

ドキュメント画像から文書検索を行うための XML 定義の提案

The proposal for XML definition to retrieve documents from the document image

松本 馨^{*1}
Kaoru MATSUMOTO

櫻田 武嗣, 中川 正樹^{*2}
Takeshi SAKURADA, Masaki NAKAGAWA

^{*1} 学校法人産業能率大学 総合研究所
Research Institute, The Sanno Institute of Management

^{*2} 東京農工大学工学部
Tokyo University of Agriculture and Technology

This paper describes XML definition to retrieve documents from document images. It has been a problem that document management systems which use Optical Character Recognition (OCR) often pose some system requirements and do not avail specifications of the data, although many products have already been invested and used. Moreover, they assume entire translation from document images to codes so that verification of OCR has been a heavy burden for operators. This paper propose an XML definition, which provides high compatibility of the data and keeps multiple OCR candidates to dispense with the verification labor of OCR.

1. はじめに

近年, 共通の仕様を定めることでデータの相互運用性を高め, 一般に普及させるための XML 形式が多く見られるようになってきた。OCR を使用したドキュメント管理システムは, 既に多くの製品が存在し, 運用されている。しかし, これらは動作する環境が限られていることや, データの詳細仕様が明らかでないこと, そのデータを利用した新たなアプリケーション開発が第三者には困難であることが問題である。また, 画像から文字列への全置き換えを前提としているものが多く, 認識結果を確認・訂正する作業が作業者にとって大きな負担となっている。そこで, ドキュメント画像の OCR 結果を仕様の明らかな XML 形式で保存し, それを文書検索に利用するための XML 定義の提案を行う。

2. 背景

電子文書の分野では Adobe 社の pdf 形式が最も有名で, よく用いられているのはいうまでもない。これは, 多くの OS 環境で動作し仕様も公開されており[1], 閲覧ソフトが無料で配布されていることから, 広く用いられている。しかし, pdf 形式は, 文字コードや文字フォント, 図表などの書式を保存し, 他環境でも元データに近い体裁で読めるようにすることを目的としており, 画像として取り込まれたデータの保存にはあまり向いていない。

pdf 形式は, 透明テキストとして画像に文字情報を付加することが可能である。しかし, これは使用する画像形式が固定的で, ほとんどユーザーが選択できず, 既に存在する多様な画像ファイル形式に対応できない。また, 仕様が公開されているとはいえ, pdf 形式自体の利用が難しく, これに対応した検索機能を持つソフトウェアは多くない。一部のファイル検索用ソフトウェアや google をはじめとする Web 検索エンジンがサポートしている程度で, 個人で手軽に利用できるものとは言い難い状況である。加えて, ユーザー主導でないバージョンアップにより, ソフトウェアのバージョンに依存した問題が発生する可能性も残っている。他にも文書画像管理のためのツールはいくつか存在するが, これらも同様の問題を抱えている。

このため, 既存の画像形式 (jpeg や png など) から独立した付加情報として, 仕様の明らかな XML 形式を用いたメタデータを

文字情報として付加することが有効であると考えた。これは, 既に存在する画像閲覧/管理ソフトウェアが利用可能であることや, 今後, 新しい画像形式が登場しても, 互換性が維持されたり, 画像変換用ソフトウェアの登場が見込まれるからである。XML により付加されるメタデータと画像ファイルが分離していることで, 画像形式の変換にも容易に対応可能となり, 付加情報のハンドリングも容易に行うことが可能となる。

3. 目的

3.1 文字認識システムの概要

文字認識システムは, 大別するとペンタブレット上で筆記情報をリアルタイムに取得するオンライン文字認識と, 紙や画像として残された筆記情報 (画像情報) をもとに文字認識処理を行うオフライン文字認識の 2 種類がある。本研究は, オフライン文字認識を用いるものである。

オフライン文字認識システムは, 大きく分けると前処理と文字認識処理, 後処理の 3 つの処理で構成される (図 1)。このうち, 前処理では, 文字認識エンジンにかけ前の画像処理 (グレースケール化, 2 値化処理など), 行切り出し, 文字切り出し, 正規化処理を行う。文字認識処理では, 切り出した文字画像に対して統計的手法による文字認識処理を行い, 辞書に登録された文字情報との距離計算を行い, 候補文字列を算出する。後処理では, 文脈による処理を行い, 候補文字列とその距離値から一部の認識結果の入れ替えを行い, より文章として妥当な文字列になるように認識結果の修正を行う。多くの OCR システムでは, この距離値による候補文字の確からしさが一定値以下の場合に, 誤認識の可能性が高いパターンとして, 作業者に認識結果の確認を促し, 目視による修正を行っている。

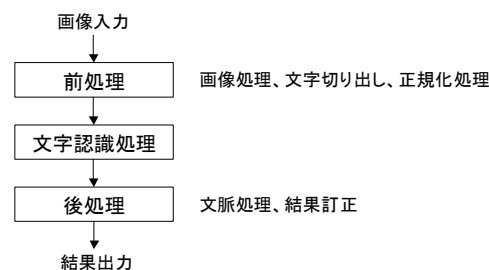


図 1 文字認識システムの構成

連絡先: 松本 馨, 学校法人産業能率大学 総合研究所,
〒158-8630 東京都世田谷区等々力 6 丁目 39 番 15 号,
MATSUMOTO_Kaoru@hj.sanno.ac.jp

