

3C-4

圧縮型手書き漢字データベース作成 のためのデータ圧縮の基礎実験

大倉 充[†] 今村 太^{††} 塩野 充[†]
 (†岡山理科大学工学部電子工学科 ††(株)三菱電機東部コンピュータシステム)

[1] まえがき

手書き漢字認識の研究では、開発した認識アルゴリズムを実際にプログラム化し、ソフトウェアシミュレーションによって認識性能を評価することが必要不可欠である。そして、その評価は認識対象となるデータの品質によって大きく左右されるため、各種のアルゴリズムを比較する場合には、共通のデータベースが必要となる。そのため、電総研によって、手書き漢字データベースETL-8⁽¹⁾やETL-9⁽²⁾が作成され、公開されている。しかし、これらのデータベースは大容量のため、磁気テープに格納されており、オープンリール磁気テープ装置の接続されたシステムでしか活用できないといった不都合な面を持つ。そこで本研究では、これらのデータベースのパソコンやワークステーション上での利用を目的とし、ETL8(2値化済みのB2タイプを用いた)の一部のデータを用いて、パソコンのOSとして一般的と思われるMS-DOS上でのファイル化及びデータ圧縮に関する基礎実験を行い^{(3),(4)}、更に、その実験結果の妥当性と最終的なデータベースの大きさの確認のために、ETL8の全データを用いて圧縮型データベースの作成を行った。

[2] 基本的なデータ圧縮法

一般の画像情報のデータ圧縮法は、可逆及び非可逆方式の2種類に大別される。本研究では、ETL8が、文字認識アルゴリズムの性能比較に用いられるデータベースであることより、前者を採用した。基本的には、ETL8が2値データであることより、通常、2値ファクシミリ信号の圧縮に用いられる手法を用いた(但し、MH符号化⁽⁵⁾のような統計量に基づく手法は用いていない)。

2-1 2ビット区切り方式⁽⁵⁾ ランレングス符号化法⁽⁵⁾の一つで、ランレングス(白または黒画素の継続する長さ)を2進数で表示した後に、下位から2ビットずつに区切り、各ブロックごとにその先頭に白領域ランレングスであれば0を、黒領域であれば1をつけて示す方式である。

2-2 複数ライン一括符号化法⁽⁵⁾ 代表的な方式として、2ラインのOR符号を用いる方式について説明を行う。まず図1に示すように、2ラインのOR符号で白区間と黒区間を定める。すなわち、第1と第2のラインがいずれも白(0,0)の場合は白区間、それ以外は黒区間とする。そして白区間では従来のランレングス符号化を、黒区間では原符号をそのまま用いる(但し、黒区間の終りは、00を挿入して区切りを明確にする)。本研究では、2ライン及び3ラインの2種類の一括符号化法を用いた。

2-3 理想圧縮率⁽⁶⁾ 一次元符号化においては、各ランレングスの生起確率より理想圧縮率を求めることができる。あるランレングスkの生起確率をPkとすると、1ラ

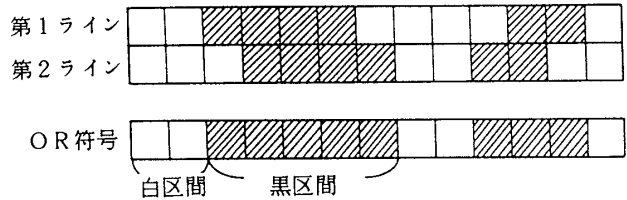


図1. 2ラインのOR符号

ンレングス当りのエントロピーHは、次式で与えられる⁽⁷⁾。

$$H = - \sum P_k \log_2 P_k \quad (1)$$

これは、一つのランレングスを2進表示するときに必要な平均ビット数の理論限界を表している。一方、平均ランレングスは $\sum k P_k$ であるから、理想圧縮率Cは、

$$C = (H / \sum k P_k) \times 100 (\%) \quad (2)$$

となる。

[3] データ圧縮実験

3-1 MS-DOSファイル化 まずETL8の図形情報に直接関係しない、ダミーレコードと識別情報部の削除を行う。次に15カテゴリずつ1枚のフロッピーディスク(1.25MB)に収納する。その結果、見本文字を含めて、ETL8の全サンプルが65枚のフロッピーディスクに収納される。使用装置は、本学情報処理センターの大型計算機FACOM/M380で、使用言語はFORTRAN77である。

