

複数特徴量の重み付け統合による一般物体認識

須賀 晃[†] 滝口 哲也^{††} 有木 康雄^{††}

[†] 神戸大学大学院工学研究科 〒 657-8501 兵庫県神戸市灘区六甲台町 1-1

^{††} 神戸大学自然科学系先端融合研究環 〒 657-8501 兵庫県神戸市灘区六甲台町 1-1

E-mail: [†]akira1234@me.cs.scitec.kobe-u.ac.jp, ^{††}{takigu,ariki}@kobe-u.ac.jp

あらまし 本稿では、複数の特徴量の重み付け統合による一般物体認識手法を提案する。近年、特定物体認識で良いパフォーマンスを得ていた SIFT 特徴を量子化し、visual word のヒストグラムとして画像を表現することで、一般物体認識へ拡張した bag-of-features を用いた手法が注目を浴びている。しかし、物体内部の輝度変化が少ないものや、他カテゴリでも共通する特徴を多く含むものについては、SIFT 特徴だけでは認識率が低いという問題も見られた。本手法では、SIFT 以外にも SIFT では用いていなかった色情報や、人検出でよく用いられている大域的領域情報の HOG 特徴を用いる。加えて、各特徴のヒストグラムインターセクションを特徴量の頻出度とし、言語処理の単語重要度抽出理論を画像分野に適用することで、識別的な特徴量を自動的に学習する手法を提案する。これにより、特徴の情報量や特徴とカテゴリの共起度・関連性が考慮され分類精度が向上した。実験の結果、従来の bag-of-features だけでは認識が難しかったカテゴリに対応し、特徴量の重み付けにより識別に有効な特徴を学習することでより精度を向上させることができた。

キーワード 一般物体認識, 画像分類, SIFT, Bag-of-features, Histograms of Oriented Gradients

GENERIC OBJECT RECOGNITION BASED ON WEIGHTED INTEGRATION OF MULTIPLE FEATURE

Akira SUGA[†], Tetsuya TAKIGUCHI^{††}, and Yasuo ARIKI^{††}

[†] Graduate School of Engineering, Kobe University 1-1, Rokkodai, Nada, Kobe, Japan

^{††} Organization of Advanced Science and Technology, Kobe University 1-1, Rokkodai, Nada, Kobe, Japan

E-mail: [†]akira1234@me.cs.scitec.kobe-u.ac.jp, ^{††}{takigu,ariki}@kobe-u.ac.jp

Abstract This paper describes a method for generic object recognition based on weighted integration of multiple features. In recent years, bag-of-features approach has become popular for generic object recognition. This approach represents the image as a histogram of visual words by extending the theory of SIFT which offers good performance on specific object recognition. However, only by SIFT, it was hard to recognize objects with a few texture or recognize objects whose appearance of the visual word changes variously. Therefore, in addition to SIFT feature, we propose to employ color feature which is not being used in SIFT, and HOG feature which is global information whereas SIFT is local information. Moreover, we propose an automatic learning of the distinctive features for each category by adopting feature importance extraction theory which has been mainly used in language processing for the improvement of the classification performance. Experimental result shows that the proposed method can deal with the above problem and obtain good classification performance by weighting distinctive features for each category.

Key words Generic object recognition, Image classification, SIFT, Bag-of-features, Histograms of Oriented Gradients

1. はじめに

近年、画像データの大容量化に伴い人手による画像データの分類や検索が困難となってきており、計算機による一般物体認識の必要性が高まっている。一般物体認

識では、同一物体カテゴリ内のバリエーションが多様であり、物体の変形などによってアピランスも変化するため困難な問題の 1 つである。一般物体認識のアプローチの 1 つとして、bag-of-features [1] が挙げられる。bag-of-features は統計言語処理における bag-of-words [2]

