

チラシ画像からの商品情報自動抽出—内容情報認識—

柴山美沙希, 高橋正信
芝浦工業大学

Automatic recognition of goods information in leaflets

-Content information recognition-

Misaki SHIBAYAMA, Masanobu TAKAHASHI

Shibaura Institute of Technology

要 旨 チラシ情報を記録し、参照、活用できるようにするため、チラシ画像から商品情報を自動認識しデータベース化する機能の実現を目指している。チラシ情報は商品の内容情報(会社名、商品名、内容量)と価格情報に分けられる。このうち、未だ実現されていない内容情報の認識機能の実現を目的とした。内容情報の認識には、複雑な背景における文字認識が必要である。そこで、Google Cloud Vision API の文字認識機能を利用したが、認識結果には誤字脱字や文字座標のずれなど多くの誤認識が含まれていた。これらの誤認識を自動修正して内容情報を認識するため、文字色と背景色の認識、それらを利用した座標の修正、自作の商品情報データベースを用いた誤字脱字の修正と会社名、商品名、内容量の分別などの処理を実現した。会社名、商品名、内容量が1つずつセットとなった154個の内容情報を用いた実験では、約半数の内容情報で文字認識結果に誤字や脱字が含まれていたが、92.9%の内容情報について会社名、商品名、内容量の全てを正しく認識できた。会社名と商品名だけであれば98.7%が正しく認識され、内容情報の認識手法として本手法は有望であると考えられる。

キーワード : チラシ, データベース, 文字認識, 商品情報, Google Cloud Vision

Abstract The purpose of this study is to automatically recognize goods information in leaflets images and to create a database in order to record and refer to leaflets information. Leaflet information is divided into the content information (company name, goods name and content) and the price information of the goods. We aimed to realize a function to recognize the content information, which has not been realized yet. In order to recognize the content information, it is necessary to recognize characters in a complex background. Therefore, characters were recognized using the OCR function of Google Cloud Vision API. In order to correct misrecognitions automatically and to recognize the content information, we realized the recognition of character color and background color; the correction of coordinates using these colors; the correction of misspellings using our own goods information database; and the separation of company name, goods name, and content amount. In the experiment, we used 154 pieces of content information, which consisted of a company name, a goods name and a content amount. Although about half of the content information contained misrecognition, 92.9% of the content information was recognized correctly. This method was shown to be effective as a recognition method of content information.

Keywords: Leaflets, Database, Character recognition, Goods information, Google Cloud Vision

1. はじめに

日本ではスーパーをはじめ多様な店舗のチラシをネットで閲覧できる。こうしたチラシ情報を記録し比較できれば最安値や販売傾向が分かり消費者にとって便利であるが、個人が全て行うのは困難である。その理由として、公開されているチラシ情報はテキストデータになっていないこと、また背景が複雑で特殊なフォ

ントが一部使用されているため OCR ソフトでも認識しづらいといったことが挙げられる。同種のサービスとして全国のチラシ情報を提供するサービス[1]もあるが、チラシ情報の認識とデータ化は人手で行っており、企業向けで高額な費用がかかる。

そこで、チラシ画像から商品情報を自動認識して日付や店舗名とともにデータベース化する機能の実現を目指している。商品情報の自動認識機能を実現でき

ば、例えば自動的にダウンロードした複数店舗のチラシ画像の商品情報を認識し、その日の最安値の店舗や特売情報など、利用者が知りたい情報を自動的に提供する機能を実現することができる。

商品情報は、商品名などの内容情報と価格で構成されている。このうち価格については埼玉県に多く店舗のあるヤオコー[2]を対象とした自動認識機能(認識成功率 99.35%)を実現し、認識された価格が税込か否かを識別する機能(識別成功率 100%)も実現した[3]。内容情報は会社名、商品名、内容量からなるが、その認識には多種多様な文字を認識したうえで、会社名、商品名、内容量の何れであるのかの認識が必要で、実現が難しかった。

本研究の目的は、残る課題である内容情報認識機能を実現することである。今回は価格認識[3]と同様にヤオコーのチラシを使用し、会社名、商品名、内容量が1つつセットとなった内容情報を対象にその認識機能の実現を目的とした。なお、ヤオコーのチラシのレイアウトなどは一般的なものであり、パラメータの調整などで多くの他店のチラシにも適用できると考える。

