

多種類特徴統合による動作認識手法の提案

野口 顕嗣[†] 柳井 啓司[†]

[†] 電気通信大学 〒182-8585 東京都調布市調布ヶ丘 1-5-1

E-mail: †{noguchi-a,yanai}@mm.cs.uec.ac.jp

あらまし 映像中の動作を認識することは、様々なアプリケーションに応用することが可能で、非常に意義のあることである。動作認識は多くの分野で扱われてきた研究だが、Youtube 動画のような制限のない環境における認識を行った研究は少ない。そこで本研究では Multiple Kernel Learning(MKL) に基づく特徴統合による動作認識フレームワークを提案し Web 動画における動作認識を行った。統合に利用した特徴は時空間特徴、視覚特徴、動き特徴である。時空間特徴は、点の周辺パターンとその点の微小時間の動きを特徴化したものを利用した。実験には KTH データセットと、二種類の Youtube データセットの計 3 種類のデータセットにおける分類を行った。結果として KTH データセットで最新手法に匹敵する分類率の 94.0%、Youtube データで最新手法を大幅に上回る分類率である 80.4% という結果が得られた。

キーワード 時空間特徴, 動作認識, KTH, Youtube data, 特徴統合, MKL

Action Recognition by Combining Features Using Multiple Kernel Learning

Akitsuugu NOGUCHI[†] and Keiji YANAI[†]

[†] The University of Electro-Communication, Choufugaoka 1-5-1, Choufu, Tokyo, 182-8585 Japan

E-mail: †{noguchi-a,yanai}@mm.cs.uec.ac.jp

Abstract Recognizing action in video is very useful, because it can be applied to many kinds of applications. Action recognition is being studied in many researches, but most research has focused on videos taken in controlled environment. In this paper, we propose action recognition method by combining features based on Multiple Kernel Learning(MKL), and we recognize action in web video shots. We combine visual, motion and spatio-temporal features. As spatio-temporal feature, we utilize features which consists point and its local tracking. In experiment, we evaluate our method using KTH dataset, and Youtube dataset. As a result, we obtain 94.0% as a classification rate for in KTH dataset which is almost equivalent to state-of-art, and 80.4% for Youtube dataset which ourperforms state-of-the-art much.

Key words spatio-temporal feature, action recognition, KTH, Youtube data, combining featrues, MKL

1. ま え が き

近年 Web 上の動画の数は爆発的に増えてきている。しかし一方でユーザが見たいシーンを探すことが非常に困難な問題となってしまう。現状の動画検索システムはタグなどによるテキストベースな手法が用いられているが、これはユーザの主観によってつけられるもので、タグのみで動画を特定することは非常に難しい。そのため動画の内容を解析するコンテンツベースな検索手法が今後求められてくると考えられる。

それを行うために動画のなかで行われている動作を解析することは極めて意義のあることである。さらに動作を認識することは動画検索だけではなく、サーベイランス、動画インデキシング、動画要約など様々なアプリケーションに応用することが

可能である。

現在の動作認識の研究の多くは、「カメラは固定」、「動作を行う人は一人のみ」のような制限のある環境で行われているが、実際の動作はそのような単純なものではない。当然、カメラが固定ではなく、カメラの動きであるカメラモーションが存在する状況下での動画も現実には多数存在し、同一映像中で同時に複数の異なる動作が行われることもある。更に、同一の動作であっても、撮影される角度によってその特徴は異なったものになる場合もある。

そこで本研究では多種類特徴統合による動作認識フレームワークを提案する。特徴統合には時空間特徴、視覚特徴、動き特徴の三種類の特徴を利用した

1.1 関連研究

近年、動画の解析のために時空間特徴が注目を集めている。時空間特徴とは、動画から抽出される特徴の一つで、動き情報と視覚情報を同時に表現することが可能な特徴である。

まず時空間特徴の抽出の考え方として、画像認識の三次元拡張の手法が挙げられる。Kobayashiらは自己相関特徴を三次元に拡張し、サーベイランスの分野に特化した cubic higher-order local auto-correlation(CHLAC)を提案した[1]。またLaptevらはハリス点検出を三次元に拡張した検出手法を提案している[2]。