のアナロジーで、まず、SIFT(Scale-Invariant Feature Transform) [3] 特徴ベクトルを k-means クラスタリングによりベクトル量子化し、keypoint を visual word として扱い visual vocabulary を作成する。そして画像を visual word のヒストグラムとして表現することで、一般物体認識を可能としている。この手法は良いパフォーマンスを得ているため注目されている。しかしながら、物体内部の輝度変化が少ないものは SIFT 特徴があまり得られず、また、どのカテゴリにも同じような visual word が頻出するものについても識別率が低く、世の中のあらゆる対象物体に対して、局所的な輝度勾配特徴である SIFT 特徴だけでは対応出来ない。そこで、SIFT ではカバーできない特徴を補うため、SIFT に加え複数の特徴量を用いて一般物体認識を行う。1 つは SIFT では用いられていなかった色特徴量、もう 1 つは SIFT では除外してしまっているエッジ成分を含む大域的形状情報である HOG(Histograms of Oriented Gradients) 特徴を用いる。HOG 特徴は、局所的な幾何学的変化と明度変化に不変で、主に人検出に用いられており、高い精度が得られていることが報告されている [4]。本研究では、これら 3 つの特徴量を用いることにより、SIFT の bag-of-features だけでは認識が難しかったカテゴリに対して認識率の向上を図る。更に、カテゴリ毎に特徴的な特徴量が違うことにも着目する。例えば、"motorbike" は種類によって色は異なるため色特徴はあまり重要ではなく、むしろ大まかな形状情報やタイヤなどの局所特徴の方が重要と思われる。一方で、"tomato" の場合、局所的な輝度変化はあまりなく、形状や色情報特徴の方が重要と考えられる。そこで本手法では、言語分野の研究で用いられている tf-idf などの重要語検出法を、ヒストグラムインターセクションを頻出度として考えることで画像の重要特徴量抽出に応用し、各カテゴリ毎に識別に重要な特徴量を学習し、それを各特徴量の重みとして与えて認識を行う。

2. 複数特徴量の重み付け統合

本手法では複数特徴量を組み合わせ、特徴重要度抽出法によりカテゴリ毎に識別的な特徴量を自動的に学習し、特徴量に重みを付けて物体認識を行う。本章では用いる特徴量と、その重要度の学習方法について述べる。

2.1 特徴量

2.1.1 局所特徴量 (Bag-of-Features)

bag-of-features では、図 1 に示すように、まず各学習画像から SIFT 特徴を抽出する。SIFT は特徴点の検出と記述を同時に行うアルゴリズムで、スケールの異なるガウシアンフィルタにより得られた平滑化画像の差分画像から極値を検出し、その周辺領域を 4×4 ブロックに分割し、各ブロックごとに 8 方向の勾配方向ヒストグラムを作成する。これにより回転・スケール変化にロバストな 128 次元の特徴量が抽出できる。次に、得られた全

SIFT 特徴に対して SIFT 特徴空間上で k-means クラスタリングを行う。このクラスタリングによって得られた各クラスタを visual word とみなし、visual vocabulary を構築する。1 枚の画像に対し、visual word のヒストグラムを 1000 次元の特徴ベクトルで表現する。

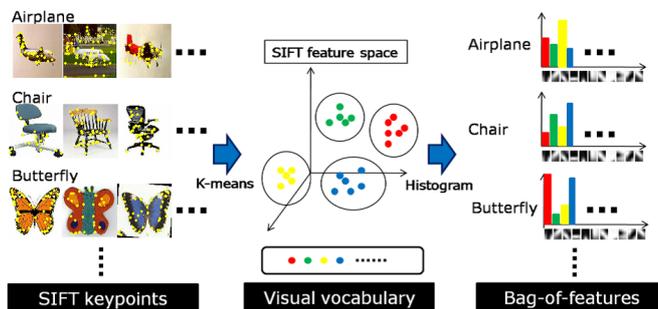


図 1 Bag-of-features

2.1.2 色特徴量

SIFT 特徴では、グレースケール画像からの輝度変化を用いているが、色情報は用いていない。しかしカテゴリによっては色情報が重要な場合が考えられる。そこで、画像から得られる RGB のヒストグラムを $256 \times 3 = 768$ 次元の色特徴量として用いる。