2. 関連研究

関連する従来技術としては、AI を利用したチラシ校正支援システム[4]がある。これは、チラシの制作を支援するもので、紙面データを個々の商品ごとにエリア分けし、入稿されたチラシ情報と比較することで校正ミスを減らすことを目的としたものである。しかし、対象とする紙面に含まれる商品情報が予め分かっていることが前提であり、本研究の状況とは異なる。また、消費者個人が利用するものでもなく、手法の詳細や精度なども公表されていない。

チラシの価格比較を目的として、画像の2値化と tesseract という光学文字認識エンジンを用いた研究[5]もあるが、チラシ特有の特殊なフォントや画数の多い漢字などの誤認識が多く精度も評価されていない。また、価格や商品名といった文字種の判定もできていないため、文字列を商品情報として認識するに至っていない。

チラシ画像の認識に関連する他の研究としては、チラシ画像中の食材名を認識してレシピを提案する手法が提案されている[6]。食材名は Google Cloud Vision API [7] の OCR 機能を用いて認識し、自作の食材データベースを利用して食材のテキストのみを抽出している。しかし、抽出される情報は食材だけであり、会社名や

内容量、そして価格は抽出されない。また、食材の抽出精度も 51.9%と、半分程度が認識できていない。

我々が調べた限り、チラシ画像中の内容情報を自動認識するという、本研究が目的とする機能を実現した報告はなかった。

3. 手法

内容情報とは、会社名、商品名、内容量のことである。その認識手順の流れを図1に示す。

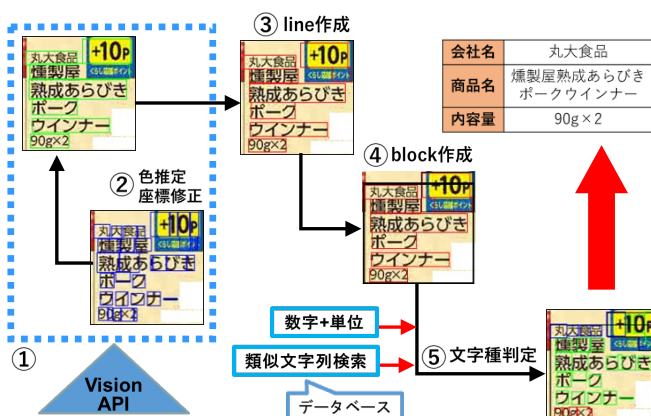


図1. 内容情報認識の手順

① Vision API による文字の認識

チラシ画像中の文字の認識には、複雑な背景での文字認識が可能な Google Cloud Vision API(Vision API) [7] の OCR 機能を利用する。Vision API により個々の文字の認識結果とその位置だけでなく、横方向に連続した文字列(text)とその位置も得られる。なお、認識結果には誤字や脱字、位置の不正確なものが含まれるため、それらを修正しながら内容情報を認識する。

② 文字の色推定と座標修正

認識された個々の文字の文字色と背景色を推定するとともに、Vision API で認識された文字の位置(文字の左右両端の座標)の間違いを修正する。

③ line の作成

line とは、1つの内容情報に含まれる1行分の文字列である。Vision API で認識された text には横方向に近接する別の内容情報や不要な文字が誤結合されている場合がある。そのような場合は text を分割して line を作成する。

④ block の作成

block とは1つの内容情報を含む文字列(line)の集合である。基本的に上下方向に近接する line をまとめて block を作成する。

⑤ 文字種の判定

1 つの block に含まれる文字列から 3 種類の文字種(会社名, 商品名, 内容量)を認識して 1 つの内容情報とする。

3.1 データベースの作成

多くの商品には JAN コードすなわちバーコードが登録されており, 会社名や商品名などの情報が記録されている。この JAN コードのようなデータベースを作成し, Vision API の認識結果に含まれている誤字脱字や余剰文字の修正に利用する。本研究では MySQL を用いて, 会社名, 商品名, 内容量, 価格からなる商品情報データベースを作成した。

