

多量 2 値文書画像の高解像度化

Enhancing Abundant Binary Document Images

佐伯 夏樹 荒井 利充 青木 恭太
SAEKI Natsuki ARAI Toshimitsu AOKI Kyota

1. はじめに

これまでの研究では、多数の低解像度画像または単一の低解像度画像から高解像度の画像を得る研究が行われてきた。[1]

大規模な公文書の電子化文書では数 10 万ページの電子文書が統一されたフォントで作成されており、多量の同一単語が存在することが期待できる。本研究の手法を用いると、このような多量に存在する低解像度文書画像の一部、または全体を高解像度化可能である。

本研究で扱う低解像度文書画像とは、各種の高解像度化手法が適用困難である 2 値画像で、文書画像中の各文字の分離が困難ではあるが、単語間では分離が可能な程度のもを想定している。

2 値画像では、既存の複雑な高解像度化アルゴリズムを適用することが難しく、また適用したとしてもその結果があまり期待できない。そこで、本研究では実験画像に出てくる多数の文字列塊を、知識として用いることにより文書画像を高解像度化する。

2. 提案方式の概要

本研究では、一度に処理する低解像度文書画像の枚数を 100 枚程度に想定しているが、実験画像より単語塊を作成する際、文字がつながってしまっているため文字単位の抽出は難しく、文字単位で切り出した場合に比べ位置決めが簡単であるため、英文書の特徴である単語間の空白を利用して単語単位での抽出を行う。

まず、低解像度文書画像の単語領域に対しラベル付けをおこなう。そしてそのラベル付けされた単語領域に対し単語領域の外接長方形領域を求める。この外接長方形座標により単語塊の抽出を行う。これらの画像片を単語塊と呼ぶ。次に、分解された単語塊を類似するもの同士でグループ分けを行い、高解像度化手法を用いて複数の低解像度単語塊から 1 つの高解像度化単語塊を得る。

高解像度化された単語塊は、元文書画像の対応する単語塊の適切な位置に貼り付けられる。

結果として、高解像度化された文書画像を得る。

図 1 に高解像度手順を示す。

3. 単語塊

最終的に高解像度化された単語塊画像の位置決めをする際、文字単位では位置決め困難であるが、単語単位ではサブピクセルの位置決めが可能であることに着目し、画像を文字単位ではなく出来るだけ単語として切り出す。

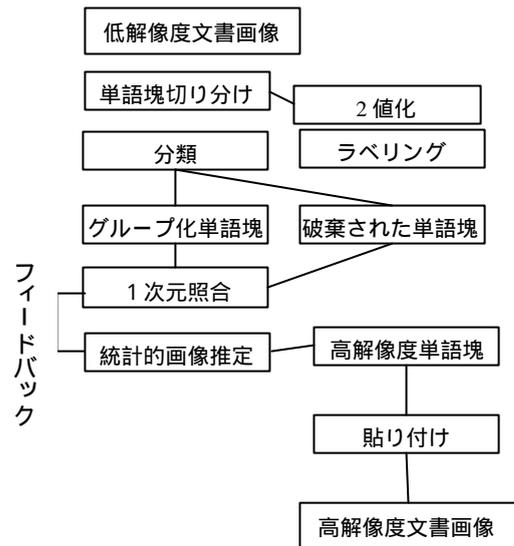


図 1. 高解像度化手順概略図

我々が対象とする低解像度画像は 2 値であり、中には文字がつながって文字ごとの分離が困難なものが存在する程度に劣悪なものであるため、文字単位の認識が出来ない場合がある。先に述べたように文字単位ではサブピクセルの位置決めが困難であるため本研究では、英文書を対象に絞り、単語間の空白を利用して実験画像を出来るだけ単語単位の塊に分け、この塊で照合を行うことを考えた。

単語塊は、同一の類似画像として複数枚収集できることを期待しており、単語塊を抽出するときに英文字単位で切り出されたか、英単語単位で切り出されたか、あるいはそのどちらでもないかどうかは問題にしない。

本研究では、100 枚程度の低解像度文書画像から類似した単語塊が複数抽出できることを期待しており、文書画像に含まれる文字列が、単語単位で十分に分離可能な程度に劣悪ならば、同一の単語塊を複数発見可能であることが十分に期待できる。分解された単語塊のうち一致しているものを探すために、これらの低解像度文書画像から得られた多量の単語塊を、類似したもの同士で分類する。分類された単語塊に対し、高解像度化手法を適用することにより高解像度化された単語塊を得るが、分類には 1 次元照合を主に用いている。

後に述べるように、照合精度に限界があるためフィードバックを組み込むことによりこれを解決する。

4. 単語塊の分類

切り出された単語塊は、その幅と高さによりグループ分けされる。最終的にグループは類似単語塊の集合となるが、位置合わせの問題や後の照合、フィードバックにかかる処理時間を短縮するために、幅と高さで単語塊をグループ分けする。次にグループ内の単語塊同時で、1次元プロフィール(図2)を用いた1次元照合で同一判定を行い、グループ内に同一の類似単語塊が集まるように候補を絞り込む。しかし、1次元照合はもともと画像同士の照合には向いていないため、その結果が良好でない場合がある。