3-2 実験データ 実験に使用したデータは、教育漢字の先頭「愛」から「角」までの100カテゴリで、1600サンプル/カテゴリゆえ、計16000サンプルとなる。また画面次数は、63×64である。圧縮実験に用いた装置はPC-9801E(パソコン)、使用言語はCである。

3-3 変換図形の作成 ランレングス符号化法には、ランレングスの値が大きい程、圧縮効率が良いという特徴があり、複数ライン一括符号化法では黒区間を符号化していないため、白区間ランレングスの値を大きくすることによって両手法での圧縮効率的向上が期待できる。そこで本研究では、原図形に対して各行あるいは各列ごとに排他的論理和をとった変換図形(圧縮時における走査方向の違いから、前者を横走査の変換図形、後者を縦走査の変換図形と呼ぶ)を作成し、その図形に対しても圧縮実験を行った。

3-4 圧縮率の定義 圧縮後のデータは、8ビット(=1バイト)の整数倍となっていない場合が多い。そのため、そのようなデータに対しては、データの終わりに1〜7個までの0を付け加えて、8ビットの整数倍となるようにしている。また、各サンプルによって圧縮後のデータの長さが異なるため、データの長さを2バイトを使って表し、データをファイルに格納する際、その情報をデータの先頭に付加している。圧縮率の計算には、この2バイトも含めており、圧縮率CRは、次式で定義した。

A Fundamental Experiment on Data Compression for Making Compressed Handprinted KANJI Character Data Base.

Mitsuru OHKURA[†], Futoshi IMAMURA^{††}, Mitsuru SHONO[†]

[†]Okayama University of Science.

^{††}Mitsubishi Electric Computer Systems(Tokyo) Co..

表1 平均圧縮率① (%)

圧縮方法	原図形	変換図形
2ビット区切り方式	46.3	37.1
2ライン一括符号化法	42.0	35.1
3ライン一括符号化法	41.2	37.2

表2 一次元符号化の理想圧縮率 (%)

圧縮方法	原図形	変換図形
横走査	33.8	27.1
縦走査	32.1	27.0

表3 平均圧縮率② (%)

横走査	縦走査	走査選択
32.6	32.5	31.7

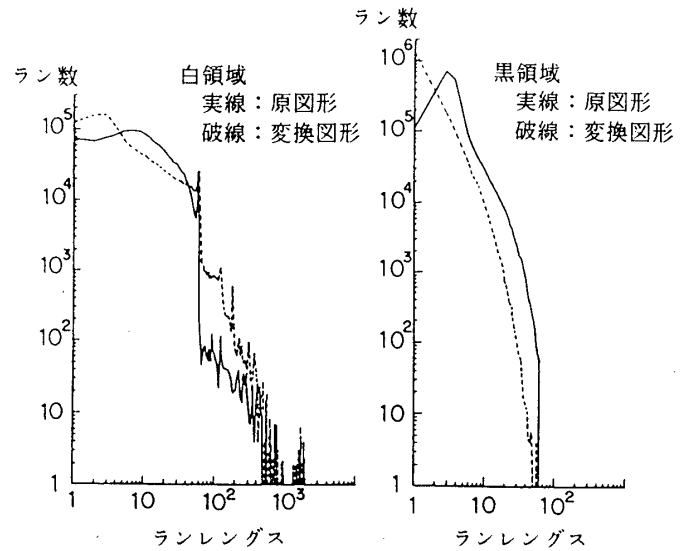


図2. ランレングス分布 (横走査)