データベースに格納する商品情報は, 対象としたヤオコーの通販ページ[8]から商品一覧(図 2)の HTML を保存して約 7000 件を抽出した。具体的には, 会社名, 商品名, 内容量, 価格に対応するクラス名が決まっているため, クラス名をキーワードとしてそれぞれの情報を自動抽出した。抽出した商品情報を csv データとしてリスト化し, MySQL のデータベースにインポートすることで商品情報データベースを作成した。

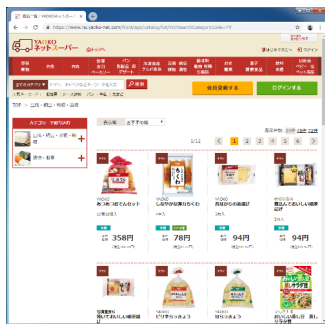


図 2. ヤオコーの通販ページ画面例[8]

3.2 Vision API による文字認識

内容情報の認識には, 複雑な背景での様々なフォントを対象とした文字認識が必要である。本研究では個人で容易に利用可能な機能の実現を目的としているため, 基本的に無料で利用可能な Vision API の OCR 機能を利用した。Vision API の推奨サイズ(1024×768 画素)に合わせて画像を分割し Vision API に入力として与えると, 1 文字ごと及び文字列ごとの認識結果(テキストデータ)と位置情報(外接長方形の 4 辺の座標)が得られる。表 1, 表 2 は図 1 内の画像に対する認識結果例である。座標は入力画像の左上を始点とした 2 次元座標である。文字列とは 1 行分の文字列であり, その認識結果を text とする。

Vision API の認識結果には誤字や脱字, 近接する別の文字を含んで 1 つの文字列として認識された text(図 3)や, 外接長方形の位置ずれ(図 4)といった誤認識が存在し, 実験で使用した内容情報の半数以上が該当する。そのため, そうした誤認識への対策が必要となる。

表 1. 1 文字ごとの認識結果例

文字	上	下	左	右
丸	276	293	1117	1137
大	278	293	1139	1151
⋮	⋮	⋮	⋮	⋮

表 2. text の認識結果例

Text	上	下	左	右
丸大食品	276	293	1117	1180
燐製屋イント	291	310	1122	1245
⋮	⋮	⋮	⋮	⋮



図 3. 近接する別の文字を含む文字列(text)



図 4. 文字の外接長方形の位置のずれ

3.3 文字の色推定と座標修正

Vision API による誤認識の 1 つである文字の外接長方形の位置ずれ(図 4)を修正する。文字色と背景色の推定を行い, その色情報を利用して座標の修正を行う。

3.3.1 文字色と背景色の推定

1 つの文字領域の画素値に対して, k-means 法を用いて代表色 3 色を抽出し, 画素値を 3 色に置換する(図 5)。画素値を置換した画像において領域外縁の最頻値色を背景色とし, 残る 2 色のうち背景色と明度の差が大きい方の色を文字色とする(図 6)。

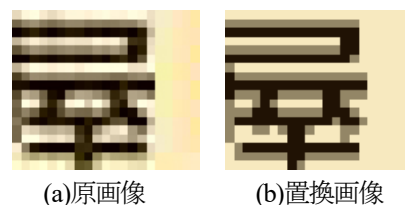


図 5. 色推定

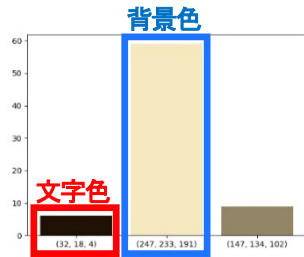


図 6. 領域外縁の色出現頻度

3.3.2 座標の修正

文字色と背景色を利用して、各文字の外接長方形の左辺と右辺の座標を修正する。文字領域を左右に R_E ずつ拡張し、拡張領域内で縦の列ごとに文字色が登場する回数を求める。文字色の出現回数 0 回の座標のうち、出現回数が 1 回以上の座標の隣接している座標は文字領域の端である可能性があるため、候補座標として抽出する。図 7 の例の場合、赤破線が文字「ん」に対して認識された外接長方形の右辺と左辺の元座標で、それを左右に拡張した領域内で抽出された候補座標が青破線で示されている。次に、抽出された複数の候補座標から「ん」の左辺と右辺の座標を 1 つずつ選択する。その選択には(1)式の値 D を用い、 D が最小となる右辺と左辺の組み合わせを選択する。この操作は同一 text の先頭から順に行い、拡張後の領域が左隣文字と重複した場合は、重複した領域を除いて、図 8 の緑破線に示すように左隣文字の領域外から候補座標の抽出を開始する。