2.1.3 大域特徴量 (Histograms of Oriented Gradients)

SIFT 特徴は局所的な領域での特徴量であり、大域的な情報が用いられていない。そこで大域的な形状を表現できる HOG 特徴を用いる。HOG 特徴は、SIFT のように局所領域の輝度の勾配方向をヒストグラム化した特徴量であるが、SIFT は特徴点に対して特徴量を記述するのに対し、HOG は大域的領域に対して特徴量を記述する。そのため、図 2 に示すように大まかな物体形状を表現することができる。

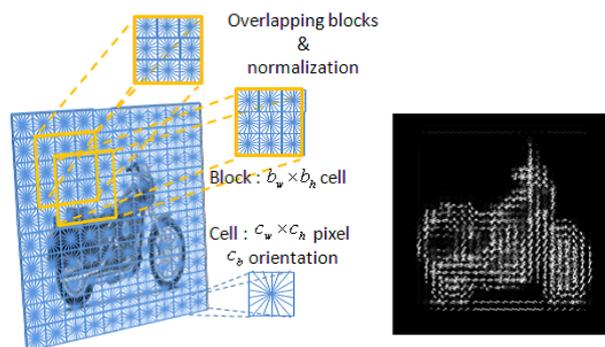


図 2 Histograms of Oriented Gradients

各ピクセルの勾配強度 $m(x, y)$ と勾配方向 $\theta(x, y)$ は以下の式で求められる。

$$m(x, y) = \sqrt{f_x(x, y)^2 + f_y(x, y)^2} \quad (1)$$

$$\theta(x, y) = \tan^{-1} \frac{f_y(x, y)}{f_x(x, y)} \quad (2)$$

$$\begin{cases} f_x(x, y) = L(x+1, y) - L(x-1, y) \\ f_y(x, y) = L(x, y+1) - L(x, y-1) \end{cases} \quad (3)$$

輝度勾配方向は次のように定義される。

$$\tilde{\theta}(x, y) = \begin{cases} \theta(x, y) + \pi, & \text{if } \theta(x, y) < 0 \\ \theta(x, y), & \text{otherwise} \end{cases} \quad (4)$$

算出された輝度勾配画像を、図2に示すようにセルと呼ばれる $c_w \times c_h$ 画素からなる小領域に分割する。それぞれの領域において、算出された勾配強度 $m(x, y)$ と勾配方向 $\theta(x, y)$ からヒストグラムを作成する。 $\theta(x, y)$ の角を c_b 方向になるよう量子化し、各方向に $m(x, y)$ を重みとして与える。すなわち、1セルあたり c_b 方向の勾配方向ヒストグラムができる。最後に、 $b_w \times b_h$ セルで構成されるブロックと呼ばれる領域を設定する。1セルずつオーバーラップさせながら正規化する。1セルあたり c_b 方向の特徴を持っているため、1ブロックあたりの特徴次元数は $d_b = b_w \times b_h \times c_b$ となる。あるブロックの特徴ベクトルを \mathbf{v} 、ブロック内で位置 (i, j) 、 $\{1 \leq i \leq b_w, 1 \leq j \leq b_h\}$ にあるセルのヒストグラムを h_{ij} としたとき、次式により正規化を行う。この処理を1セルずつオーバーラップさせながら実行する。

$$h'_{ij} = \frac{h_{ij}}{\sqrt{\|\mathbf{v}\|_2^2 + \epsilon}} \quad (\epsilon = 1) \quad (5)$$

2.2 特徴重要度の抽出による重み付け

識別的な特徴量に重みを付けるために、各カテゴリにおいて各特徴量の重要度を算出する、特徴重要度を用いることにより、特定のカテゴリのみで類似性の高い特徴量の重要度を上げたり、特徴量とカテゴリの共起関係を用いることができる。本章では、tf-idf、相互情報量、カイ二乗値という言語処理分野で用いられている3つの理論を、ヒストグラムインターセクションを特徴の頻出度として用いて画像分野に応用し、各カテゴリごとの重要な特徴量を求める手法について述べる。