例えば、(a o 0) や (l i 1) などの形状が非常に似ている文字を含む単語塊の比較は1次元照合では上手く処理できない可能性があり、これらの単語塊が同じグループに所属することになってしまう。

また逆に、照合の際に生じる単語塊画像同士の位置のずれによって、本来同一のグループに含まれるべき単語塊が別のグループに属してしまい、高解像度化が上手く行われないうえに冗長になってしまう。

これらを解決する方法として、1次元照合の他に2次元照合など、他の照合を用いてこれらを区別する方法も考えられるが、多量に存在する単語塊に対し、時間のかかるアルゴリズムを適用すると、その単語塊の量により処理時間が膨大に増加してしまう。そこで、本研究ではできるだけ処理時間が短い1次元照合アルゴリズムを主に用い、加えて高解像度化手法から導かれる誤差値をフィードバックとして画像同士の違いを検出することにより、照合結果を更に厳密にしようと試みるものである。

本研究では、高解像度化単語塊を同じグループに含まれる同一単語塊候補数枚、または数十枚から1つの高解像度化された単語塊を生成する。この時グループ内に、先に述べたような形状的に似ている文字を含む単語塊が存在すると、誤差の差が著しく大きい場合がある。

この値をフィードバックとして用い、グループ内の全ての単語塊の誤差がある閾値以下になるまで単語塊の分類を繰り返すことで、グループ内の単語塊の類似性を高める。この結果、グループに含まれる単語塊からの高解像度化精度を上げることが出来る。

5. 照合のフィードバック

照合はグループに対してそのグループ内の単語塊同士の誤差がある閾値以下に収束するまで行われる。

これにより、グループ内の単語塊は同一のものだけに分類されるようになり、結果として良好な高解像度化単語塊を得ることができる。さらに、グループ内の候補同士の誤差

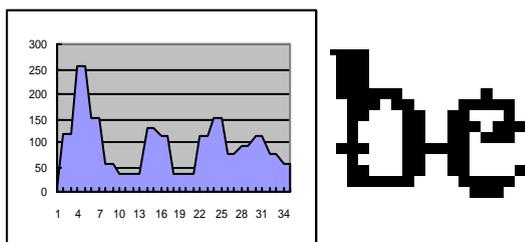


図2. 1次元プロフィール (be)

情報をフィードバックに用いることで1次元照合では検出が難しい単語塊の違いを検出することができるようになり、グループ内の単語塊同士の類似性が増大する。

この手法によって1つのグループから高解像度化された単語塊が得られると、これらの単語塊は人間が読める程度に高解像度化されており、OCR等を用いて完全に高解像度化された単語塊を得ることも可能である。

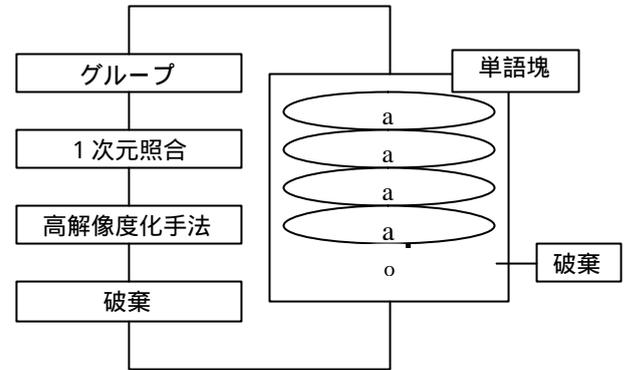


図3. フィードバック

6. むすび

単語単位で分類できるものはサブピクセルまでの位置合わせができる可能性があるが、文字単位で分類された画像に対してはサブピクセル単位での位置合わせは難しい。画像同士の位置合わせには画像の輝度重心を用いて位置合わせを行い、それを基準に照合を行うことを予定している。

低解像度文書画像より単語塊を抽出した結果、ほとんどの類似単語塊は1次元照合で分類できることが分かった。しかし、極めて形状の似ている単語や文字を含む単語塊では、それらを複数抽出できたとしても、画像内での対象文字の位置や大きさがまばらであると、1次元照合では十分ではない。本提案手法は、単語塊同士の照合の結果を改善するような、フィードバックをシステム内に含むことによってこの問題を解決する。

これによって、十分に高解像度化された単語塊はOCRを用いて認識することができ、文書の一部にとどまらず文書全体を高解像度化することができるようになる。

本研究では、フィードバックとして用いている誤差値の情報を1つのグループ内から単語塊を破棄するためにしか用いられておらず、本来同一のグループに所属すると思われる単語塊を判別し、同一グループに所属させるような機能はない。

今後の課題として、単語塊を高解像度化する際の冗長性を排除するために、グループから単語塊を破棄するだけでなく、複数のグループから同一の単語塊を収集することが望まれる。

参考文献

[1] 中野 康明: 文字認識・文書理解の最新動向 [] 電子情報通信学会誌 Vol83 No.6 pp.467-471 2000年 6月

[2] Sean Borman, Robert Stevenson : Spatial Resolution Enhancement of Low-Resolution Image Sequence.

A Comprehensive Review with Directions for Future Research.