

イメージPDFを検索対象とする全文検索システムの運用

佐藤 郁・坂本正秀・大西みどり
(大分県農業技術センター)Kaoru Sato, Masahide Sakamoto and Midori Onishi:
Employment of the full-text search system by Image PDF

1. はじめに

大分県農業技術センターでは、試験研究成果の情報を公表しており、その一部は、ホームページ上で公開している。

しかし、多くの情報は印刷物として存在するため、全文を対象とした内容の検索や原文の利用ができなかった。そこで、紙文書の文献についてもインターネット上で検索できる全文検索システムを構築したので報告する。

2. データファイルの作成

パソコンで作成された文書であれば Adobe 社の PDF (Portable Document Format) に変換することで、検索が可能な原文と同じイメージを再現できる。しかも、OS に関係なく画面表示や印刷ができるという特徴がある。

一方、紙文書の文献情報を利用するには、文献を画像情報に変換する方法と、文献を OCR ソフト等で文字情報に加工する方法がある。

そこで、データファイルの作成には、それぞれの方法の特徴を組み合わせ、文献の原文イメージを表側の情報として保持し、また、文字情報を裏側の情報として保持することが可能な、イメージPDF (透明テキスト) 形式を使用した。

第1図にデータファイルの作成方法を示す。イメージスキャナを使用し、解像度は400dpi、色階調は白黒2値の条件で原稿を読み取り、TIFF形式の画像ファイルとして保存する。次に日本語OCRソフトを使用し、この画像ファイルの文字認識を行いテキスト形式のデータを抽出する。このTIFF形式の画像ファイルとテキスト形式のデータを融合し、イメージPDF (透明テキスト) 形式として保存する。

3. 全文検索システムの改良

全文検索システムは、文書管理、文書フィルタ、インデкса、検索エンジン、検索クライアントの各機能を持つプログラムの集合体である。日本語全文検索システムとしては、フリーソフトである Namazu が多くの WEB サイトで使用されている。

大分県農業技術センターでは、ホームページ内の情報を検索するために Namazu 1.3 を使用してきたが、検索できるものは、HTML 等のテキスト形式のファイルだけであった。しかし、Namazu 2.0 以降は、文書フィルタの種類が増えており、PDF 形式のファイルを検索対象とする事ができるようになったので、Namazu 2.0.12 を採用した。ただし Namazu では、事前に検索用インデックスファイルを作成する必要があるため、PDF 形式のファイルから検索に必要なテキスト情報を抽出するため、フリーソフトである Xpdf に含まれる pdftotext、pdfinfo というプログラムを使用した。

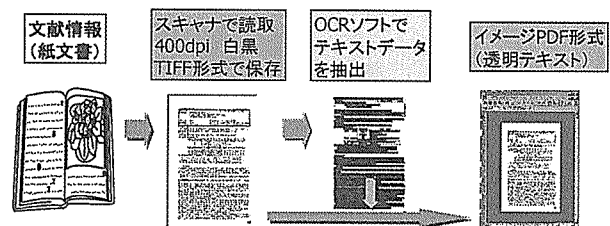
4. 全文検索システムの運用と課題

全文検索システムは、WEB サーバ上で運用しており、インターネット経由でブラウザを利用して検索することができる。ホームページ上の検索ボタンをクリックすると検索式の入力画面が表示される。検索式では複数の単語をスペースで区切って記述すると AND 検索ができる。第2図に検索結果の表示例を示す。

現在は、著作権や作業性の問題があり、この全文検索システムで利用できる文献は、大分県（農業技術センターおよび関係機関）が発行したものとなっている。また、利用者も、ホームページで公開したもの以外は関係機関職員に制限している。

今後は、データファイル作成手順を改善し、効率的に文献登録を行う必要がある。

また、データファイルの容量を小さくし、利用者が扱いやすくする必要がある。



第1図 データファイルの作成手順

Namazu による全文検索システム

現在、4,826 の文書がインデックス化され、177,140 個のキーワードが登録されています。
インデックスの最終更新日: 2003-07-28

検索式: [検索方法]

表示件数: 表示形式: ソート:

検索結果

参考ヒット数: [低コスト: 388] [技術: 2532] [水稲: 789]
検索式にマッチする 205 個の文書が見つかりました。

10. [B011-1998.pdf](#) (スコア: 206)
著者: 不明
日付: Mon, 30 Jun 2003 13:52:37
農業技術センター 第1号 水稲の低コスト栽培技術「低コスト栽培技術 II 灌水管理技術 平成11年3月 大分県農業技術センター発行」について 水稲の低コスト栽培技術の確立 稈白含量、収量向上のための措置
[/member/book/B01-1998.pdf](#) (6,063,037 bytes)

11. [seizer-h14-000.pdf](#) (スコア: 204)
著者: 不明
日付: Thu, 24 Jun 2003 13:52:10
平成14年度農研機構研究費助成事業「農産物の高品質・安全生産技術の開発 I 中山間地における転作大豆の高品質多収技術の開発(1)加工適性の高い高品質大豆栽培技術の確立」 稈白含量、収量向上のための措置
[/member/book/h14/seizer-h14-000.pdf](#) (818,773 bytes)

12. 大分県農業技術センター(久住試験地の施設内型) (スコア: 198)
著者: 不明
日付: Mon, 08 Jul 2003 15:25:01
大分県農業技術センター各種技術情報(各研究部のコーナー)久住試験地 秋作 久住試験地から九重連山を望む 大分県は山が多く、水稲の作付け地帯は標高0mの平地地帯から、800mの高冷地までと大きな幅広さがあります
[/data/tech/tech_center/030801](#) (5,816 bytes)

第2図 全文検索結果表示例