2.2.1 TF-IDF

tf-idfの理論[5]では、tf(特徴量の出現頻度)とidf(逆出現頻度)の2つの指標から、そのカテゴリを認識する上でより識別的な特徴量に対し大きな重みを与える。カテゴリ c における特徴量 f の重要度 $tfidf_f^c$ は、次の式で求められる。

$$tfidf_f^c = tf_f^c \cdot idf_f^c \quad (6)$$

$$tf_f^c = \frac{\sum_{n=1}^{N_c} HI_f^c(n)}{N_c} \quad (7)$$

$$= \frac{\sum_{n=1}^{N_c} \sum_{i=1}^{DIM_f} \min(h_i^f(n), h_i^f(\bar{n}))}{N_c \cdot DIM_f} \quad (8)$$

$$idf_f^c = \log \frac{D}{d_f^c} \quad (9)$$

ここで、 N_c はカテゴリ c 中の学習画像枚数、 DIM_f は特徴量 f の次元数である。 $h_i^f(n)$ はカテゴリ c の画像 n 、特徴量 f のヒストグラムにおける i 番目の次元の値、 $h_i^f(\bar{n})$ は画像 n 以外に対する特徴量 f のヒストグラムの値であり、また、 D は画像データ総数、 d_f^c は特徴量 f のヒストグラムインターセクションがカテゴリ c 内の平均ヒストグラムインターセクション値 θ 以上の画像枚数である。そのカテゴリ内でよく類似しているほど tf_f^c は大きな値となる。また、そのカテゴリと他のカテゴリとの類似性が低いほど idf_f^c は大きな値となる。この重み付けにより、識別性の高い特徴量ほど大きな重みとなる。

2.2.2 相互情報量

相互情報量[6]は、2つの要素が共起して現れやすい度合いを統計的に数値化したものである。特徴量 f とカテゴリ c の共起度の尺度として、この相互情報量 $I(f, c)$ を用いる。本稿では、 $I(f, c)$ を以下のように定義する。

$$I(f, c) = H(f) + H(c) - H(f, c) \quad (10)$$

$$= -\log P(f) - \log P(c) + \log P(f, c) \quad (11)$$

$$= \log \frac{P(f, c)}{P(f)P(c)} \quad (12)$$

ここで、

$$P(f) = \frac{\sum_{c=1}^K \left[\frac{\sum_{n=1}^{N_c} \sum_{i=1}^{DIM_f} \min(h_i^f(n), h_i^f(\bar{n}))}{N_c \cdot DIM_f} \right]}{\sum_{f=1}^L \sum_{c=1}^K \left[\frac{\sum_{n=1}^{N_c} \sum_{i=1}^{DIM_f} \min(h_i^f(n), h_i^f(\bar{n}))}{N_c \cdot DIM_f} \right]} \quad (13)$$

$$P(c) = \frac{\sum_{f=1}^L \left[\frac{\sum_{n=1}^{N_c} \sum_{i=1}^{DIM_f} \min(h_i^f(n), h_i^f(\bar{n}))}{N_c \cdot DIM_f} \right]}{\sum_{f=1}^L \sum_{c=1}^K \left[\frac{\sum_{n=1}^{N_c} \sum_{i=1}^{DIM_f} \min(h_i^f(n), h_i^f(\bar{n}))}{N_c \cdot DIM_f} \right]} \quad (14)$$

$$P(f, c) = \frac{\sum_{n=1}^{N_c} \sum_{i=1}^{DIM_f} \min(h_i^f(n), h_i^f(\bar{n}))}{\sum_{f=1}^L \sum_{c=1}^K \left[\frac{\sum_{n=1}^{N_c} \sum_{i=1}^{DIM_f} \min(h_i^f(n), h_i^f(\bar{n}))}{N_c \cdot DIM_f} \right]} \quad (15)$$