3.2 文字認識システムの問題点

文字認識システムを利用するにあたって問題となる点として下記が挙げられる。

(a) 画像の2値化処理

隣接する文字がつながらなく、かつ、1つの文字が細切れにならないようにグレースケール画像を2値画像に変換すること

(b) レイアウト認識

文字のレイアウト(縦/横書き、段組など)を認識し、1つの文章としてつながるようにすること

(c) 文字切り出し

文字画像を1文字ずつ切り出して、文字認識エンジンに渡すこと

(d) 認識結果の決定と修正

候補文字列(認識スコアの上位1位~10位程度を使用)の中から妥当なものを選択すること

この中で(b)は非常に難しい問題であるが、認識結果を文章として扱うのではなく、画像上の座標に文字を対応付ける場合、大きな問題ではなくなる。(a),(c)も重要な問題であるが、利用者が実際に介入できるようにしているシステムは少ない。

利用者にとって大きな負担となっているのは(d)である。認識結果は、基本的に1位候補を結果として用いることが多いが、認識スコアが1位候補と他の候補で大きく変わらない場合、文脈後処理により順位の入替えが行われ、2位以下の候補が認識結果として採用されることがある。また、スコアが悪く、結果が妥当でないと判断される場合、利用者に認識結果が正しいかどうか確認を求める処理を入れることで、認識結果訂正を行うことがある。

認識結果訂正については、認識結果を1つに特定しないことで訂正を不要にできると考える。これは、文字認識システムの認識率を考えるときに、1位認識率について着目すると必ずしも認識率が高くないが、例えば10位認識率(1~10位候補までに正解が含まれる確率)で考えれば、その認識率は、より100%に近づくからである[2]。つまり、1つの文字に複数の候補文字を付加しておくことで、その中のどれかに正解が含まれている状態にするのである。

この方法は、文字認識システムに対する負担を減らす代わりに、その後の利用で負担が増える危険がある。つまり、本来の正しい文字と関係のない情報が付加されていることで、検索の効率が落ちてしまう危険である。

しかし、実際に検索を行う場面では、複数文字で構成された単語を用いる場合がほとんどであり、1文字での検索を行うのはそれほど多くないと思われる。つまり、2文字以上の組み合わせで検索をかける場合、多少、間違った文字が含まれていても、その間違った文字が組み合わせられた形でその検索キーワードと一致する確率は低いであろうと考えられるのである。

3.3 本定義の目的

本定義では、XML形式を用いて次に示す情報(OCRによる文字認識結果)を記録し、これをもとに画像ファイルにどのような文字が書かれているか検索できるようにすることを目的とした。

- (a) 候補文字列情報
- (b) 文字認識スコア情報
- (c) 文字の画像上の座標情報
- (d) 画像ファイルの場所情報
- (e) 文字認識エンジン名称、バージョン情報

通常、OCRでは文章化したものをデータとして保存するが、そのためには認識率を100%に近づける必要があり、加えて、レイアウト認識や文字列切り出しを正確に行わせる必要がある。これには多大なコストや手間がかかるため、OCR結果は検索のみに使用し、元画像データは破棄しないで併用する設計とする。つまり、データの表示は元々の画像データを表示し、文字列検索をかける時はOCR結果を利用する。

このような検索方式は、既に一部で試行されている。例えば、目録カードの検索である[3]。目録カードは、一定のサイズ、一定の書式で書かれているため、OCRによる処理に向いている。しかし、万単位の数量ある目録カードを完全に文字コードに置き換えるのは、内容確認、修正の手間を考えると現実的ではない。かといって、単に画像情報としてデータを残すだけでは省スペース性や、同時に多人数が閲覧できるといったメリットはあるものの、検索性向上に活かされない。

このため、画像情報はそのまま残しておき、文字列検索を行う部分にOCR結果を使用することで、画像を文字に置き換える手間をなくし、認識結果の間違いによる情報の損失を防ぐことが可能になると考えた。

4. XML 定義

4.1 定義項目

ここでは、OCR結果を文書検索に使用するためのXML定義を考える。まず、使用する項目を次の通りに定めた。

(a) DocOcr 定義

- Id: 識別子
- Date: データを出力した日時
- SystemName: OCRシステムの名称
- SystemVersion: OCRシステムのバージョン情報
- Image: 画像データへのリンク情報

(b) OcrResult 定義

- LocationX: 画像上の文字のx座標
- LocationY: 画像上の文字のy座標
- SizeX: 文字画像の横幅
- SizeY: 文字画像の縦幅
- CandidateNumber: 候補文字数
- CandidateString: 候補文字列
- CandidateScore: 候補文字スコア

DocOcr定義では、文字認識システムについての定義と、認識にかけた画像データの情報を格納する。ここで、OCRシステムの名称やバージョン情報を格納するのは、文字認識システムによって出力される結果の傾向が異なる可能性があるからである。例えば、候補文字のスコアなどは、文字認識システムによって出てくる値の範囲が異なったり、値が小さい方が正解に近いのか、遠いのか、などが異なるのである。

OcrResult定義では、1文字ごとの文字認識結果を格納する。ここでは、Locationで画像のどこにその文字が書かれているかの座標(左上を原点としたx, y座標系)を記録し、Sizeでその文

