

ジオタグ画像認識における周辺テキスト情報の有効性の検証

八重樫恵太^{†1} 柳井啓司^{†2}

本稿では、我々が提案した航空写真を用いる位置情報付き画像の認識に、撮影位置に関する周辺情報のテキストおよび時間情報を追加した場合の画像カテゴリ分類における分類精度を比較検討するとともに、それぞれの特徴量の有効性を Multiple Kernel Learning(MKL) による重みを以て推定する。結果として、28 種類のコンセプトについて、その平均適合率の平均値が、航空写真と画像の組み合わせの場合と比較して 75%から 80%へ向上した。

Verification of Effectiveness of Using Location-related Textual Information for Geotagged Image Recognition

KEITA YAEGASHI^{†1} and KEIJI YANAI^{†2}

For geotagged image recognition, we have already proposed the method to combine raw image features and features of corresponding aerial photos. In this paper, we propose adding time information and location-related information as features for geotagged image recognition. In the experiments, we evaluate not only recognition rates but feature fusion weights estimated by Multiple Kernel Learning (MKL). As a result, the mean average precision for 28 categories increased up to 80% by the proposed method, compared with 75% by the baseline.

1. はじめに

今日では、デジタルカメラの普及により、撮影位置の情報を持つ位置情報付き写真は WWW(World Wide Web) 上に大量に存在している。また、Flickr などの写真共有を行うソーシャルサイトの普及により、個人が多く位置情報付き写真を収集、整理し、共有する

ことが容易になった。一方で、Google Maps などのオンラインマッピングサービスも普及し、地図情報の検索機能と共に日々高度化している。地図情報から、デジタル写真や航空写真、Google Street View に代表されるような周辺写真情報へアクセスすることはネットユーザーにとってはや当然のこととなった。

利用者は自宅から地図を利用した位置の検索だけでなく、航空写真の閲覧など高度なサービスを利用できるようになった。また、写真の収集に関しても、位置情報を積極的に活用しようとする試みが行われてきた。その例として Flickr では、写真に位置情報を付加して投稿し、地図 (Flickr Map) を用いて写真を検索、整理できるようになっている。

大量の写真情報の普及にもかかわらず、それらを自動で整理・分類し、ユーザーの手間を省くことは未だ困難な課題であり続けている。単純な画像分類へのソリューションとして、現状ではタグ (内容を表現する複数の単語の集合) やタイトルなど、テキストベースのメタデータが主流になっている。写真を分類するための一般画像認識の基礎技術が高度化する中で、認識精度の向上を図るにあたり、写真と関連する多様な情報を、特に実世界と関連の深い情報をいかに効率的に組み合わせるかが求められる。

我々は、位置情報付き写真の一般画像認識において写真の撮影位置に対応する航空写真を付加的な画像特徴量として利用する研究を行ってきた^{9),10)}。画像認識において、位置情報の有効性が高い認識カテゴリと、そうでないカテゴリを明確に区別し、デジタル写真と実世界情報との対応付けにおける有効な手段を提案することを目標として位置付けている。

一般に、認識対象と位置は密接な関係があり、例えば、海岸は海と陸の境目にしかなく、海や陸の真ん中には存在し得ない。しかしながら、海岸の写真の認識において位置情報を役立てるには、世界中の海岸の位置を学習データに持つておかないといけない。これには、膨大なデータを用意する必要がある。そこで、我々は、写真の撮影位置の地理的な状況を表す情報源として、航空写真に加えて、位置の周辺を記述するテキスト情報に注目している。地図に掲載されているような、駅や建物の名前などのテキスト情報のことである。写真の画像特徴量に合わせて、航空写真から抽出した画像特徴量とテキストに由来する特徴量を認識に用いることで、写真の撮影場所の地理的なコンテキストを反映した認識が可能となると考えている。

本研究では、分類カテゴリによって、どの程度位置情報由来の特徴量の有効性に違いがあるか分析を行う。認識実験にあたり、Flickr から収集した位置情報付き写真と、それぞれの位置情報に対応する複数の縮尺 (レベル) の航空写真、そして周辺テキストから抽出した特徴量を用いる。各特徴量の有効性を明確に判定するにあたり、機械学習の段階においてマルチカーネル学習 (MKL) を用いて、画像から抽出した特徴量と、位置情報に基づく航空写真、周辺テキスト、緯度経度の各特徴量、さらに撮影時間についての特徴量の認識における重要度の重みを推定し、本稿で追加した周辺テキスト、緯度経度、撮影時間についての特徴の有効性を評価する。