ここで、 K はカテゴリの数であり、 L は特徴の数である。以上のように、特徴量 f とカテゴリ c の相互情報量を用いて重み付けを行う。相互情報量を用いることで、各カテゴリと共起する特徴量ほど大きな重みが得られる。

2.2.3 χ 二乗値

2組の定性的変数間の関連の強さの有無を統計的に検討するための手法として χ 二乗値による手法[7]がある。本稿では、表1のような 2×2 分割表を各特徴量、カテゴリ毎に作成し、特徴量とカテゴリ間の関連性を調べる。

2×2 分割表の場合の χ^2 値は、以下の式により求められる。

表 1 分割表

	Category	!Category
Feature	O_{fc}	$O_{f\bar{c}}$
!Feature	$O_{\bar{f}c}$	$O_{\bar{f}\bar{c}}$

$$\chi_{fc}^2 = \frac{N(O_{fc}O_{\bar{f}\bar{c}} - O_{f\bar{c}}O_{\bar{f}c})^2}{(O_{fc} + O_{f\bar{c}})(O_{fc} + O_{\bar{f}c})(O_{\bar{f}\bar{c}} + O_{\bar{f}c})(O_{\bar{f}c} + O_{\bar{f}\bar{c}})} \quad (16)$$

ここで、 O_{fc} はカテゴリ c における特徴量 f の平均ヒストグラムインターセクションを表し、 \bar{f} は f 以外の特徴量、 \bar{c} は c 以外のカテゴリを表している。 N はその総数である。特徴量 f とカテゴリ c の関連の強さを表す χ^2 値を重みとして用いることで、特徴量とカテゴリの関連度が強いほど大きな重みとなる。

2.3 認識

本章では、前章で述べた特徴重要度抽出法により得られた重みを用いて認識結果を統合する手法について述べる。認識の一連の流れを図 3 に示す。まず、学習データから bag-of-feature, Color, HOG の各特徴を抽出し、各特徴ごとに SVM(Support Vector Machine) を用意し境界面を学習する。同時に、tf-idf, 相互情報量, χ 二乗値により各特徴の重みを算出しておく。入力画像が与えられると、同様にまず各特徴を抽出する。そして各特徴量ごとに学習された SVM を用いて各カテゴリの SVM-Score を算出する。この SVM-Score S は次の式により 0 から 1 に正規化する。特徴量 f の SVM によって算出されるカテゴリ c の正規化されたスコア S_f^c は、

$$S_f^c = \frac{1}{1 + \exp(-S_f^c)} \quad (17)$$

となる。このように算出されたスコア S_f^c に対し、tf-idf によって得られた重み $tfidf_f^c$, 相互情報量によって得られた重み $I(f, c)$, χ 二乗値によって得られた重み χ^2 をそれぞれ以下の式のように統合し、最終的に最も高い値が得られたカテゴリを認識結果 c' として出力する。

$$c' = \arg \max_c \sum_{f=1}^3 W_f^c \cdot S_f^c \quad (18)$$

ここで、 W_f^c はそれぞれの重みを表している。

3. 実験

3.1 実験条件

実験はデータベースに Caltech101 [8] を用いた。その中から 10 カテゴリをランダムに選び、それを 1 データセットとして 3 セットのデータセットを作り、各重み付け手法による認識率の比較を行った。学習データは 1 カテゴリあたり 10 枚である。テストデータは 100 枚で各カテゴリ 10 枚の画像を含んでいる。学習データから bag-of-features, Color 特徴, HOG 特徴を抽出し、

Multiclass-SVM を用いて、各カテゴリの画像を正データとして学習し、他カテゴリのデータを負データとして各特徴ごとに学習する。SVM により得られる SVM-Score に対し、各重み付け手法により得られた重みを統合した結果が認識結果となる。

