	複合的な文書構造の特徴	13
5.2.1	図表表現における複合構造	13
5.2.2	表構造における複合構造	14
5.2.3	課題ラベル体系による複合構造の記述	14
第6章	おわりに	16
	参考文献	17

第1章 はじめに

デジタル社会の進展に伴い、自治体においても行政手続きのデジタル化や情報公開の電子化が進められている。とりわけ、住民への説明責任や行政の透明性確保の観点から、予算書や決算書、統計資料、計画書といった行政文書がPDF形式で広く公開されるようになった。これらのPDF文書には、数値情報を整理した表、施策や事業の関係性を補足する図、統計データを視覚化したグラフなど、多様な図表要素が含まれている。

一方で、PDF文書に含まれる図表は、人間による閲覧や理解を前提として設計されていることが多く、その構造や表現形式は多様である。例えば、セル結合や余白を用いた表現、複数の図表を同一ページ内に配置する構成、図表外への注釈の書き出しなど、多様な構造的および表現的工夫が見られる。このような構造的および表現的特徴は、人間にとっては理解しやすい場合がある一方で、機械による自動的な情報抽出や構造化を困難にする要因となり得る。図1.1は、本研究で対象とした自治体PDF文書の一部を例として示したものであり、複数の図表が同一ページ内に配置されている構成や、図表外に注釈が書き出されている表現など、文書構造の多様性が確認できる。

また近年は、大規模言語モデル（LLM）を用いた文書理解や図表解析に関する研究が進展している。しかし一方で、図表を含む文書に対しては、テキストのみの場合と比較して解析精度が低下することが報告されている [1] [2]。加えてこれらの研究は、個別の応用や性能評価に主眼を置いたものが多く、文書側に内在する構造的および表現的な特徴を、図表解析の観点から体系的に整理した研究は限られている。

もし、PDF文書に含まれる表や図を自動的に解析し、構造化されたデータとして扱うことが可能になれば、自治体文書に含まれる情報の再利用性は大きく向上すると考えられる。例えば、年度や部署をまたいだ数値の比較や、複数文書に分散した情報の横断的な分析が容易になるほか、質問応答、文書要約、検索といった下流タスクにおける前処理としての活用も期待できる。このような利点を実現するためには、PDF文書に含まれる図表の構造を適切に把握し、図表解析の観点からどのような構造的特徴が障壁となり得るのかを整理する必要がある。そのため、自治体文書の利活用を高度化する上で、図表解析は重要な基盤的課題の一つと位置づけられる。

以上の背景を踏まえ、本研究では実際の自治体PDF文書を対象として、図表解析において障壁となり得る文書構造や表現形式を人手による観察に基づいて分類し、整理する。具体的には、小樽市が公開する行政文書を対象とし、表、図およびテキストの外形的および構造的特徴に着目した課題ラベル体系を設計し、その出現傾向や複合的な出現状況を分析する。本研究の目的は、特定の解析手法やモデル性能を評価することではない。自治体PDF文書に内在する構造的特徴を記述可能な形で整理し、将来的な自動図表解析や文書解析手法を検討するための基礎的知見を提供することを目的とする。

以下、本論文の構成を示す。第2章では、図表を含む文書理解および図表解析に関する関連研究を概観する。第3章では、本研究で用いる対象データの概要を示すとともに、課題ラベル体系の設計および分析手法について説明する。第4章では、設計した課題ラベル体系を実データに適用した結果を示す。第5章では、第4章で得られた結果を踏まえ、課題ラベルの複合出現に着目しながら、自治体PDF文書における文書構造の特徴について考察する。

第2章 関連研究

本章では、図表を含む文書理解および図表解析に関する先行研究を概観する。特に、文書理解においてレイアウトや構造情報が重要であると認識されてきた経緯と、それにもかかわらず文書構造が解析上の障壁として残り続けている点に着目し、本研究が扱う課題が先行研究の議論の中でどのような点を対象としているかを整理する。

2.1 レイアウト・構造情報を用いた文書理解と残る課題

近年の文書理解研究では、テキスト内容のみならず文書中のレイアウトや視覚的配置を考慮する必要性が広く指摘されている。文書は単なる文字列の集合ではなく、段落構造、表や図の配置、また項目間の空間的關係といった構造的要素を通じて情報の意味が補完される媒体であると考えられている。

Nishida らは、表構造の理解においてセル間の関係性や配置情報が意味理解に重要であることを示し、構造情報を考慮した表理解モデルを提案した [3]。ただし、対象は主として Web 上の表であり、公的 PDF 文書のように図表配置や周辺テキストを含む文書全体への適用可能性については別途検討が必要であると考えられる。

Hsu らは、科学論文中の図を対象として、画像とキャプションの対応関係を学習する手法およびデータセット (SciCap) を提案し、図が文書理解において重要な役割を果たすことを示した [4]。一方で、生成されるキャプションの内容や図の理解に関しては、なお課題が残ることも報告されている。Yang らはこれを拡張し、本文中の言及情報や OCR トークンを統合することで性能向上を示したが、図の構造的な理解そのものについては依然として困難であることが示唆されている [5]。

Lee らは、文書画像から直接構造化テキストを生成する視覚言語モデルを提案し、レイアウト理解が文書理解性能に寄与することを示した [6]。ただし、日本語の公的 PDF 文書のような多様なレイアウトを持つ文書への適用については十分な検証が行われているとは言い難い。

これらの研究は、文書理解において構造情報が不可欠であることを示している一方で、多様なレイアウトや人間向けの表現を含む実文書における構造的困難さを体系的に整理する段階には至っていない。

2.2 OCR 依存手法と OCR 非依存手法の限界

多くの文書理解手法では、光学文字認識 (OCR) により抽出されたテキストを入力として解析を行う手法が採用されてきた。しかし、OCR 誤認識や文字列と文書構造との対応付けの不整合が解析精度の制約となることが指摘されている。Xu らは、学術文書中の表を対象として、本文との対応関係を活用したキャプション生成を試みたが、表内の数値構造や階層的配置の理解には課題が残ることを報告している [7]。

これに対し、Kim らは OCR を介さずに文書画像から直接構造化表現を生成する手法 (Donut) を提案し、OCR に起因する誤差伝播を回避しつつエンドツーエンドでの文書理解を可能にした [8]。

一方で同論文においても、複雑な表構造や図表配置を一貫して扱うことの難しさが示唆されている [8].