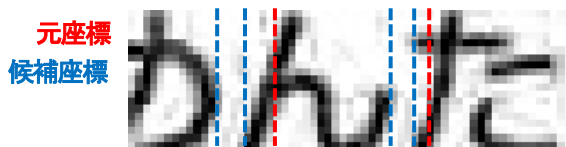


図 7. 拡張領域と候補座標

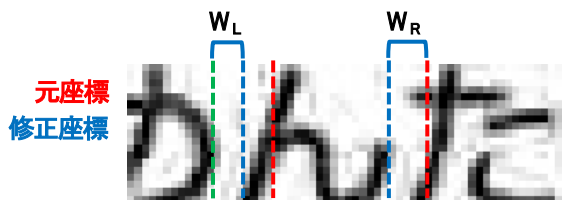


図 8. 元座標と修正座標

$$D = W_L + W_R + \alpha \times A \quad (1)$$

$$W_L = |\text{左辺候補座標} - \text{左隣文字の右辺座標}|$$

$$W_R = |\text{右辺候補座標} - \text{元の右辺座標}|$$

$$A = |\text{見本画像の縦横比} - \text{候補座標の縦横比}|$$

(1)式において、 W_L は左隣文字の右辺座標(図 8 緑破線)と左辺候補座標の差であり、左辺座標が正しければ小さな値となる。なお、text の先頭文字については左隣文字が存在しないため、元の左辺座標との差を利用する。 W_R は右辺候補座標と VisionAPI で抽出された元の右辺座標との差であり、右辺座標が元座標の近くで選択される働きを持つ。 A は見本画像と候補座標の縦横比の差である。係数 α は $\alpha \times A$ が W_L , W_R と同じオーダーになるように文字ごとに自動設定される。

見本画像とは、内容情報に使用されているフォントと似た「BIZ UD ゴシック」で作成した 1 文字ずつの画像のことで、英数字やひらがな、カタカナ、漢字まで含めた約 9000 字が保存されている。図 9 の例のように、文字の縦横比は大きく異なる場合があるため、縦横比の差 A を利用することで誤った選択を防ぐ。以上の処理により修正された結果例を図 10 に示す。

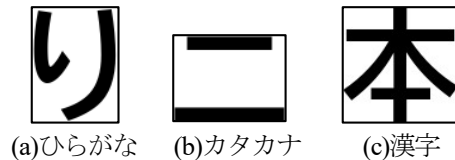


図 9. 見本画像例



(a)元座標



(b)修正後座標

図 10. 座標の修正結果例

3.3.3 背景色の再推定

1 つの内容情報内の背景色は同じであり、それが後述する line の作成(3.4)でも利用されているが、図 11 の「ズ」のように文字の周囲に商品画像が入ると商品画像の色を背景色と誤認識してしまい、同一 text 内であっても隣接文字との背景色が異なってしまう。この問題を改善するため、背景色を再推定する。

背景に商品画像が入る場合、強調するために図 12(a)のように文字を縁取りしていることが多い。この縁取りは隣接文字の背景色と同じ色が使用されているため、座標修正後の文字領域において文字色の画素の 8 近傍画素(図 12(b))を抽出し、その平均色を背景色として再設定する。表 3 は図 12(a)の原画像に対する再推定結果であり、再推定によって正しい色に修正されている。



図 11. 背景色の異なる text



(a)原画像 (b)文字色の周囲画素

図 12. 背景色の異なる text

表 3. 背景色の再推定

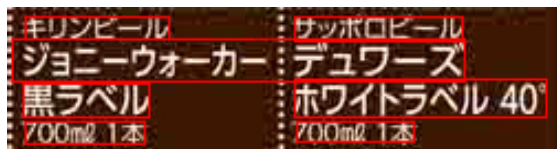
	背景色	
前	(211, 207, 182)	
後	(226, 230, 68)	