次に主要な手法として、cuboidと呼ばれる立方体を抽出し、それを特徴化する手法がある。Dollarらは空間軸にガウシアンフィルタ、時間軸にガボールフィルタを適用することでこのcuboidを抽出する手法を提案した[3]。またcuboidの記述子として、LaptevやDollarらはHistogram of Orient Gradient(HoG)やHistogram of Orient Flow(HoF)で表現することを提案している。一方でKläserらは、記述子として三次元的なHoGを利用することを提案している[4]。

しかしこのようなcuboid主体の手法は計算コストが高くなる傾向があり、本論文で行うような大量データを扱うのに向いていない。更にcuboidのサイズを求めることは非常に困難な問題である。そこで本研究では、特徴的な点と、その点の微小時間の動きを特徴化する新たな時空間特徴を提案する。この特徴は計算時間が少なく、かつ正確な分類が期待できる。

Web動画に対する動作認識を行った研究は少ない。CinbisらはWeb上から動作を自動学習する手法を提案し、Youtubeデータを動作認識に利用している[5]。この研究では、まずquery wordをもとにして画像を収集し、その画像から特徴を抽出することで動作モデルを構築し、実際の映像において認識を行っている。しかしこの手法の学習データはWeb上から収集された静止画像であり、動作の記述も全て動画像ではなく静止画像ベースで行われている。本研究では画像の視覚特徴のみではなく、動き特徴も考慮して動作を分類することを考える。

最も本研究の手法と類似している手法として、Liuらの研究が挙げられる[6]。この手法は特徴量として[3]で提案された時空間特徴を、視覚特徴としてSIFT記述子[7]を利用し、Adaboostに基づき統合する手法を提案している。またこの研究では、Page Rankに基づく重要な特徴の選択を行っている。本研究では、動作を認識するために新たな時空間特徴を提案し、利用するが、この手法では既存の特徴をどのように扱うかに重点が置かれている。

2. 提案手法

2.1 時空間特徴抽出

本研究では、Web動画分類に適した新しい時空間特徴抽出手法を提案する。提案する手法は[8]を拡張したものである。[8]を含め、多くの動作認識の研究は、カメラモーションに対する対応がない。しかしWeb動画を分類するためには「カメラモーション」をどのように扱うかは非常に重要な問題となってくる。そこで本研究では、特徴を抽出する前にカメラモーションの検出を行う。

また既存の抽出手法は計算コストが高く、本研究のような大量のデータから抽出することに向いていない。本研究では、その問題を解決するために、点の周辺パターンと、その点の動きで動画を特徴化する手法を提案する。この手法は既存手法に比べ計算コストが低く、大量のデータから特徴を抽出することが可能である。

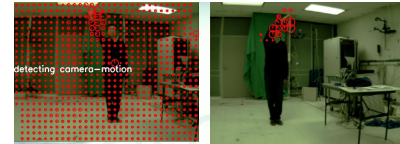


図1 カメラモーション検出例(左)、カメラモーションが検出されない例(右)

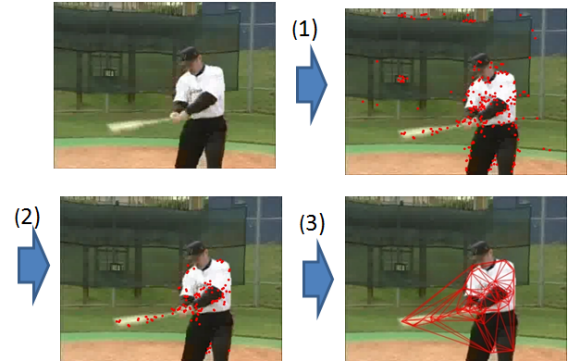


図2 視覚特徴抽出の様子

手順1に提案する特徴抽出手法の流れを示す。

手順1 時空間特徴抽出手法の流れ

step1	: カメラモーション検出
step2	: 時空間特徴における視覚特徴抽出
▷step2-1	: SURF抽出
▷step2-2	: 時空間特徴点の決定
▷step2-3	: Delaunay三角分割
step3	: 時空間特徴における動き特徴抽出
▷step3-1	: Lucas-Kanade法による動き抽出
▷step3-2	: SURFのdominant rotationによる方向の正規化
step4	: 視覚特徴ベクトルと動き特徴ベクトルの結合

本手法は大きく4つに分けることができる。まず、カメラモーション検出部(step 1)、二つ目に視覚特徴抽出部(step 2)、三つ目が動き特徴抽出部(step 3)、最後に特徴統合部(step 4)である。以下ではそれぞれのステップ毎に詳しく述べていく。