このように、OCR に依存するか否かに関わらず、文書中に内在する構造的特徴そのものが解析上の困難を生み出していることが示唆される。

2.3 Document VQA ベンチマークが示す構造理解の要求

文書理解能力を評価する枠組みとして、文書画像を対象とした質問応答タスク (Document VQA) が提案されている。Mathew らは、文書画像を対象とした質問応答データセット (DocVQA) を構築し、文書理解においてテキスト情報だけでなく配置や構造を踏まえた理解が必要であることを示した [9]。しかし、文書構造の把握を要する質問においてはモデル性能が低下する傾向が報告されている。

このような結果は、視覚情報と言語情報を統合したモデルであっても文書構造の把握が依然として困難であることを示しており、図表解析における主要な障壁が文書構造にあることを示唆している。

2.4 日本語文書および公的文書における構造的課題

日本語文書を対象とした研究においても文書構造の理解が重要な課題であることが報告されている。Onami らは、日本語 PDF 文書を対象とした質問応答データセット (JDocQA) を構築し、図表を含む質問に対してテキストのみの場合と比較して性能が大きく低下することを示した [1]。この結果は、日本語文書に特有のレイアウトや表現形式が文書理解の難易度を高めている可能性を示唆している。

Aida らは、有価証券報告書の表を対象として、視覚的レイアウト情報を導入することで性能向上を示したが、表構造の正確な解釈にはなお課題が残ることを指摘している [2]。

さらに、杉山らは、PDF から CSV への変換過程において生じる課題を分析し、セル結合や非スカラ値といった視覚的表現が構造理解を困難にしていることを報告している [10]。

これらの研究は、日本語文書や公的文書において文書構造の多様性が解析上の障壁として顕在化しやすいことを示しているが、その構造的特徴を文書横断的な観点から整理し、どのような構造が解析の障壁となり得るのかを体系的に示す枠組みは依然として十分に提示されていない。

第3章 図表解析に向けたデータ分析

本章では、自治体が公開する PDF 文書を対象として、図表解析の観点から生じる構造的および表現的な課題を整理し、その後の定量分析に向けた分析枠組みを構築する。本研究は、実際に大規模言語モデル（LLM）を用いた自動解析を行うものではなく、将来的な図表解析や文書変換を想定した際に障壁となり得る文書構造を、文書側の特徴として把握することを目的としている。

本章ではまず、分析対象とした自治体 PDF 文書の概要を示す。次に、本調査に先立って実施した予備調査について述べ、自治体文書に見られる図表要素の構造的特徴を整理する。その上で、予備調査の結果を踏まえて設計した課題ラベル体系および分析手法について説明する。

3.1 対象データの概要

本研究では、小樽市公式ウェブサイトにおける「市政の透明化」ページで公開されている PDF 文書を分析対象とした。これらの文書は、財政状況や施策の進捗を住民に説明することを目的として作成されており、表や図、補足説明が多く含まれている。そのため、自治体文書における図表解析の検討対象として適していると判断した。

2025年5月4日時点で取得可能であった PDF ファイルは計 929 件であり、その中から「財政」、「電源立地地域対策交付金事業」、および「事務執行状況説明書」に該当する文書を抽出した。複数年度の文書が存在する場合には最新年度の文書を優先し、最終的に 63 件の PDF 文書、合計 1,682 ページを対象とした。

文書ごとのページ数にはばらつきがあるため、4 ページ以上からなる文書については、各文書からランダムに 3 ページを抽出することで、分析対象の偏りを抑えつつ、作業負荷との両立を図った。その結果、本研究における最終的な分析対象は 139 ページとなった。

3.2 予備調査による図表要素の観察

本調査に先立ち、課題ラベル体系を設計するための予備調査を実施した。この予備調査は、課題の出現頻度を定量的に把握することを目的とするものではなく、自治体 PDF 文書に見られる図表要素の構造的および表現的特徴を把握し、分析観点を整理することを目的として行った。

予備調査では、「市政の透明化」関連 PDF および議会関連 PDF の中から、内容や構成の異なる 8 件の文書を対象として、ページ単位で表、図およびテキスト要素を概観した。観察は人手によって行い、図表解析の観点から特徴的と考えられる構造をメモとして記録した。この段階では、数値化や件数の集計は行っていない。

観察の結果、自治体 PDF 文書には、視覚的な配置や図形要素を多用した構成が見られた。例えば、1 ページ内に複数の表を配置して情報量を集約する構成や、セル結合や余白を用いて項目間の階層関係を示す表現が確認された。また、縦書きと横書きを併用した見出し配置や、表中の同一項目を矢印で示す表現も見られた。

表 3.1: 課題ラベル体系の分類および定義.