評価実験は、まず bag-of-features, Color, HOG のそれぞれ単独で認識を行う。次に、3 つの特徴を同じ重みで統合したものと、更に識別的な特徴に重み付けを行い統合したもので実験を行い、その精度を比較する。最後に、各カテゴリへの重み付けの効果を検証する。

3.2 実験結果

まず、bag-of-features, Color 特徴, HOG それぞれ単独で認識実験を行った結果を表 2 に示す。表に示すように HOG 特徴が最も良い結果となった。HOG は回転・スケール変化に非不変であるが、今回用いた Caltech のデータベースはカテゴリごとに向きやスケールが揃っていたため、大域的な形状情報による識別がうまくいったものと考えられる。

表 2 各特徴量による認識率 (%)

特徴量	BoF	Color	HOG
DatasetA	42	32	57
DatasetB	46	28	54
DatasetC	51	35	60

次に、従来の bag-of-features だけを用いた手法と、3 特徴を用いた手法、3 特徴に更に重み付けを行った手法の認識結果の比較を表 3 にまとめる。表を見ると、複数の特徴量を用いて重み付けを行うことにより、分類精度が向上していることが分かる。全体的には相互情報量の重み付け手法が最も良い結果となった。DatasetC において、一部精度が下がったものがみられるが、これは適切な重みが学習段階で得られなかったためと考えられる。適切な重みが得られなかった原因としては、実験に用いた画像には多くの背景が含まれており、これらから得られた特徴が特徴重要度の抽出に悪影響を与えたものと考えられる。

表 3 実験結果 (%)

特徴量	(従来手法)	(提案手法)			
	BoF	BoF+Color+HOG			
重み	-	同重み	TF-IDF	MI	Chi-square
DatasetA	42	47	66	65	60
DatasetB	46	48	54	59	52
DatasetC	51	59	60	62	56

最後に、最も精度の高かった DatasetA において、tf-idf により得られた重みを図 4 に示す。またその重み付けによる認識精度比較を図 5 に示す。ここで、”3feat/wo”は 3 つの特徴に同じ重みを付けて認識実験を行った結果で