(Step 1) カメラモーション検出部 動画は撮影者の意図によって、ズームやパンなどのカメラモーションを含むことがある。しかしWeb動画においてのカメラモーションは手振れなどの、撮影者の意図しないカメラモーションが多く含まれている。またWeb動画は解像度が低い。これらのことから収集した動画の正確なカメラモーションを求めることは非常に難しい問題である。

Liuらの手法ではカメラモーションを検出した場合そのフレームを破棄することで対応をしている[6]。本研究でも、これを参考にし、カメラモーションを検出した場合、その特徴を破棄する。しかしカメラモーションを全て破棄しては、「動画全てにカメラモーションを含んでいるような動画から特徴が抽出されない」、「カメラモーション中に存在する重要な特徴を抽出できない」、などの問題点もある。

検出手法として、図1のようにグリッド上にLucas-Kanadeアルゴリズムに基づいて、動き情報を計算する。その中で動きのあった領域が一定以上だった場合カメラモーションの検出とする。

(Step 2) 時空間特徴における視覚特徴抽出部 図2は視覚

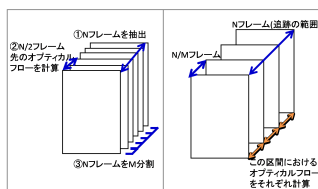


図3 全体の動き特徴抽出概要(左), 局所追跡の概要(右)

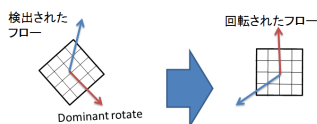


図4 オプティカルフローの回転

特徴の抽出の様子を示している．各ステップ毎に見ていくと，(1) フレーム画像から SURF を抽出する．本研究では動画としての特徴を抽出したいので，動きがない点は時空間特徴として適していない．(2) よってそれぞれの特徴点の動きを計算し，動きがなかった点を削除する(時空間特徴の決定)．(3) 最後に残った点について Delaunay 三角分割を行う．この時空間特徴を三角分割することは本手法独自のものである．これを行うことで，その点のみではなく隣接する特徴も考慮した特徴を構築することが可能となる．

視覚特徴には三角形の頂点を構成する三点の SURF 記述子を使用する．それぞれの三角形の頂点の点を SURF 記述子のスケール毎に整列させる．ただし同一スケールの特徴が存在した場合は， x 座標の最も小さい特徴点から見て右周りに特徴を配置する．SURF の次元数は 64 次元なので，結果として視覚特徴は $64 \times 3 = 192$ 次元で表現される．

(Step 3) 時空間特徴における動き特徴抽出部 動き情報として，本手法では三角形の頂点を構成する，それぞれの点の動きと，三角形の面積の変動を特徴化する．点の動きは視覚特徴の時と同様で SURF のスケールによって整列させる．

図3は動き特徴の抽出の概要を示している．まず図3(左)が時空間特徴点の決定の様子を示している．最初に SURF を抽出したフレームから N フレーム先までのフレームをとり，以降この区間から動き特徴を抽出する．次に，抽出された SURF の点に対して， $N/2$ フレームの動き情報を計算する．この動き情報に基づいて時空間特徴は決定される(手順1の step 2-2)．

次に動き特徴の抽出に関して説明する，図3(右)がその様子を示している．選択された N フレームを， M 分割する．そして分割したそれぞれの区間で Lucas-Kanade アルゴリズムによって，特徴点のオプティカルフローを計算する．ただし，インターバル区間 i の特徴点の座標 L_i は，前の区間 $i-1$ で推定された，動き情報によって求められる．この分割数 M が N の値に近いほど詳細な追跡が可能になり，逆に M が 1 に近くなると簡略な動き特徴が抽出される．

各区間の動き情報は x^+, x^-, y^+, y^- ，及び動きなし，の 5 次元で表現される．ただし x^+ はオプティカルフローの x 成分の正の方向を， x^- は x 成分の負の方向を示している．実験において特徴を抽出する区間 $N=5$ ，分割数 $M=5$ に設定しているので，それぞれの点における次元数は $(M-1) \times 5$ で 20 次元で，三角形の面積変化は 5 次元で表現される．よって動きの次元数は $20 \times 3 + 5 = 65$ 次元となる．