表4 全サンプルの平均圧縮率 (%)

ひらかな 75カテゴリ	教育漢字 881カテゴリ	見本文字 956サンプル	全サンプル 153916サンプル
23.4	32.0	27.3	31.3

(圧縮後、実際にファイルに格納されたときのバイト数)

$$CR = \frac{504}{504} \times 100(\%) \quad (3)$$

ここで分母の504は原図形のバイト数を表している。

3-5 圧縮実験結果 表1に、得られた漢字100カテゴリ全体の平均圧縮率を示す。これらは、図形の左上から右下へ水平方向に走査(横走査)した結果である。図2に実験に用いた100カテゴリについての、ランレングスの分布を示す。横軸はランレングスを、縦軸はそのランレングスの出現した回数(ラン数)を表しており、紙面の都合上、横走査の結果のみを示す(縦走査の場合もほぼ同様の結果となる)。ここで、黒領域の分布で変換図形において分布のピークが、ランレングス1のところへ移動し、ランレングスの増加と共にラン数が減少していることが分かる。これは、黒領域においてHuffman符号化⁽⁵⁾がなされていることに相当する。また、この結果より求めた一次元符号化の理想圧縮率を表2に示す。これより変換図形においては、横走査及び縦走査で圧縮率に差がないことがわかるが、文字形状によっては、両走査の結果に大きな差異の生じることが推測される。以上述べたことより、2ビット区切り方式に対して、変換図形の黒領域には原符号を割り当てるという方法で圧縮実験を行い、得られた結果を表3に示す。表中、走査選択というのは、サンプルごとに両走査を行い、圧縮率の良い方を選択するという方法である。表1、表3より本実験においては、変換図形に対しての、2ビット区切り方式(黒領域原符号)走査選択で最も良好な圧縮率が得られた。

3-6 圧縮型データベースの作成 上述の実験結果の妥当性及び最終的なデータベースの大きさの確認のために、最も良好な圧縮率の得られた方法で、ETL8の全サン

ルのデータ圧縮を行い、圧縮型データベースを作成した。表4に得られた全サンプルの平均圧縮率を示す。この結果より、圧縮実験に用いた漢字100カテゴリに対して得られた結果と大差ないことが分かる。また、この圧縮によって、ETL8の全サンプルは、復元用プログラムを含めて、20枚のフロッピーディスクに収納された。なお圧縮されたデータの復元時間は、PC9801Eを用いて、1カテゴリ(160サンプル)当たり約120秒である。

〔4〕 おわりに

手書き漢字データベースの、パソコンやワークステーション上での利用を目的とし、基礎的なデータ圧縮実験を行い、最も良好な圧縮率の得られた手法を用いて、ETL8を基にした圧縮型データベースの作成を行った。しかし実用レベルで考えたとき、そのデータベースはまだ大きすぎるので(特に、ETL9を考えた場合)、更に圧縮率を向上させるデータ圧縮法についての検討が必要と考えている。

参考文献

- (1) 斎藤, 山田, 山本: "手書文字データベースの解析(V)", 電総研彙報, 45, 1, 2, 1981.
- (2) 斎藤, 山田, 山本: "手書文字データベースの解析(VIII)", 電総研彙報, 49, 7, 1985.
- (3) 大倉, 今村, 塩野: "手書き漢字パターンデータの圧縮に関する一考察", 昭62中国連大, 092109.
- (4) 大倉, 今村, 塩野: "手書き漢字データベースのデータ圧縮実験", 昭62関西連大, G8-54.
- (5) 安居院, 中嶋: "画像工学の基礎", 昭晃堂(1986).
- (6) 吹抜敬彦: "画像のデジタル信号処理", 日刊工業新聞社(1981).
- (7) 瀧 保夫: "通信方式", コロナ社(1985).