3.4 line の作成

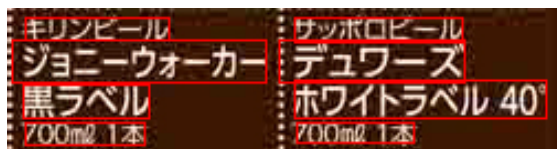
line とは、1 つの内容情報に含まれる 1 行分の文字列である。Vision API で認識された text には左右に近接する別の内容情報、見出しや背景の商品上のラベルなど不要な文字が誤結合されている場合がある。そのため text を分割し、line を作成する。

3.4.1 左端位置による分割

図 13(a)に示すように、内容情報の中間部分の text が、別々の内容情報に属する line が結合されたものになっている場合がある。



(a)分割前



(b)分割後

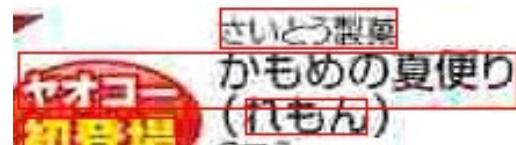
図 13. 左端位置による分割

同じ内容情報に属する line の左端位置はほぼ同じであるため、上下に近接する text (上 text の下端 - 下 text の上端 $< T_L$) どうしを比較し、上 text よりも下 text が左に存在するとき、上 text の先頭文字と最も近い下 text の文字色が同じであれば (RGB 全てにおいて文字色の差 $< D_c$)、該当位置で分割を実行する(図 13(b))。

3.4.2 文字による分割

図 14(a)のように見出しや商品画像中の文字と誤結合している場合がある。これらの多くは隣接する内容情報とは異なる文字色又は背景色であるため、隣接文字の色を比較して図 14(b)となるように分割を行う。

具体的には、同一 text 内で隣接文字の文字色の差 $> D_c$ 、又は背景色の差 $> D_B$ となる場合、該当文字間を分割候補地点として抽出する。その後、分割ありと分割なしの 2 通りに対しデータベースを用いて文字列検索を行う。このとき MySQL の N-gram 全文検索機能を用いて対象を絞り込む。全文検索は高速化の研究が多くなされており、N-gram についても効率的な手法[9]が実現されている。この文字列検索により得られた複数の内容情報に対して編集距離[10]を算出し、最小となる値を分割あり、なしそれぞれのコストとして記録する。編集距離とは、2 つの文字列がどの程度異なっているかを示す指標である。1 文字の挿入、削除、置換を 1 回として、一方の文字列をもう一方の文字列に変換するために必要な手順の最小回数として定義される。分割ありについては前後文字列のそれぞれで編集距離を算出し、加算したものをコストとする。分割ありと分割なしのコストを比較し、前者の値が小さい、又は表 4 のように編集距離が同値でも前後文字列のどちらかが完全一致であれば分割を実行する。また、隣接文字どうしの横の間隔 $> T_w$ の場合も分割する。



(a)分割前



(b)分割後

図 14. 文字の色による分割

表 4. 分割ありと分割なしの比較

	分割あり	分割なし
原文	ヤオコー／ かもめの夏便り	ヤオコーかもめ の夏便り
データベース の内容情報	“該当なし”／ かもめの夏便り	かもめの夏便り
編集距離	4 + 0	4

3.4.3 重複領域の削除

line の領域の高さが実際の文字の高さよりかなり大きく抽出され, 図 15(a)のように line の領域が重なっている場合がある。そこで, 領域が重なっている場合は, 図 15(b)に示すように高さが大きい方から重複領域を削除する。

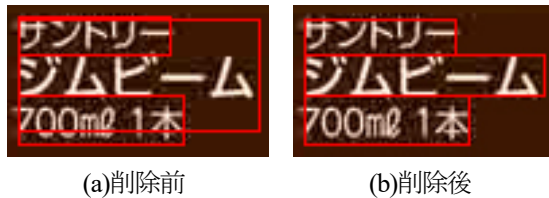


図 15. 重複領域の削除

3.4.4 内容情報以外の文字列の削除

漢字を除いて 1 文字のみで構成されている line や, チラシで頻繁に用いられる「精肉コーナー」, 「お 1 人様 2 点まで」といった定型文は内容情報ではないため削除する。

3.5 block の作成

上下に近接し, かつ左端の位置が近い line どうしを結合する。具体的には, 上 line の下端 - 下 line の上端 $< T_L$ かつ, 左端座標の差 $< T_l$ となる上下の line どうしを結合し, 1 つの内容情報を含む block を作成する。