カテゴリ	課題ラベル	定義
表	階層構造	階層化された項目の位置づけを、セルの中での改行によって表記
	セル結合・非スカラ値	表のセル結合や、一つのセルに複数の項目や補足説明、空白セルを利用した階層構造の表現など
	複数の表	1 ページ中に形の異なる複数の表が存在
	重複項目の指示	表中の同じ内容を示す数値などについて、セル同士を矢印で繋げることで表現
図	論理構造の説明	複雑な論理構造や論理展開を、図や表、矢印を用いて整理
	外部への書き出し	図やグラフの中に入りきらない項目名を外部に書き出して表記
	視覚的データ表現	帯グラフなどを用いて、各項目における数値の大小関係を視覚的に表現
	グラフ読み取り困難	折れ線グラフにおける数値や変化をテキストとして変換できない表記
テキスト	段落表現	段落の先頭記号が複数の行に跨って表記されている状態
	多列レイアウト	ページ中に列数の異なる文章が混在して配置されている状態
	縦書き・横書き	縦書きの文章で構成されている、あるいは縦書きと横書きが混在している状態
	画像説明	写真とその説明文が対になって配置されている表現
	箇条書き・リスト	箇条書きとして整理可能な内容が、文章中で連続して羅列されている状態
	表的表現	改行やインデントではなく、表のような配置によって文章構造を表現している状態
	レイアウト調整	目次中でページ番号を縦に揃えるために、「…」などを用いて調整している表現

図やグラフにおいては、帯グラフを用いた数値関係の表現や、注釈や項目名を図外に書き出し、線によって対応付ける表現が確認された。このような構成では、数値情報や項目名が図形要素と分離して配置されており、図表要素を単純に抽出し、構造化することが困難となる場合がある。

このように、自治体 PDF 文書には多様な図表要素が存在し、それらが必ずしも単一の整理原理に基づいて構成されていないことが確認された。これらの観察結果を踏まえ、次節以降では、問題となり得る構造的特徴を体系的に整理する。

3.3 課題ラベル体系の設計

本研究では、自治体 PDF 文書に含まれる文書要素の外形的および構造的特徴を整理するため、表、図およびテキストの 3 カテゴリからなる課題ラベル体系を設計した。

本研究で用いる課題ラベルは、図表解析や PDF からの自動変換において問題となり得る構造的特徴を記述することを目的としており、文書の意味内容や主題の解釈を直接の対象とはしない枠組みとして位置づけている。このため、表、図およびテキストに現れる配置や形状、構成といった外形的および構造的特徴に着目してラベル設計を行った。

課題ラベルの設計にあたっては、予備調査において記録した観察メモをもとに、個別事例として現れていた特徴を整理し、統合し、内容の重複や粒度のばらつきを調整した上で分類を行った。その結果、課題ラベルは文書要素の性質に基づいて分類され、各カテゴリ内に複数のラベルを設定した。

また、本研究で設計した課題ラベルは、互いに排他的なものとして定義していない。これは、自治体 PDF 文書において、単一のページ内に複数の構造的特徴が同時に現れる場合があることが、予備調査の段階で確認されたためである。そのため、同一ページに対して複数の課題ラベルが同時に付与されることを想定した設計としている。

各課題ラベルの名称および定義は、表 3.1 に示す通りである。

3.4 分析手法

本研究では、設計した課題ラベル体系を用いて、ページ単位での人手アノテーションを行った。アノテーションは、日本語ネイティブであり、小樽市の市政に関する事前知識を持たない1名のアノテータが担当した。

各ページについて、表、図およびテキスト要素を観察し、該当する課題ラベルを付与した。ラベル付与後は、各ラベルの出現件数を集計するとともに、同一ページ内で複数のラベルが同時に出現する状況について定量的に分析した。

これにより、単一の課題の出現傾向だけでなく、課題が複合的に出現する文書構造の特徴を把握することを可能とした。

第4章 課題ラベル体系の適用結果

本章では、第3章で定義した課題ラベル体系を、小樽市が公開するPDF文書139ページに適用した結果について述べる。各課題ラベルの出現件数や、複数の課題ラベルが同一ページ内で出現する状況を集計し、自治体文書における課題要素の出現傾向を示す。

4.1 図表ラベルの出現傾向

表4.1は、小樽市が公開するPDF文書139ページにおける課題ラベルの出現件数を示している。最も出現件数が多かったのは「セル結合・非スカラ値」(85件)であり、次いで「複数の表」(36件)、「縦書き・横書き」(22件)、「視覚的データ表現」(20件)、「多列レイアウト」(19件)、「グラフ読み取り困難」(18件)などが続く。

4.2 ラベルの複合出現パターン

課題ラベルは1ページに1件だけではなく、複数のラベルが同時に出現する場合もあった。表4.2は同一ページにおける課題ラベルの出現数の分布を示している。PDF139ページ中、132ページで1件以上のラベルが出現した。また、2件以上のラベルが出現したページ数は78ページであり、今回調査したページの過半数に達した。

表4.3および表4.4は複数の課題ラベルが同一ページ内で同時に出現する事例について、その組み合わせと件数を示している。表4.3は、2件以上の課題ラベルが同一ページ内で出現した事例において、同時に出現することが多かった2ラベルの組み合わせを示している。例えば、「表：セル結合・非スカラ値」と「表：複数の表」は26件で同時出現しており、当該2ラベルのみで構成される場合のほか、他のラベルを加えた形で出現する場合も含まれている。

表4.4は、3件以上の課題ラベルが同一ページ内で出現した事例におけるラベルの組み合わせの件数上位を示している。

表 4.1: 課題ラベルの出現件数.