本稿は以下のように構成される。まず実験の全体的手順と本研究の方針について第 3 節で触れる。後述するように、実験には画像収集、特徴抽出、機械学習の手順を踏むが、特徴抽出で適用する手法の詳細は、第 4 節で説明する。また、機械学習に用いる手法については

^{†1} 電気通信大学大学院 電気通信学研究科 情報工学専攻
Department of Computer Science, The University of Electro-Communications

^{†2} 電気通信大学 電気通信学部 情報工学科
Department of Computer Science, The University of Electro-Communications

第6節で述べる。実験方法と評価方法、結果を第7節で考察し、第8節で結論付ける。

2. 関連研究

我々の扱う位置テキストと航空写真を伴う画像認識について、それぞれ先に研究がおこなわれている。

位置テキスト情報を画像認識に取り入れた例として、Dhirajら⁶⁾の研究が挙げられる。彼らの研究では、位置情報からGISデータベースの逆ジオコーディングで地名などを求め、オッズ比に基づく確率的手法で分類し、これと画像特徴量との統合にはNaive fusionとConfidence-based fusionの2通りの手法を用いている。

Luoら⁷⁾はGoogle Earthから手作業で取得した853枚の比較的詳細な航空写真を用いた(以下この研究をThird Eyeと呼ぶ)。画像はFlickrから収集した981枚の写真と、被験者に旅行させることで収集した720枚の写真を使用している。画像特徴量はSIFTに基づくbag of visual words表現と、HSV空間によるカラーヒストグラムを単純ベクトル結合したものであり、画像の分類にはSVMを、航空写真の分類にAdaBoostを用いている。画像と航空写真の認識精度を融合するにあたっては、Qiら⁸⁾の手法に基づくSVMのメタ分類器を用いている。実験の結果、多くの分類カテゴリで、画像よりも航空写真の精度が良いことを示している。

本研究では一般物体認識を視野に入れ、多量の画像と航空写真を扱う。また、テキスト情報はYahoo!ローカルサーチを用いて抽出する。特徴量の融合にあたり、Multiple Kernel Learning(MKL)を導入する。

3. 手順と方針

我々の認識実験は全体的に、図1に示す要領で行われる。全体の流れとしては、Flickrより収集した位置情報付き画像から特徴抽出したものを、機械学習することで認識精度を検証するものである。

機械学習の段階において、各特徴量がどのように有効に利用されているかどうかを考察す

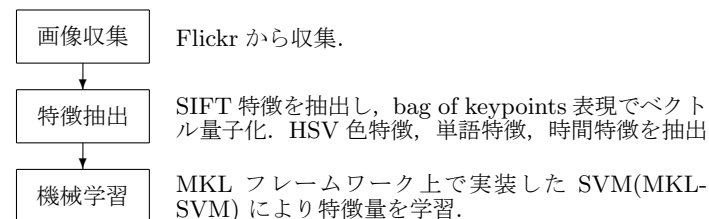


図1 研究の全体的手順

必要がある。本研究では、第6.2節で説明するMKL-SVMを用いて、認識精度のほかに本稿で追加した周辺テキスト、緯度経度、撮影時間についての特徴の有効性をMKLによって推定された重みに基づいて評価する。各手順において用いる手法の詳細については後述する。

4. 画像からの特徴抽出

本節では、画像から認識に必要な特徴を抽出する方法について述べる。画像の特徴を記述する手法としては、本実験では、特徴抽出のために局所特徴の一種であるSIFT特徴を用いる。また、この局所特徴を簡潔に記述するためにBag of Keypoints手法を用いてデータをベクトル量子化する。画像からはBoKとHSV色特徴、航空写真からはBoKのみを抽出する。

4.1 SIFT特徴

SIFT(Scale Invariant Feature Transform)²⁾とは、David Loweによって提案された特徴点とそれに付随する特徴ベクトルの抽出法である。特徴点周りの局所画像パターンを勾配方向ヒストグラムから成る128次元特徴ベクトルで表現する。特徴点の抽出は、自動で行う方法(ガウシアン差分による選定、特定の物体同士の対応点検出に有用)と、手動で指定する方法(主に格子状あるいはランダム点指定)がある。本実験では、特定の物体を認識することに固執しないので、10ピクセルの格子状に点を抽出する方法を採用する。

4.2 Bag of keypoints表現

Bag of keypointsモデル¹⁾とは、画像を局所特徴の集合と捉えた手法である。局所特徴の特徴ベクトルをベクトル量子化し、visual wordsと呼ばれる特徴ベクトルを生成する。これらの集合をコードブックと呼び、それを記述子として画像の特徴ベクトルを生成する。これにより画像をvisual wordsの集合(bag)として表現する。ベクトル量子化は、ユークリッド距離を尺度とし、k-Means法を用いて行う。本研究では、処理速度や結果における精度の差異を考慮した上で、クラスタ数を1000に固定した。