3.6 文字種の判定

3.6.1 内容量の抽出

1 つの block 内の文字を左上から順番に接続して 1 つの文字列とする。この文字列は会社名, 商品名, 内容量の順番で構成されている。

3 つの文字種のうち, 初めに内容量の抽出を行う。内容量は以下の(a)~(c)のパターンで表記される(図 16)。そこで, block を末尾から検索して, この(a)~(c)のパターンの表記を抽出する。「単位」は「当り」や「当たり」, 全角と半角など表記ゆれのパターンを含んだ辞書を保持し, 使用する。

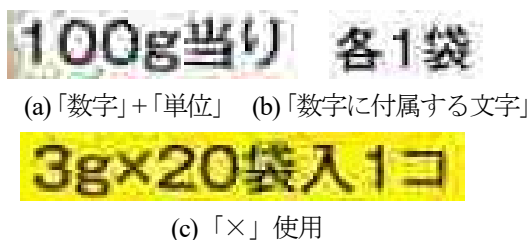


図 16. 内容量の表記パターン

(a) 「数字」+「単位」

(b) 「数字に付属する文字」(「各」, 「約」など) + (a)

(c) (a) or (b) + 「×」 + (a)

また, 誤字脱字により内容量を抽出できなかった場合は, 会社名, 商品名の認識後に末尾の余剰文字を内容量とする。

3.6.2 会社名と商品名の認識

内容量を除いた block 内の文字列は会社名と商品名からなる。そこで, 内容量を除いた文字列の適切な分割位置を求める。具体的には, 文字列の文頭から分割位置を 1 文字ずつずらして分割し, 文字列の前者を会社名, 後者を商品名の候補としてデータベースの検索を行う。

検索の際に問題となるのが半数近くの内容情報にある文字の誤字脱字である。誤り検出を行う場合, N-gram や形態素解析などを様々な識別方法が提案されている[11][12]。文であれば単語同士の接続規則や接続情報などが制約として考えられるが, 内容情報は固有名詞の並びであるため, 自作の商品情報データベース(3.1)を辞書として活用する。会社名, 商品名のそれぞれで編集距離を求め, 編集距離の合計値が最小となる内容情報をデータベースから引用する。これにより, 誤字脱字や余剰文字に頑強な検索と正しい名称への修正が可能となる。

表 5 に会社名と商品名の候補, およびその編集距離の例を示す。この例では誤字脱字が無い場合, 正しい分割位置で編集距離の合計値は 0 となる(図 17)。

表 5. 分割候補と編集距離

会社名	商品名	編集距離(合計)
久	原醤油あごだしつゆ	6
久原	醤油あごだしつゆ	4
久原醬	油あごだしつゆ	2
久原醤油	あごだしつゆ	0
⋮	⋮	⋮

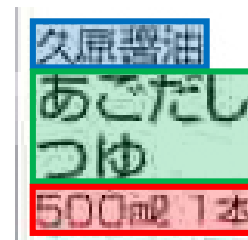


図 17. 文字種の判定

4. 実験

4.1 実験データ

ヤオコーのホームページ[4]からダウンロードした2018年6月から2020年11月までのチラシ画像12枚を実験に使用した。チラシ画像1枚は3700×2600画素以上あるため、Vision APIの推奨サイズ(1024×768画素)に合わせて画像を分割し、計53枚の分割画像を使用した。内容情報の認識実験は、チラシ画像中で会社名、商品名、内容量が1つずつセットとなった内容情報154個を対象として行った。表6にVision APIによる認識結果に誤字や脱字のあった内容情報の内訳を示す。実験対象の内容情報のうち81個(53%)に誤字か脱字が含まれていた。

- (a) 誤字も脱字もない
- (b) 脱字はないが、誤字を含む
- (c) 脱字がある(誤字はありとなしがある)

表6. 実験データの内訳

(a)	(b)	(c)	合計
73 個	61 個	20 個	154 個

4.2 実験結果

手法中の諸パラメータは結果を目視で確認して最適化し、 $R_E = 25$, $T_L = 13$, $D_C = 52$, $D_B = 50$, $T_W = 18$, $T_I = 30$ と設定した。結果の評価は以下の①～④の各段階で行った。