カテゴリ	課題ラベル	件数
表	階層構造	15
	セル結合・非スカラ値	85
	複数の表	36
	重複項目の指示	3
図	論理構造の説明	8
	外部への書き出し	16
	視覚的データ表現	20
	グラフ読み取り困難	18
テキスト	段落表現	3
	多列レイアウト	19
	縦書き・横書き	22
	画像説明	1
	箇条書き・リスト	13
	表的表現	8
	レイアウト調整	4

表 4.2: 同一ページ内における課題ラベル同時出現数の分布.

同時出現ラベル数	ページ数
1	54
2	40
3	23
4	10
5	3
6	1
7	1

表 4.3: 2 ラベルの複合出現パターン (件数上位).

複合ラベルの組み合わせ	件数
表：セル結合・非スカラ値 × 表：複数の表	26
表：セル結合・非スカラ値 × テキスト：縦書き・横書き	18
表：階層構造 × 表：セル結合・非スカラ値	15
図：グラフ読み取り困難 × 図：視覚的データ表現	14
図：グラフ読み取り困難 × 図：外部への書き出し	13
図：外部への書き出し × 図：視覚的データ表現	11

表 4.4: 3 ラベルの複合出現パターン (件数上位).

複合ラベルの組み合わせ	件数
図：グラフ読み取り困難 × 図：外部への書き出し × 図：視覚的データ表現	11
表：セル結合・非スカラ値 × 表：複数の表 × テキスト：縦書き・横書き	6

第5章 課題ラベル体系の考察

本章では、第4章で示した課題ラベルの適用結果を踏まえ、本研究において定義した課題ラベル体系が、どのような前提のもとで分析結果を与えているかを整理するとともに、それらの結果から読み取れる自治体 PDF 文書の文書構造上の特徴について考察する。具体的には、前者については5.1節においてラベル体系の射程と限界を整理し、後者については5.2節において課題ラベルの複合出現に着目した文書構造の特徴を検討する。本章の目的は、第4章で得られた分析結果をどのように理解し位置づけるべきかを明確にすることであり、具体的な解析手法の設計や解決策の提示ではなく、結果の解釈と研究全体における意味づけに主眼を置く。なお、本章に示す図は、いずれも小樽市が公開する行政文書を原資料とし、課題ラベルの複合出現に対応する文書構造の特徴を説明する目的で、当該ページの一部を抜粋して示したものである。

5.1 ラベル体系の射程と限界

本研究の結果を適切に解釈するためには、本研究で用いた課題ラベル体系が何を対象とし、どの範囲までを記述するものかを明確にしておく必要がある。

まず、本研究で提示した分析結果は、自治体 PDF 文書に含まれる図表やテキストの外形的および構造的特徴を人手により整理し、分類したものであり、文書内容の意味理解や数値の正確性、政策的妥当性を評価するものではない。また、本研究は実際に大規模言語モデル（LLM）を用いた図表解析を行ったものではなく、将来的に自動解析を行う際に障壁となり得る文書構造を事前に整理することを目的としている。したがって、本研究の結果は特定の解析手法やモデル性能を直接的に示すものではない。

次に、課題ラベルの付与は、人手による観察と判断に基づいて実施されており、主観的判断が完全に排除されるわけではない。しかし一方で、本研究では各課題ラベルについて名称のみを用いるのではなく、具体的な説明を定義し、その意味に基づいてラベル付与を行っている。判断基準を明示することで、ラベルの解釈が恣意的に揺れることを抑制し、分類結果に一定の一貫性を持たせるよう設計している。

さらに、対象データについては、小樽市が公開する行政文書全体を網羅的に分析したものではなく、計算量と分析の実現可能性を考慮し、文書およびページ単位でランダムサンプリングを行った。このため、本研究の結果をもって自治体文書一般の傾向を直接的に断定することはできない。しかし、対象集合内では特定の形式に偏らないよう配慮した抽出を行っており、同一自治体内の文書であっても、表や図、テキストの表現形式に多様性が存在することを確認している。

以上より、本研究で提示した課題ラベル体系は、自治体 PDF 文書に含まれる図表やテキストの構造的特徴を記述するための探索的かつ設計的な枠組みであり、本研究の範囲内で得られた観測結果を解釈するための前提として有効に機能するものである。

5.2 複合的な文書構造の特徴

本節では、第 4 章で示した課題ラベルの複合出現に着目し、自治体 PDF 文書における文書構造の重層性について考察する。第 4 章の分析結果では、半数を超えるページにおいて 2 件以上の課題ラベルが確認され、また特定のラベルの組み合わせが繰り返し観測された。このことから、自治体 PDF 文書における表現は単一の構造的特徴のみを前提として把握することが難しく、複数の構造的要素が同時に存在する文書表現の実態を反映していると解釈できる。本節では、まず図表表現における複合構造を取り上げ、続いて表構造における複合構造を検討する。

5.2.1 図表表現における複合構造

第 4 章の分析結果より、とくに図に関する「視覚的データ表現」「グラフ読み取り困難」「外部への書き出し」の組み合わせが高い頻度で観測された。これらは、(1) 数値の大小関係や項目の該当範囲の図示において、テキストや矢印の大小による記述や、線による囲い込みを使用した記述といった表現方法に加えて、(2) 帯グラフや円グラフ等における面積の小さい項目に対し、その図内に収まらない数値や項目名等のテキストをグラフ外周部に書き出すといった表現方法が併存する構成である。