4.3 HSV色特徴

本研究では、配色の偏りに配慮するため、位置情報付き画像を5×5に分割し、各面について4³=64次元を、画像1枚について計1600次元のHSVカラーヒストグラムを抽出する。

5. メタデータからの特徴抽出

5.1 地理テキストの特徴

本研究では、位置情報の近傍を周辺検索することで得られるテキスト情報を利用する。具体的には、Yahoo!ローカルサーチAPIに撮影位置の緯度経度を与え、その周囲500m以内にあるランドマーク等に関する情報(建物・施設・公園など地図上に掲載され得る情報)を取得する。テキスト情報を、機械学習への入力に対応する特徴ベクトルに変換するにあたり、茶筌(chasen)¹¹⁾で単語に分割し、抽出した全名詞の出現頻度の上位2000語に対してヒストグラムを作成し、地理テキスト特徴とする。

5.2 時間・季節情報の利用

本研究では、メタデータとして撮影時刻に基づく時間情報も用いる。特徴量を得るにあたり、時間帯の類似度を判別できる形式とするため、時刻情報をヒストグラムに準ずる表現に変換する。具体的には、月について12ビンと時間について24ビンを用意し、該当する月・時刻とその隣接にそれぞれ0.5, 0.25を投票し、ヒストグラムに準ずる表現とする。

5.3 緯度経度

位置情報から抽出しうる特徴のみならず、メタデータとして緯度と経度の2次元ベクトルを直接学習・分類に用いる。測位系はFlickrで使用されているWGS84を想定する。ヒストグラムに準ずる表現ではないため、MKL-SVMに入力する際はRBFカーネルを用いる。

5.4 航空写真

航空写真を利用するにあたり、それぞれの画像に対応する位置情報を地図サイトで表示し、スクリーンキャプチャしたものを特徴抽出に利用する。特徴抽出の方法はデータセットの画像と同様である。写真の内容によっては、より詳細な縮尺を持つ航空写真の方が検証に際し有効であると考えられるが、撮影位置によっては、詳細な航空写真を入手できないこともある。広範なデータセットに対応するため、図2に示すような4種類のスケールを採用する。図2に示すように、航空写真は撮影位置が中央にくるような1辺256ピクセルの正方形に加工したものを利用する。

6. 学習と分類

本節では、機械学習で用いる手法の詳細を説明する。認識精度の検証に当たっては、我々の従来の手法と同様にサポートベクタマシン(SVM)を用いる。画像と航空写真の有効性について判定するために、マルチカーネル学習を導入する。

6.1 サポートベクタマシン

サポートベクタマシン(SVM)はニューロンのモデルとして最も単純な線形しきい素子を用いて、2クラスのパターン識別器を構成する手法である。カーネル学習法と組み合わせると非線形の識別器になる。本実験では、カーネル関数として非線形の χ^2 カーネルを用いるが、緯度経度に関してのみ、RBFカーネルを用いる。

6.2 マルチカーネル学習

本研究では、特徴を統合して画像を認識するために、複数の特徴量のカーネルを線形結合することにより統合カーネルを作成し、それをサポートベクタマシン(SVM)に適用して特徴統合による画像認識を実現する。最適なカーネル(カーネルを重みつきで線形結合したカーネル)のサブカーネルに対する重み j を学習する。これはマルチカーネル学習(Multiple Kernel Learning, MKL)³⁾問題と呼ばれ、統合カーネルは以下で定式化される。

$$K_{\text{combined}}(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^K \beta_j k_j(\mathbf{x}, \mathbf{x}') \quad \beta_j \geq 0, \sum_{j=1}^K \beta_j = 1 \quad (1)$$