ただしこのままの計算では動き特徴は回転に関して敏感な特徴になってしまう．同じ「歩く」という動作においても，動作の方向によって異なる特徴が抽出されてしまう．そこで本研究では視覚特徴で得られる dominant rotation を利用し，オプティカルフローを回転させる．特徴は三点で一組になっているが，ここではそれぞれの dominant rotation を利用して，動き情報の回転を行う．その様子を図4に記載する．



図5 動き特徴が有効な動作の例(左)，視覚特徴が有効な動作の例(右)

特徴点の座標を (x_1, y_1) ，オプティカルフローが検出された座標を (x_2, y_2) とした，SURF の dominant rotation を θ とした場合，回転されたフローの座標 (x, y) は式1で定義されるものとなる．

$$\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} \cos\theta & -\sin\theta & x_2 \\ \sin\theta & \cos\theta & y_2 \end{bmatrix} \begin{bmatrix} x_1 - x_2 \\ y_1 - y_2 \\ 1 \end{bmatrix} \quad (1)$$

最後に二つの特徴を単純に結合することによって時空間特徴を構築する．

2.2 特徴統合

動作によって重要となってくる特徴は変わってくる．図5(左)のような“jogging”と“running”を視覚情報だけで判断することは難しいが，動き特徴を利用することで分類が可能となってくる．逆に図5(右)は boxing と clapping という二つの動作を示しているが，これは動き特徴だけで判断しても，共に腕を左右に振っているのが特徴的に類似したものになってしまう．しかし視覚特徴を用いることで，簡単に分類が出来る．

よって本研究では Multiple Kernel Learning(MKL) によって特徴の重みを自動で推定する手法を利用することで，これらのショットの分類を行った．統合には視覚，動き，時空間特徴の3種類の特徴を利用した．

2.2.1 Multiple Kernel Learning

本研究では動作認識を行うために Multiple Kernel Learning(MKL) を利用する．これは複数のサブカーネルを線形結合することで，新たな最適なカーネルを求める手法で，統合カーネルは以下の式で求められる．

$$K_{combined}(x, x') = \sum_{j=1}^K \beta_j k_j(x, x') \quad (2)$$

with $\beta_j \geq 0, \sum_{j=1}^K \beta_j = 1.$

各サブカーネルの重み β_j の設定によって，統合された新たなカーネルの出力は変化する．そのためこの最適な重みをどのように求めるかが問題となり，この問題は MKL 問題と呼ばれている[9]．この MKL 問題は全ての β_j の組み合わせを実際に試すことで解くことが出来る．しかし，特徴やカーネルの数が増えるにつれ， β_j の組み合わせも爆発的に増えてしまう．実時間内にこの問題を解決するために，近年 MKL 問題を凸面最適化問題として解く手法が提案されている[10]．Sonnenburg らは単一カーネルの SVM 学習を反復することによって最適なカーネルの重み β_j を求める手法を提案している[10]．

2.2.2 特徴抽出

以下では統合に利用する動き特徴，視覚特徴，及び特徴の表現手法について述べる．

動き特徴 時空間特徴は，動き情報も内包しているが，それだけでは局所的な点の動きしか表現できない．よって本研究ではフレーム全体から動き特徴を抽出する．この動き特徴は，全体的な動きを表現出来るので，時空間特徴点の局所点の動き情



図 6 抽出された動き特徴



図 7 視覚特徴の抽出

表 1 特徴表現手法

特徴	次元数	表現手法	コードブックサイズ
時空間特徴	257	BoSTF	5000
視覚特徴	24	BoK	5000
動き特徴	56	BoFr	3000



図 8 Our Youtube データセット

報とは、異なる識別能力が期待できる。

動き特徴として、本研究ではグリッド点におけるオプティカルフローを利用する。それぞれの検出されたフローは、8方向、7段階の大きさからなるヒストグラムに投票されていく。図6は実際に抽出された動き特徴を示している。

視覚特徴 視覚特徴には6方向、4周期のガボール特徴を使用する。ただし、フレーム画像を 20×20 の局所領域に分け、それぞれの領域から特徴を抽出し、それぞれのパッチをBag-of-Keypoints 表現で認識に利用する。

この二つの特徴を時空間特徴と統合することで動作認識を行う。しかしフレーム一枚で認識しようとした場合、選ばれたフレームによって結果が大きく変わってしまうことがある。特に動き特徴の場合この傾向が顕著に現れる。そして、最も重要なキーフレームを選択することは非常に難しい問題である。よって本研究では bag of frames 表現を導入することで、この問題を解決する。