- ① line の作成
- ② block の作成
- ③ 会社名と商品名の認識(②で正解したもののみ)
- ③' 会社名と商品名の認識(全ての内容情報について)
- ④ 内容情報(会社名, 商品名, 内容量)の認識

表7は正しく処理された内容情報の個数であり、誤字や脱字があるかどうか(表6(a)～(c))で内容情報を分けて集計した結果である。③はblockを正しく作成できた内容情報のみについて会社名と商品名を正しく認識できた個数である。図18(a)のように文字の過不足なくblockを正しく作成できたもの(②, 全体の87.7%)については、会社名と商品名も全て正しく認識できた。③'はlineやblockを正しく作成できなかった場合も含めた全ての内容情報について、会社名と商品名を正しく認識できた個数である。図18(b)のような脱字などで文字の過不足があるものも含まれるが、商品情報データベースを用いた修正機能により98.7%で正しく認

識できた。

最終的に会社名, 商品名, 内容量の全てを正しく認識できた内容情報(④)は, 154 個中 143 個(92.9%)であった。誤認識された11個のうち6個はVision APIの認識結果に脱字がある場合で, 主に内容量が抽出できなかったことが原因である。その他はVision APIによる文字座標の誤差が, 座標修正ができないほど大きく, textを誤分割した場合であった。

表7. 実験結果

	①	②	③	③'	④	全体数
(a)	66	65	65	73	72	73 個
(b)	58	54	54	61	59	61 個
(c)	19	16	16	18	12	20 個
合 計	143 92.9%	135 87.7%	135 87.7%	152 98.7%	143 92.9%	154 個



(a)文字の過不足なし (b)脱字あり(「一」)

図18. 認識成功例

5. 考察

内容情報認識において92.9%の認識正解率を得たが, 誤認識の多くは脱字によるものであった。特に内容量は文字サイズが小さいため, 図19のようにVision APIでも全く認識されない場合がある。このように内容量が全て脱字となってしまったものについては, 商品情報データベースに記録されている内容量を参照したり, 価格と紐づけし最も近い価格の内容量を引用したりすることで解決できると考える。

また図20は「ポ」の初期座標が背景の商品画像に重なるように大きくずれているため, 表8に示すように背景色と文字色を誤って推定し, 背景色の再推定(3.3.3)でも補正されなかった場合である。その結果, 文字の左端座標が修正されず, lineの左端座標のずれが上下のlineを結合する閾値(T_I)より大きくなり, blockの作成に失敗している。これを改善するため, 座標修正のより良い手法の実現が必要となる。

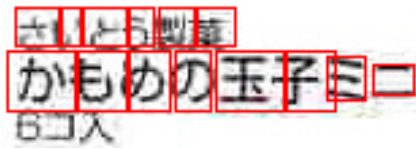
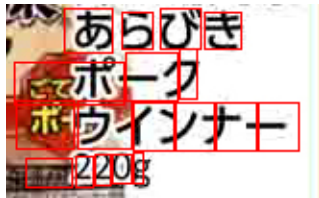
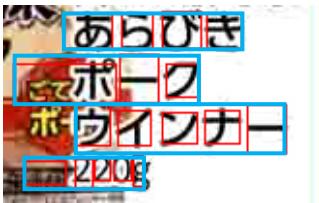


図 19. 内容量が全て脱字







(a)元座標



(b)座標修正失敗

図 20. 大きな座標ずれ

表 8. 「ポ」の色の誤推定

文字	背景色	文字色
ポ		
—		

実現した手法のヤオコー以外のチラシへの適用可能性を調べるため、内容情報が同様なレイアウトであるイトーヨーカドー[13]のチラシ1面分を用いて39個の内容情報について実験を行った。手法中のパラメータをヤオコーのチラシ画像に対する実験(4.2)と同じにした場合、会社名と商品名を正しく認識できた内容情報は39個中37個(94.9%)、会社名、商品名、内容量の全てを正しく認識できた内容情報は39個中32個(82.1%)であった。イトーヨーカドーのチラシに合わせてパラメータの一部を $D_c=50$, $T_l=23$ と変更し、単位の表記ゆれのパターンを追加したところ、会社名と商品名を正しく認識できた内容情報は変わらず39個中37個(94.9%)であったが、会社名、商品名、内容量の全てを正しく認識できた内容情報は39個中36個(92.3%)に改善された。図 21 に認識に成功した内容情報の例を示す。この結果から、実現した自動認識機能は、パラメータを調整することで同様なレイアウトを持つ他店のチラシに適用できる可能性が示された。