このような構成のうち、(1) の視覚的データ表現では、数値の大小関係や項目の該当範囲といった情報が、テキストや矢印の配置、線による囲い込みなどの視覚的手段によって示される。これにより、数値や項目は図形的要素と一体化した形で配置され、図全体の中で関係づけられた形で構成される構造を持つ。一方、(2) の外部への書き出しを伴う構成では、帯グラフ等において面積の小さい項目に対応する数値や項目名が、図内部には配置されず、外周部に書き出される。この際、外部に配置されたテキストは線によって該当するグラフ要素と結び付けられており、グラフ要素とテキストとの対応関係自体は維持されている。すなわち、同一の意味関係が、図形内部への直接配置と、線による参照関係という二つの異なる配置手段によって表現されている。図 5.1 に示すように、一部の項目では数値や項目名が図形内部に直接配置され、一方で、面積の小さい項目については、図外に書き出されたテキストが線によって対応付けられている。これら二つの表現が同一の図内で併存する場合、項目によって情報の配置規則が異なる構造となる。その結果、図表は単一の表現様式に基づいて構成されるのではなく、視覚的強調と空間的制約への対応という異なる要請に基づく表現が重なり合った構造を持つことになる。

以上より、本研究においてこれらの課題ラベルが複合して出現したことは、自治体 PDF 文書における文書構造の一つの特徴を示している。具体的には、情報を視覚的に強調するための表現と、空間的制約に対応するための配置上の工夫といった、異なる設計原理に基づく表現が、同一ページ内で同時に適用されている構造が確認された。

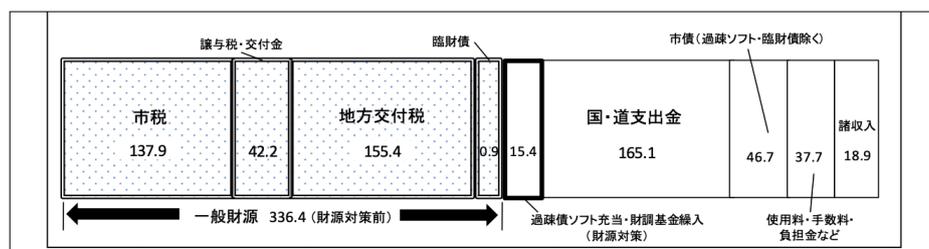


図 5.1: 図形内部への直接配置と外周部への書き出しが併存するグラフ表現の例。

5.2.2 表構造における複合構造

第 4 章の分析結果より、表に関する課題ラベルとして「セル結合・非スカラ値」「複数の表」「縦書き・横書き」が同一ページ内で複合して出現する事例が多く確認された。これらは、表構造において単一の整理原理が用いられているのではなく、複数の異なる構造的工夫が同時に適用されている状況を示している。

まず、「セル結合・非スカラ値」が付与された表では、項目名や数値、補足説明など、複数の意味単位が同一セル内または結合セル内にまとめて配置される構造が観察された。このような表では、情報の整理単位がセル境界と必ずしも一致せず、セルの形状や配置によって意味的なまとまりが表現されている。その結果、表の構造は単純な行列構造というよりも複数の意味単位を内包した複合的な配置構造を持つことになる。

次に、「複数の表」が付与されたページでは、同一の事業や施策に関する複数の表が、1 ページ内に集約して配置されている事例が多く観察された。図 5.2 上段の例では、目的や項目構成の異なる複数の表が、同一ページ内に並置されている。また、図 5.2 下段の例では、現行収支や改善目標、対策後の収支不足額といった異なる段階の情報を示す複数の表が縦方向に配置され、項目名セルや数値セルの幅を揃えることで、ページ全体として一体的な構成が取られていた。このように、「複数の表」は、関連する情報を一括して提示するための配置上の工夫として用いられており、結果として、粒度や構造の異なる表が同一ページ内に共存する構成を生じさせている。

さらに、「縦書き・横書き」が付与された表では、縦方向にセル結合されたカテゴリセルや見出しセルの形状に対応して、結果的に縦書きが採用されている事例が多く確認された。すなわち、縦書きは独立した表現選択というよりも、セル結合によって生じた表構造と連動して現れていると解釈できる。この場合、横書きを前提とした数値や説明文と、縦書きで配置されたカテゴリ表記とが同一表内に混在する構成となる。

なお、本節で例示した図 5.2 は、複数の表配置に加え、セル結合や記述方向の混在といった他の表構造上の特徴も併せて含んでおり、課題ラベルの複合出現を視覚的に確認できる事例である。

以上のように、表構造に関する課題ラベルの複合出現は、セル内部への情報集約、異なる目的を持つ表の集約配置、および記述方向の混在といった複数の構造的工夫が同一ページ内で同時に適用されている状況を反映している。本研究においてこれらのラベルが複合して観測されたことは、自治体 PDF 文書における表構造が単一の整理規則に基づくものではなく、複数の配置原理や表現方針が重ね合わされた構造を持つことを示している。

5.2.3 課題ラベル体系による複合構造の記述

自治体 PDF 文書における課題ラベルの複合出現は、解析上の問題が単に重なっている状態というよりも、人間の参照や補完を前提として設計および整形されてきた文書構造が、結果として複合的な特徴を持つことを示していると位置づけられる。本研究で整理した課題ラベル体系は、このような文書構造の重層性を記述し、どのような特徴が同時に現れやすいかを把握するための枠組みとして、本研究の目的と範囲において機能している。