特徴表現 表1に特徴の表現手法についてまとめた。抽出された257次元の時空間特徴はBoSTFで表現される。視覚特徴は全てのフレームから抽出される、24次元のパッチを bag of keypoints(BoK)で表現する。時空間特徴、視覚特徴におけるコードブックサイズは5000に設定した。

動き特徴はフレーム毎に抽出される、57次元の特徴を用いて bag of frames(BoFr)で表現する。これはフレームから抽出された動き特徴をベクトル量子化して、その特徴の出現頻度で動画を表現する手法である。ただしこの時のコードブックサイズは3000に設定した。

カメラモーションが検出された場合、時空間特徴と動き特徴では特徴が破棄されるが、視覚特徴は抽出が行われる。ショット全体がカメラモーションを含んでいるショットでは、時空間特徴や動き特徴が検出されない場合がある。この場合、時空間特徴、動き特徴のベクトルは0ベクトルで表される。



図 9 Wild Youtube データセット

表 2 KTH の分類結果

	walking	jogging	running	boxing	waving	clapping
walking	0.98	0.02	0	0	0	0
jogging	0.07	0.84	0.09	0	0	0
running	0	0.17	0.83	0	0	0
boxing	0	0	0	0.92	0	0.08
waving	0	0	0	0	0.98	0.02
clapping	0	0	0	0.05	0.04	0.91

時空間特徴 (91.0%)

	walking	jogging	running	boxing	waving	clapping
walking	0.99	0.01	0	0	0	0
jogging	0.03	0.94	0.03	0	0	0
running	0.01	0.14	0.85	0	0	0
boxing	0.01	0	0	0.96	0.01	0.02
waving	0	0	0	0	0.98	0.02
clapping	0	0	0	0.05	0.02	0.93

MKL(94.0%)

3. 実験

3.1 データセット

動作認識には標準データセットである KTH データセット、本研究で構築した Our Youtube データセット、Liu らが構築した Wild Youtube データセットの三種類を使用した。

KTH データセットには「カメラモーションを含まない」「動作をしている人は一人だけ」など様々な制約が存在する。しかし Web 動画のような実世界の映像はそのような制約は存在せず、そのような制約のない動画を分類することはより難しい問題となってくる。

そこで本研究では Youtube から収集したショットから独自の Youtube データセットを構築した(図8)。このデータセットは“batting”, “running”, “walking”, “shoot”, “jumping”, “eating”の動作から成り、カメラモーション、視点変化、背景ノイズなどを含む、より制約の少ないデータセットとなっている。

次に Web 動画分類における他の手法と比較するために、Liu らが構築した、データセット(図9)における動作分類を試みた。このデータセットは11の動作(“basketball shooting”, “volleyball spiking”, “trampoline jumping”, “soccer juggling”, “horse riding”, “cycling”, “diving”, “swinging”, “golf swinging”, “tennis swinging”, “walking_with_dog”)を含むデータセットで、Our Youtube データセットと同様に様々なノイズを含む挑戦的なデータセットである。本論文ではこのデータセットを Wild Youtube データセットと呼ぶ。

4. 実験結果

本論文では多種類特徴の統合による動作認識フレームワークについて提案をした。ここでは3つのデータセットを使用することで、提案手法の有効性を実証する。

4.1 KTH データセット

KTH データセットは動作認識の研究において最も広く利用されているデータセットである。本研究でもこのデータセットを使用して、5-fold cross validation で分類を行うことで特徴統合の有用性を示し、最新手法との比較を行う。