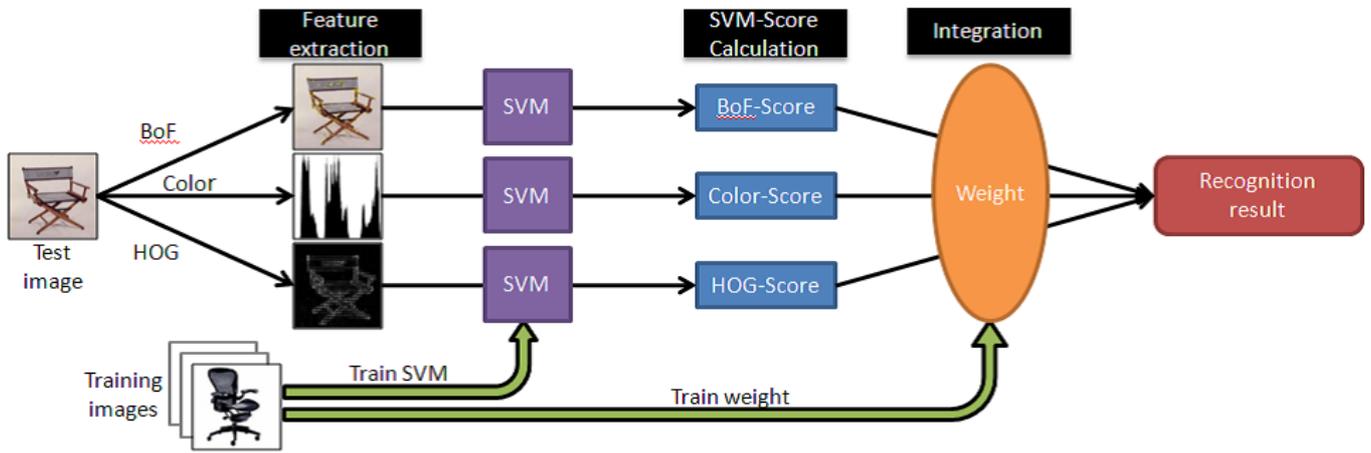


図 3 認識の流れ

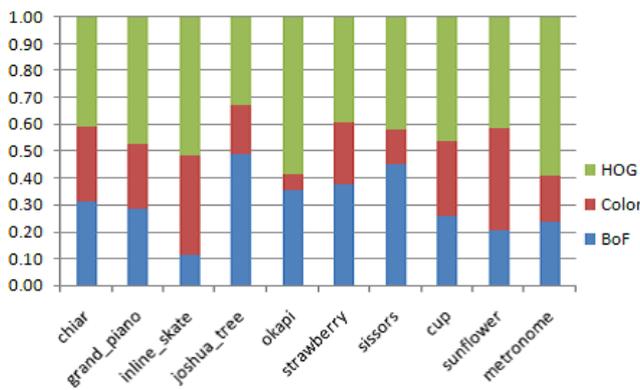


図 4 TF-IDF による重み

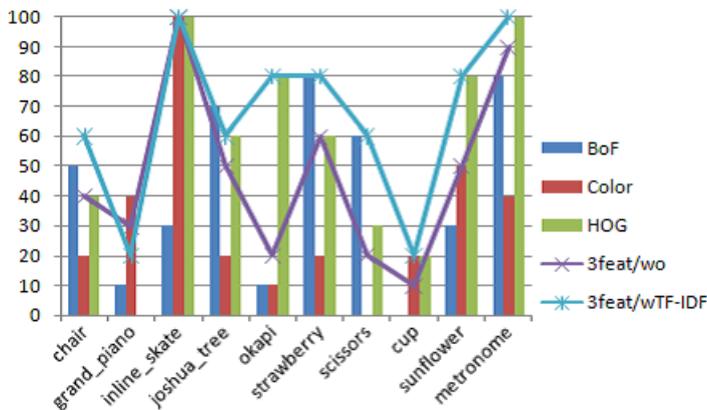


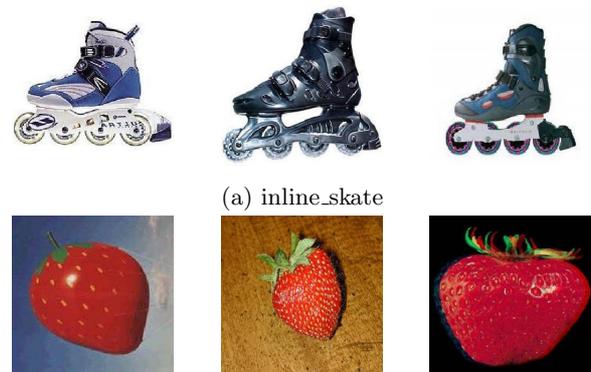
図 5 認識率の比較

ある。”3feat/wTF-IDF”は TF-IDF によって重み付けを行って認識実験を行った結果であり、同じ重みを用いた場合よりも高い値が取られていることが分かる。これは、カテゴリ毎に識別に重要な特徴量を予め学習しておくことにより、識別精度の高い特徴量が認識の際に優先されたためである。これにより、例えばカテゴリ”okapi”では、bag-of-features だけでは認識率は 10% 程度だったが、特徴量の重みを用いることで 80% まで向上していることが分かる。まだカテゴリ毎に認識率にばらつきはあるが、bag-of-features だけでは認識の難しかったカテゴ

リに対しても、複数の特徴量で補うことで認識できるようになった。また、カテゴリ毎に識別に重要な特徴量を学習しておくことで、各カテゴリの識別精度の高い特徴量が優先され、さらに認識率を向上させることができた。

3.3 考察

実験の結果、識別的な特徴量へ重み付けを行うことで全体的に認識率が向上した。今後の課題としては、背景の問題への対処が挙げられる。図 6 に実験に用いた学習データの一例を示す。