(2) 工事等				(3) 業務																																																																							
イ 建設改良費の概況【税込】				イ 業務量																																																																							
<table border="1"> <thead> <tr> <th>名称</th> <th>内容</th> <th>金額</th> <th>施工年月日</th> <th>資産科目</th> </tr> </thead> <tbody> <tr> <td rowspan="5">水 道 メ ー タ ー</td> <td>電子検漏式メーター φ13mm 5個</td> <td rowspan="5">347,270</td> <td rowspan="5"></td> <td rowspan="5">機械及び装置 水道メーター</td> </tr> <tr> <td>φ20mm 3個</td> </tr> <tr> <td>φ25mm 4個</td> </tr> <tr> <td>φ40mm 3個</td> </tr> <tr> <td>φ70mm 1個</td> </tr> <tr> <td>計 16個</td> <td></td> <td></td> <td></td> </tr> <tr> <td colspan="2">営業設備費計</td> <td>347,270</td> <td></td> <td></td> </tr> <tr> <td colspan="2">合 計</td> <td>347,270</td> <td></td> <td></td> </tr> </tbody> </table>				名称	内容	金額	施工年月日	資産科目	水 道 メ ー タ ー	電子検漏式メーター φ13mm 5個	347,270		機械及び装置 水道メーター	φ20mm 3個	φ25mm 4個	φ40mm 3個	φ70mm 1個	計 16個				営業設備費計		347,270			合 計		347,270			<table border="1"> <thead> <tr> <th>事項</th> <th>令和5年度</th> <th>令和4年度</th> <th>比較</th> <th>備 考</th> </tr> </thead> <tbody> <tr> <td>給水業者数(件)</td> <td>68</td> <td>64</td> <td>▲4</td> <td></td> </tr> <tr> <td>年間給水量(m³) A</td> <td>319,642</td> <td>346,119</td> <td>△26,477</td> <td></td> </tr> <tr> <td>年間取水量(m³) B</td> <td>178,538</td> <td>177,719</td> <td>▲819</td> <td></td> </tr> <tr> <td>有収率(%) B/A</td> <td>55.9</td> <td>51.3</td> <td>▲4.6</td> <td></td> </tr> <tr> <td>一日最大給水量(m³)</td> <td>1,569</td> <td>1,387</td> <td>▲182</td> <td>令和5年7月26日</td> </tr> <tr> <td>一日最大取水量(m³)</td> <td>347</td> <td>348</td> <td>▲1</td> <td>令和5年1月1日</td> </tr> <tr> <td>一日平均給水量(m³)</td> <td>873</td> <td>948</td> <td>▲75</td> <td></td> </tr> </tbody> </table>				事項	令和5年度	令和4年度	比較	備 考	給水業者数(件)	68	64	▲4		年間給水量(m ³) A	319,642	346,119	△26,477		年間取水量(m ³) B	178,538	177,719	▲819		有収率(%) B/A	55.9	51.3	▲4.6		一日最大給水量(m ³)	1,569	1,387	▲182	令和5年7月26日	一日最大取水量(m ³)	347	348	▲1	令和5年1月1日	一日平均給水量(m ³)	873	948	▲75	
名称	内容	金額	施工年月日	資産科目																																																																							
水 道 メ ー タ ー	電子検漏式メーター φ13mm 5個	347,270		機械及び装置 水道メーター																																																																							
	φ20mm 3個																																																																										
	φ25mm 4個																																																																										
	φ40mm 3個																																																																										
	φ70mm 1個																																																																										
計 16個																																																																											
営業設備費計		347,270																																																																									
合 計		347,270																																																																									
事項	令和5年度	令和4年度	比較	備 考																																																																							
給水業者数(件)	68	64	▲4																																																																								
年間給水量(m ³) A	319,642	346,119	△26,477																																																																								
年間取水量(m ³) B	178,538	177,719	▲819																																																																								
有収率(%) B/A	55.9	51.3	▲4.6																																																																								
一日最大給水量(m ³)	1,569	1,387	▲182	令和5年7月26日																																																																							
一日最大取水量(m ³)	347	348	▲1	令和5年1月1日																																																																							
一日平均給水量(m ³)	873	948	▲75																																																																								
<p>ロ 保存工事の概況</p> <p>(イ)水道メーター取替及び修理 既設水道メーターの有効期限の経過又は故障等により、取替・修理した水道メーターの個数は次のとおりである。</p> <p>水道メーター取替個数 6個 故障等によるもの 2個 水道メーター修理個数 -個 故障等によるもの -個</p>				<p>(ウ)口探別水道メーター設置個数</p> <table border="1"> <thead> <tr> <th>口径</th> <th>令和5年度</th> <th>令和4年度</th> <th>備 考</th> </tr> </thead> <tbody> <tr> <td>φ1.3mm</td> <td>24</td> <td>22</td> <td rowspan="6">臨時メーターを除く</td> </tr> <tr> <td>φ2.0mm</td> <td>16</td> <td>14</td> </tr> <tr> <td>φ2.5mm</td> <td>17</td> <td>15</td> </tr> <tr> <td>φ4.0mm</td> <td>22</td> <td>22</td> </tr> <tr> <td>φ5.0mm</td> <td>16</td> <td>15</td> </tr> <tr> <td>φ1.0mm</td> <td>1</td> <td>1</td> </tr> <tr> <td>計</td> <td>96</td> <td>89</td> <td></td> </tr> </tbody> </table>				口径	令和5年度	令和4年度	備 考	φ1.3mm	24	22	臨時メーターを除く	φ2.0mm	16	14	φ2.5mm	17	15	φ4.0mm	22	22	φ5.0mm	16	15	φ1.0mm	1	1	計	96	89																																										
口径	令和5年度	令和4年度	備 考																																																																								
φ1.3mm	24	22	臨時メーターを除く																																																																								
φ2.0mm	16	14																																																																									
φ2.5mm	17	15																																																																									
φ4.0mm	22	22																																																																									