チラシ画像から情報を抽出する報告としては、食材

名を認識する手法の報告[6]はあるが、会社名、商品名、内容量という内容情報全ての自動認識機能を実現した例は我々の知る限り他にない。



図 21. 認識成功例 (イトーヨーカドー)

6. おわりに

内容情報認識において Vision API の文字認識結果に含まれている文字や文字位置の誤認識を改善するため、文字の色情報とデータベースの文字列検索を活用した。各文字の文字色と背景色を利用し、座標のずれの修正や余剰文字の分割を行う。また、自作の商品情報データベースから編集距離を用いて類似の内容情報を抽出することで、誤字脱字を修正しながら文字列を会社名、商品名、内容量に分別することが可能である。

会社名、商品名、内容量が1つずつセットとなった内容情報を対象として実験を行い、全ての内容情報(会社名、商品名、内容量)については92.9%、脱字の多い内容量を除いた会社名と商品名であれば98.7%の精度が得られた。

今後の課題としては、誤認識が多い内容量への対策による認識正解率の改善と、複数の商品名を含む箇条書きなどへの対応が挙げられる。また、他店のチラシへの適用についても進めていきたい。

参考文献

- [1] 株式会社ドゥ・ハウス, “全国チラシ情報サービスセンター”, <https://www.dohouse.co.jp/>, (参照 2020-11-01)
- [2] 株式会社ヤオコー, “ヤオコー MARKETPLACE”, <https://www.yaoko-net.com/>, (参照 2020-11-10)
- [3] 染谷謙太郎, 高橋正信: チラシ画像からの商品情報自動抽出—価格認識—, 電子情報通信学会総大会学生ポスターセッション, ISS-SP-250 (2007).
- [4] 方正株式会社, “AI による校正支援でチラシ制作コストを大幅削減”, <https://www.founder.co.jp/archives/1994>, (参照 2020-12-28)

- [5] 片桐圭祐, 田村仁: チラシの価格表示比較支援システム, 情報処理学会第 78 回全国大会, 5N-06 (2016).
- [6] 釜田祐哉, 伊東聖矢, 金子直史, 鷺見和彦: 食品チラシ画像を用いたレシピ推薦システム, 精密工学会誌, Vol.85, No.12, pp.1127-1135 (2019).
- [7] Google Cloud, “Cloud Vision API”, <https://cloud.google.com/vision?hl=ja>, (参照 2021-01-08)
- [8] 株式会社ヤオコー, “YAOKO ネットスーパー”, <https://www.ns.yaoko-net.com/front/app/common/index>, (参照 2020-11-10)
- [9] 小川泰嗣, 松田透: n-gram 索引を用いた効率的な文書検索法, 電子情報通信学会論文誌 D, Vol. J82-D1, No.1, pp.121-129 (1999).
- [10] Vladimir I. Levenshtein.: Binary codes capable of correcteing deletions, insertions, and rever-sals, Soviet physics doklady, vol.10, No.8, pp.707-710 (1966).
- [11] 伊東伸泰, 丸山宏: OCR 入力された日本語文の誤り検出と自動訂正, 情報処理学論文誌, Vol.33, No.5, pp.664-670 (1992).
- [12] 河田岳大ほか: 両方向 n-gram 確率を用いた誤り文字検出法, 電子情報通信学会論文誌 D, Vol.J88-D2, No.3, pp.629-635 (2005).
- [13] 株式会社イトーヨーカ堂, “イトーヨーカドー”, <https://www.itoyokado.co.jp/>, (参照 2021-02-24)

著者紹介

柴山美沙希

芝浦工業大学大学院・理工学研究科・システム理工学専攻に所属。画像処理を用いたチラシ画像認識に関する研究に従事。

高橋正信

1986 年大阪大学大学院・工学研究科修士課程修了。同年, 三菱電機株式会社入社。2001 年芝浦工業大学システム工学部電子情報システム学科助教授。現在, 同大学システム理工学部電子情報システム学科教授。画像処理, 画像認識および応用システムの研究に従事。博士(工学)。正会員。