表2は時空間特徴と MKL の混合行列を示している。この結果から“running”と“jogging”を区別することが難しいことが

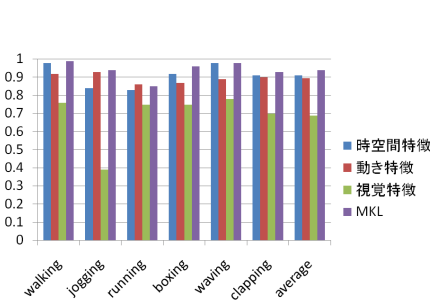


図 10 KTH データセットによる分類結果

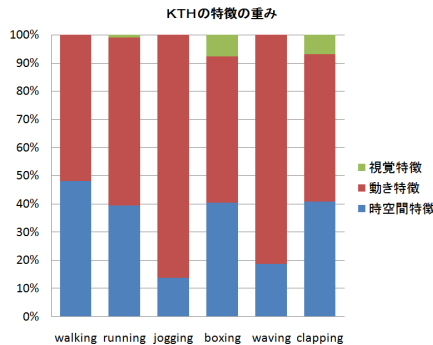


図 11 KTH データセットの特徴毎の重み

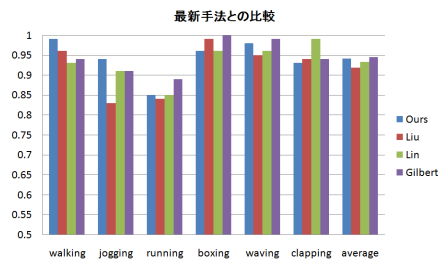


図 12 KTH における最新手法との比較

分かる．図 10 は KTH データにおける動作ごとの分類率を示している．特徴を統合しない場合，最も良かったものは時空間特徴の 91.0%であったが，特徴を統合することで 94.0%まで精度が向上した．

特徴統合の重みは図 11 に示されている．KTH データセットは視覚的变化に乏しいので，視覚特徴に重みはほとんど割り当てられていない．一方で動き特徴は“running”，“jogging”，“waving”で重みの値が高くなっている．実際“running”と“jogging”を区別するためには動き特徴が重要であることを示している．

図 12 に最新手法との比較を示している．比較に用いた手法は Liu らの研究 [6]，Lin らの研究 [11]，Gilbert らの研究 [12] の 3 つである．本提案手法が 94.0%，Liu らの手法で 91.8%，Lin らは 93.3%，Gilbert らは 94.5%の分類率であった．これらのことから，提案した分類手法は非常に有用であることが分かる．

動作ごとに見ていくと，全ての研究で“running”と“jogging”を分類することは難しい問題であることが分かる．特に“running”の分類においては分類率は極端に下がっている研究が多い．本研究は“walking”の精度が他より特に高く，他の動作においても，同精度の結果を残すことが出来ている．

4.2 Youtube データセット

実際の映像は KTH の映像のように単純ではない．そこで本論文では Youtube から集められたデータを利用して，より実用的な映像における動作分類を行う．

本研究では 2 種類の Youtube データセットで動作分類を行った．一つ目が独自に構築した Our Youtube dataset．もう一つが Liu らが構築した Wild Youtube dataset である．

a) Our Youtube dataset

表 3 は 5-fold cross validation で分類したときの時空間特徴と MKL の混合行列を示している．時空間特徴のみでは“running”と“walking”を混同してしまっているが，MKL で特徴を統合することにより，比較的分類が出来るようになっている．

図 13 は動作ごとの分類結果を示している．最も良かった視覚特徴で 88.1%であったが，MKL で特徴を統合することで 94.5%まで向上した．統合する際の特徴の重みは図 14 に示されているが，“batting”，“shoot”，“walking”などのその動作が発生する場面が限定される動作においては，視覚特徴の重みが比較的高くなる傾向があった．

b) Wild Youtube dataset

次に Wild Youtube データセットにおいて動作分類を行い，他手法との比較を行う．ただし分類は 5-fold cross validation で行った．

表 4 に各特徴の混合行列を記載した．全体的に basket_shooting の結果が良くない．これはこの動作が多くのカメラ

表 3 Our Youtube データの分類結果

	batting	running	walking	shoot	jumping	eating
batting	0.69	0.01	0.05	0.08	0.02	0.15
running	0.02	0.89	0.08	0	0	0.01
walking	0.03	0.12	0.68	0.03	0	0.14
shoot	0.05	0.01	0.05	0.82	0.05	0.02
jumping	0.04	0	0.01	0.06	0.85	0.04
eating	0.06	0.04	0.09	0.01	0.01	0.79

時空間特徴 (77.8%)

MKL(94.5%)