φ5.0mm	16	15																																																																									
φ1.0mm	1	1																																																																									
計	96	89																																																																									
<p>(水道メーター取替及び修理個数口探別内訳表)</p> <table border="1"> <thead> <tr> <th>口径</th> <th>取替個数</th> <th>修理個数</th> <th>備 考</th> </tr> </thead> <tbody> <tr> <td>φ1.3mm</td> <td>2</td> <td>-</td> <td></td> </tr> <tr> <td>φ2.0mm</td> <td>2</td> <td>-</td> <td></td> </tr> <tr> <td>φ2.5mm</td> <td>2</td> <td>-</td> <td></td> </tr> <tr> <td>φ4.0mm</td> <td>2</td> <td>-</td> <td></td> </tr> <tr> <td>φ5.0mm</td> <td>-</td> <td>-</td> <td></td> </tr> <tr> <td>φ1.0mm</td> <td>-</td> <td>-</td> <td></td> </tr> <tr> <td>計</td> <td>8</td> <td>-</td> <td></td> </tr> </tbody> </table>				口径	取替個数	修理個数	備 考	φ1.3mm	2	-		φ2.0mm	2	-		φ2.5mm	2	-		φ4.0mm	2	-		φ5.0mm	-	-		φ1.0mm	-	-		計	8	-		<p>(ハ)水道料金徴収及び収納状況【税込】</p> <table border="1"> <thead> <tr> <th rowspan="2">区 分</th> <th colspan="2">件数</th> <th colspan="2">金額</th> <th rowspan="2">収納率 B/A</th> </tr> <tr> <th>A</th> <th>B</th> <th>A</th> <th>B</th> </tr> </thead> <tbody> <tr> <td>納付額</td> <td>481件</td> <td>481件</td> <td>31,938,141円</td> <td>33,480,328円</td> <td>100.0%</td> </tr> <tr> <td>口座振替</td> <td>612件</td> <td>605件</td> <td>45,564,589円</td> <td>44,913,268円</td> <td>98.6%</td> </tr> <tr> <td>計</td> <td>1,093件</td> <td>1,086件</td> <td>77,502,730円</td> <td>78,393,596円</td> <td>101.1%</td> </tr> </tbody> </table>				区 分	件数		金額		収納率 B/A	A	B	A	B	納付額	481件	481件	31,938,141円	33,480,328円	100.0%	口座振替	612件	605件	45,564,589円	44,913,268円	98.6%	計	1,093件	1,086件	77,502,730円	78,393,596円	101.1%								
口径	取替個数	修理個数	備 考																																																																								
φ1.3mm	2	-																																																																									
φ2.0mm	2	-																																																																									
φ2.5mm	2	-																																																																									
φ4.0mm	2	-																																																																									
φ5.0mm	-	-																																																																									
φ1.0mm	-	-																																																																									
計	8	-																																																																									
区 分	件数		金額		収納率 B/A																																																																						
	A	B	A	B																																																																							
納付額	481件	481件	31,938,141円	33,480,328円	100.0%																																																																						
口座振替	612件	605件	45,564,589円	44,913,268円	98.6%																																																																						
計	1,093件	1,086件	77,502,730円	78,393,596円	101.1%																																																																						
<p>ロ 事業収入に関する事項</p> <p>(イ)収入状況【税込】</p> <table border="1"> <thead> <tr> <th rowspan="2">科目</th> <th colspan="2">令和5年度</th> <th colspan="2">令和4年度</th> <th rowspan="2">比較</th> </tr> <tr> <th>収 益</th> <th>構成比</th> <th>収 益</th> <th>構成比</th> </tr> </thead> <tbody> <tr> <td>事業収益</td> <td>70,687,800</td> <td>51.9%</td> <td>71,458,890</td> <td>60.5%</td> <td>▲771,090</td> </tr> <tr> <td>給水収益</td> <td>70,457,400</td> <td>51.7%</td> <td>71,107,040</td> <td>60.2%</td> <td>▲649,640</td> </tr> <tr> <td>その他営業収益</td> <td>230,400</td> <td>0.3%</td> <td>351,850</td> <td>0.3%</td> <td>▲121,450</td> </tr> <tr> <td>営業外収益</td> <td>65,615,217</td> <td>48.1%</td> <td>66,647,716</td> <td>39.5%</td> <td>▲10,032,500</td> </tr> <tr> <td>計</td> <td>136,533,017</td> <td>100.0%</td> <td>118,106,606</td> <td>100.0%</td> <td>▲18,426,411</td> </tr> </tbody> </table>				科目	令和5年度		令和4年度		比較	収 益	構成比	収 益	構成比	事業収益	70,687,800	51.9%	71,458,890	60.5%	▲771,090	給水収益	70,457,400	51.7%	71,107,040	60.2%	▲649,640	その他営業収益	230,400	0.3%	351,850	0.3%	▲121,450	営業外収益	65,615,217	48.1%	66,647,716	39.5%	▲10,032,500	計	136,533,017	100.0%	118,106,606	100.0%	▲18,426,411																																
科目	令和5年度		令和4年度		比較																																																																						
	収 益	構成比	収 益	構成比																																																																							
事業収益	70,687,800	51.9%	71,458,890	60.5%	▲771,090																																																																						
給水収益	70,457,400	51.7%	71,107,040	60.2%	▲649,640																																																																						
その他営業収益	230,400	0.3%	351,850	0.3%	▲121,450																																																																						
営業外収益	65,615,217	48.1%	66,647,716	39.5%	▲10,032,500																																																																						
計	136,533,017	100.0%	118,106,606	100.0%	▲18,426,411																																																																						