(a) inline_skate
(b) strawberry
図 6 学習データの例

(a) の inline_skate は、形状がとてもよく似ているため、大まかな形状情報が重要であり、色情報については種類によって異なるためあまり重要ではないと思われる。しかし、図 4 の重みを見てみると、形状情報である HOG は高い重みが得られているが、色情報も高い値になってしまっている。これは、このカテゴリの学習データが全て背景が白色で統一されていたため、色情報のカテゴリ内類似度が高くなり、重みが高くなってしまったものと考えられる。

(b) の strawberry は、逆に色情報が重要なカテゴリと考えられるが、このカテゴリの学習データは背景の色が様々で、色情報の重みが小さくなってしまっている。

重み付けとしては、図5に示すように、背景の混じったデータで学習した識別器で色情報をもとに識別を行ったものは認識率が低かったため、小さな重みが得られたという結果はむしろ良いのではあるが、真に物体上の識別に有効な特徴を学習したとは言えない。

そこで今後は、Saliency Map [9] を用いて画像中の視覚注意を引く領域を抽出し、それを物体の事前確率として、Graph Cuts [10] により物体領域を自動で抽出し、背景の影響のない状態で学習・認識する手法を検討中である。

4. ま と め

本稿では、複数の特徴量の重み付け統合による一般物体認識手法について述べた。SIFT 特徴に基づく bag-of-features 用いた手法において、輝度変化の少ない物体などカテゴリによって認識率が極端に下がる問題があったが、SIFT では用いていない情報を組み合わせることで認識率が向上した。さらに、主に言語処理分野において用いられている重要特徴抽出理論を画像特徴量のヒストグラムインターセクションを特徴の頻出度として画像分野に適用した。これにより、カテゴリ毎に識別に有効な特徴量を重視して認識が可能となった。実験の結果、bag-of-features だけでは認識の難しかったカテゴリの精度が上がり、さらに重み付けを行ったことにより精度が向上し、本手法の有効性が確認できた。今後は、Saliency Map を用いて物体領域を自動的に抽出し、背景の除去を行った上で物体領域から識別的な特徴量を学習し、認識精度の向上を図る予定である。

文 献

- [1] G. Csurka, C.R. Dance, L. Fan, J. Willamowski, C. Bray, Visual categorization with bags of keypoints, Proc. of ECCV Workshop on Statistical Learning in Computer Vision, 1-22, 2004.
- [2] C.D. Manning, H. SchFutze, Foundation of Statistical Natural Language Processing, The MIT Press, 1999.
- [3] D.G. Lowe, Distinctive image features from scaleinvariant keypoints, Journal of Computer Vision, vol.60, 2, 91-110, 2004.
- [4] N. Dalal, B. Triggs Histograms of oriented gradients for human detection, Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 886-893, 2005.
- [5] G. Salton, C. Buckley, Term-weighting approaches in automatic text retrieval, Information Processing & Management, 24, 5, 513-523, 1988.
- [6] H. C. Peng, F. Long, C. Ding, Feature selection based on mutual information: criteria of max-dependency, max-relevance, and min-redundancy, IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 27, No. 8, pp.1226-1238, 2005.
- [7] C. D. Manning, H. Schutze, Foundations of Statistical Natural Language Processing, MIT Press, 1999.
- [8] Caltech 101, http://www.vision.caltech.edu/Image_Datasets/Caltech101/
- [9] L. Itti, C. Koch and E. Niebur, A Model of Saliency-based Visual Attention for Rapid Scene Analysis,

- IEEE Trans. Pattern Analysis and Machine Intelligence, Vol.20, No.11, pp.1254-1259, 1998.
- [10] Y. Boykov, M. P. Jolly, Interactive graph cuts for optimal boundary & region segmentation of objects in N-D images, In IEEE International Conference on Computer Vision and Pattern Recognition, Vol.2, pp.731-738, 2004.