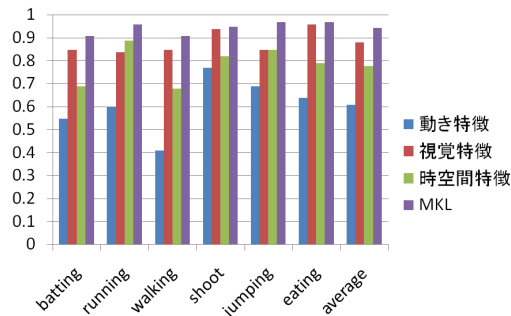


図 13 Our Youtube データセットによる分類結果

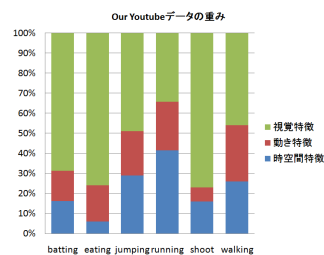


図 14 Our Youtube データセットの特徴毎の重み

ラモーションを含んでいるためと考えられる．本手法ではカメラモーションは除去してしまうので動画全てがカメラモーションの場合は特徴を抽出することが出来ないという問題がある．

このデータセットにおいても特徴を統合することにより，時空間特徴のみでは 63.4%，視覚特徴で 69.2%の分類率であったが，80.4%まで向上した．その際の重みは図 15 に示すようになった．全体的に視覚特徴が大きな重みを得ていることがわかる．特に basket_shoot はそのカメラモーションにより，動き特

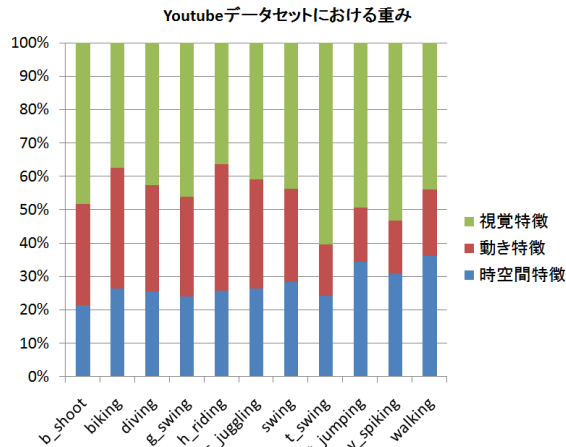


図 15 wild Youtube データセットの特徴毎の重み

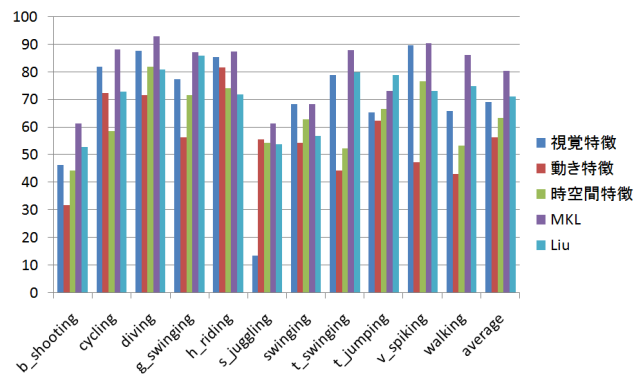


図 16 wild Youtube データセットの他手法との比較

表 4 Wild Youtube データの分類結果

	b_shooting	cycling	diving	g_swimming	h_riding	s_juggling	swinging	t_swinging	t_jumping	v_spiking	walking
b_shooting	61.4	1.4	9.3	14.3	2.1	5.7	0.7	1.4	0	3.6	0
cycling	0.88	2.2	0.14	7.1	1.1	11.6	2.2	2.2	1.4	5.8	0
diving	2.7	0	81.9	3.4	1.3	3.4	2	2	0	0.7	2.7
g_swimming	8.7	0	14	71.7	0.7	5.1	2.2	6.5	0	2.2	1.4
h_riding	2.1	5.3	2.1	0.5	74.2	1.6	1.1	2.1	0.5	2.6	7.9
s_juggling	2	0	7.3	5.3	1.9	54.3	4	8.6	10.6	6	0.7
swinging	2.3	12.9	1.5	1.5	0.8	0.8	62.8	0	7.6	1.5	8.3
t_swinging	12.5	9	3.8	10.6	4.4	5.6	0.6	52.5	0	4.4	0.6
t_jumping	0	1.8	0.9	2.6	3.5	12.3	7.8	1.8	66.7	0	2.6
v_spiking	6.2	2.7	2.7	1.8	1.8	1.8	1.8	3.6	0	76.8	0.9
walking	4.2	1.1	4.2	4.2	12.7	0.8	0.8	1.7	1.7	5.1	53.4