-12-

財政再建推進プラン実施計画収支試算 (一般財源ベース)						
単位:億円						
	H18	H19	H20	H21	累計	
現行収支	単年度収支不足額 A	0.0	▲439	▲407	▲379	▲1225
繰越収支	繰越収支不足額 H17年度決算見込 ▲20億円	▲20.0	▲639	▲1046	▲1425	▲1425
改 善 目 標	繰出削減対策 B		36.1	32.3	32.8	101.2
	1. 人件費の抑制		22.9	21.7	24.6	69.2
	(1) 退職者の不補充		1.1	0.8	2.0	3.9
	(2) 職員給与等の削減 H19以降10% (地域間格差相当4.8%+独自削減5%)		14.3	14.1	13.9	42.3
	(3) その他 退職手当償の導入など		7.5	6.8	8.7	23.0
	2. 事業の見直し		13.2	10.6	8.2	32.0
	(1) 管理経費の圧縮		0.6	0.6	0.9	2.1
	(2) 特別会計 企業会計の収支改善 (繰出金の削減)		6.9	5.5	5.0	17.4
	(3) 市債の償還		0.7	0.7	0.7	2.1
	(4) その他		4.8	3.6	1.6	10.0
3. 繰入増の取組 C		0.9	0.9	1.2	3.0	
(1) 入湯税課税免除の見直し		0.4	0.4	0.4	1.2	
(2) 使用料・手数料等の改定					0.3	
(3) その他		0.5	0.5	0.5	1.5	
財源対策 D				4.0	4.0	
小 計 E=B+C+D	0.0	37.0	38.2	38.0	108.2	
対策後単年度収支不足額 A+E	0.0	▲69	▲7.5	0.1	▲14.3	
対策後繰越収支不足額	▲20.0	▲259	▲344	▲343	▲343	

図 5.2: 同一ページ内に複数の表が配置されている例。上段は、異なる業務項目や管理観点に属する表が同一ページ内に集約して配置された事例を示す。下段は、同一テーマに関する表が段階的に縦方向へ配置された事例を示す。

第6章 おわりに

本研究では、自治体が公開するPDF文書を対象として、図表解析において障壁となり得る文書構造や記述形式を明らかにすることを目的とし、人手による観察に基づいて文書構造上の課題を分類し、整理した。具体的には、小樽市が公開する行政文書を分析対象とし、表、図およびテキストの外形的および構造的特徴に着目して、計15種類の課題ラベルからなる分類体系を設計した。

設計した課題ラベル体系を139ページの実文書に適用した結果、78ページにおいて2件以上の課題ラベルが同時に付与される状況が確認された。この結果は、自治体PDF文書に含まれる図表やレイアウト上の課題が単一の構造的特徴として現れるのではなく、複数の異なる構造的要因が同一ページ内に併存していることを示している。すなわち、図表解析においては個別の課題を独立に扱うだけでは不十分であり、複数の構造的特徴が同時に存在する文書構造を前提として整理および解析する必要があることが明らかとなった。

本研究で示した課題ラベルは、特定の大規模言語モデルや解析手法の性能を評価することを目的としたものではなく、文書側に内在する構造的および表現的特徴を記述可能な形で整理するための枠組みである。そのため、本研究の結果は自治体PDF文書における図表解析が直面する課題を文書構造の観点から把握するための基礎的知見として位置づけられる。

今後は、本研究で整理した課題ラベル体系を基に大規模言語モデルを用いた図表解析や文書理解を実施し、どのような文書構造がどの処理段階で解析上の障壁となるのかを検証することが求められる。また、本研究では小樽市の行政文書を対象としたが、他自治体の文書や手書き文字を含む資料などへ対象を拡張することで、課題ラベル体系の汎用性や限界を検討する余地がある。

以上より、本研究は、自治体PDF文書に含まれる図表およびレイアウト上の構造的特徴を整理するための基盤的な整理枠組みを提示したものであり、今後の自動図表解析手法やアノテーション基盤の設計に向けた出発点の一つとして意義を有する。

参考文献

- [1] Eri Onami, Shuhei Kurita, Taiki Miyanishi, and Taro Watanabe. Jdocqa: Japanese document question answering dataset for generative language models, 2024.
- [2] Hayato Aida, Kosuke Takahashi, and Takahiro Omi. Enhancing large vision-language models with layout modality for table question answering on japanese annual securities reports, 2025.
- [3] Kyosuke Nishida, Kugatsu Sadamitsu, Ryuichiro Higashinaka, and Yoshihiro Matsuo. Understanding the semantic structures of tables with a hybrid deep neural network architecture. Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 31, No. 1, Feb. 2017.
- [4] Ting-Yao Hsu, C. Lee Giles, and Ting-Hao 'Kenneth' Huang. Scicap: Generating captions for scientific figures, 2021.
- [5] Zhishen Yang, Raj Dabre, Hideki Tanaka, and Naoaki Okazaki. Scicap+: A knowledge augmented dataset to study the challenges of scientific figure captioning, 2023.
- [6] Kenton Lee, Mandar Joshi, Iulia Turc, Hexiang Hu, Fangyu Liu, Julian Eisenschlos, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. Pix2struct: Screenshot parsing as pretraining for visual language understanding, 2023.
- [7] Junjie H. Xu, Kohei Shinden, and Makoto P. Kato. Table caption generation in scholarly documents leveraging pre-trained language models, 2021.
- [8] Geewook Kim, Teakgyu Hong, Moonbin Yim, Jeongyeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park. Ocr-free document understanding transformer, 2022.
- [9] Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawahar. Docvqa: A dataset for vqa on document images, 2021.
- [10] 杉山陽菜乃, 阿部瑞稀, 中村彩乃, 前多陸玖, 坂口遥哉, 佐藤栄作, 木村泰知, 小林暁雄, 大友将宏, 石原潤一, 桂樹哲雄, 川村隆浩. 農林業基準技術に含まれる表を対象とした pdf から csv へ変換する際の課題分析. 言語処理学会第 31 回年次大会 (NLP2025), pp. 3016–3020, 3 2025.