時空間特徴 (63.4%)

MKL(80.4%)

特徴による分類が困難になるため、視覚特徴に大きな重みがつけられている。

図 16 は特徴ごとの結果と、他手法との比較を示している。ここでは Liu の手法 [6] との比較である。結果として本手法の分類率が 80.4%であったのに対し、Liu らの手法は 71.1%と Liu らの手法を大きく上回っている。これらのことから MKL による特徴統合は Web 動画分類において有効であることが分かる。

5. おわりに

本研究では Web 動画分類のための時空間特徴抽出手法について提案し、その特徴を利用した Web 動画分類を行った。提案した時空間特徴は視覚的な局所特徴と、その点の動きを特徴化することで抽出を行った。また複数特徴を統合による動作認識フレームワークについての提案も行った。統合する特徴には、提案した時空間特徴、視覚特徴、動き特徴の 3 種類を用いた。

動作認識の実験として、KTH データセットの他に、二種類の Youtube データセットを使用した。その結果 KTH データセットでは 94%、Youtube データセットで 80.4%と最新手法を上回る精度であった。

時空間特徴の改良点として、現状ではカメラモーションを検出した場合、その特徴を排除することで対応してきた。しかし、それでは有用な特徴を抽出出来ない場合がある。よってカメラモーションをより精巧に取り扱うため、動き補正を入れることでモーションを検出した場合も正しい特徴を抽出するシステムを構築していく必要がある。

また提案手法では一つのショットから大量の特徴が抽出されてしまっている。より良い分類を行うために、有用な特徴を選択して抽出するようにならなければならない。

Youtube のような動画には複数の人間が動作を行っている。「歩いている人」の隣で「走っている人」がいるときもある。しかし現在のシステムでは、このように同時に起こる、異なる動作を検出することができない。このような動作を検出するために、人検出とトラッキングに基づく、抽出手法を構築していくことも、精度を向上させるために有効なことである。

文 献

- [1] T. Kobayashi and N. Otsu. A three-way auto-correlation based approach to human identification by gait. In *Proc. of IEEE Workshop on Visual Surveillance*, pp. 185–192, 2006.
- [2] I. Laptev and T. Lindeberg. Local descriptors for spatio-temporal recognition. In *Proc. of IEEE International Conference on Computer Vision*, 2003.
- [3] P. Dollar, G. Cottrell, and S. Belongie. Behavior recognition via sparse spatio-temporal features. In *Proc. of Surveillance and Performance Evaluation of Tracking and Surveillance*, pp. 65–72, 2005.
- [4] A. Kläser, M. Marszałek, and C. Schmid. A spatio-temporal descriptor based on 3d-gradients. In *British Machine Vision Conference*, pp. 995–1004, sep 2008.
- [5] R. I. Cinbins, R. Cinbins, and S. Sclaroff. Learning action from the web. In *Proc. of IEEE International Conference on Computer Vision*, pp. 995–1002, 2009.
- [6] J. Liu, J. Luo, and M. Shah. Recognizing realistic action from videos. In *Proc. of IEEE Computer Vision and Pattern Recognition*, pp. 1–8, 2009.
- [7] D. Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, pp. 91–110, 2004.
- [8] A. Noguchi and K. Yanai. Extracting spatio-temporal local features considering consecutiveness of motions. In *Proc. of Asian Conference on Computer Vision (ACCV)*, 2009.
- [9] G. R. G. Lanckriet, N. Cristianini, P. Bartlett, L. E. Ghaoui, and M. I. Jordan. Learning the kernel matrix with semidefinite programming. *Journal of Machine Learning Research*, Vol. 5, pp. 27–72, 2004.
- [10] S. Sonnenburg, G. Rätsch, C. Schäfer, and B. Schölkopf. Large scale multiple kernel learning. *Journal of Machine Learning Research*, Vol. 7, pp. 1531–1565, 2006.
- [11] Z. Lin, Z. Jiang, and L. S. Davis. Recognizing action by shape-motion prototype trees. In *Proc. of IEEE International Conference on Computer Vision*, pp. 444–451, 2009.
- [12] A. Gilbert, J. Illingworth, and R. Bowden. Fast realistic multi-action recognition using mined dense spatio-temporal features. In *Proc. of IEEE International Conference on Computer Vision*, pp. 925–931, 2009.