最近の研究では、このMKL問題を凸最適化問題として効果的に解く方法が提案されてい

表1 カテゴリの種類と選定基準

カテゴリ	選定基準
ディズニールゾート 東京タワー	ディズニールゾートの敷地内で撮影されている。建造物や人物など 東京タワーが目立つように映っている。東京タワーの内部からの眺めは含まない。
橋	橋の構成要素が目立つように映っている。橋の上からの眺めは含まない。
神社 建物	神社の境内や鳥居などの建造物が明確に映っている。 高層建築または建物の全体が目立つように映っている。撮影位置から建物までの距離は遠すぎない。
城	城の建築が明確に映っている。
鉄道	鉄道車両が至近でまたは目立つように映っている。
公園	公園内の風景。遊具など公園特有の設備が映っている。
庭園	庭園内で撮影されている。水流や木々など庭園を象徴する人工構成要素が明確に映っている。
風景	周囲に遮るものがなく、遠方の物体が多く明確に映っている。
湖畔	湖畔で撮影されている。水面が目立つように映っている。
川	川岸で撮影されている。水面が目立つように映っている。
海岸	海岸で撮影されている。水面または砂浜が目立つように映っている。
像・モニュメント	仏像・銅像・モニュメントなどの不動の物体が明確に映っている。
自動車	自動車が明確に映っている。自動車からの眺めは含まない。
自転車	自転車が明確に映っている。自転車が小さいものや、自転車からの眺めは含まない。
落書き	落書きが目立つように映っている。
自動販売機	1台または複数台の自動販売機が明確に映っている。
紅葉	紅葉が目立つように映っている。
桜・花見	花の咲いている桜の木が目立つように映っている。または、花見の情景が映っている。
夕日	夕日が目立つように映っている。
コスプレ	派手な衣装をまとい1人または複数人で正面を向き顔・上半身・全体のいずれかが明確に映っている。
祭	神輿、提灯、縁日など祭りを演出する構成要素が十分に存在する。
猫	猫が目立つように映っている。
鳥	鳥が目立つように映っている。
花	花が目立つように映っているか、写真のほとんどが花で占められている。
ラーメン	食べられる状態にあるラーメンが明確に映っている。
寿司	食べられる状態にある1個または複数個の寿司が明確に映っている。

る⁴⁾。マルチカーネル学習はSVMのみを前提としたものではないが、SVMのフレームワークで解く方法が一般的で、MKL-SVMと呼ばれることもある。本研究では、SHOGUN⁵⁾ツールキットを用いて実装したMKL-SVMを使用して実験を行う。

7. 実験方法

実験は、各画像から抽出した特徴量をサポートベクタマシン(SVM)で学習させ、分類結果により精度を判定することにより行う。学習と分類は正解画像と不正解画像の2クラスで行う。正解画像については、画像の内容に応じていくつかのカテゴリのデータセット用意し、不正解画像は正解画像に用いないデータをランダムに選択することで構成される。また、SVMのフレームワークの下で、マルチカーネル学習で種類ごとの重みを推定する(MKL-SVM)ことにより、画像と航空写真のどちらがどれだけ有用であるかを判断する。

7.1 データセット

後述する複数のカテゴリを準備し、カテゴリ毎の精度検証を行う。Flickrでは大量の画像

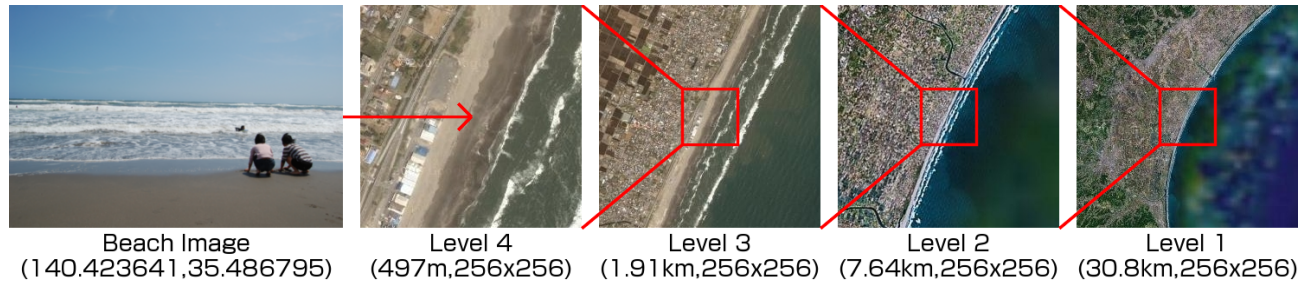


図 2 位置情報付き写真と航空写真の対応付け

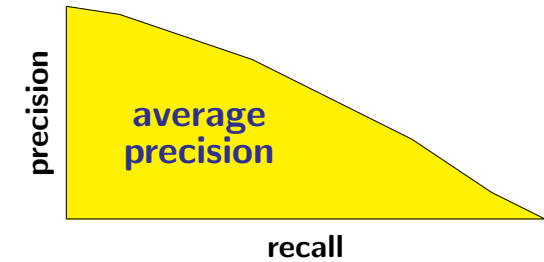


図 3 平均適合率.

を収集可能であり、位置情報付き写真は少なくとも 40 万枚は存在する。ただし、それらの写真は Flickr ユーザーの主観的な観点から撮影されたものであり、あるカテゴリにおいてそのカテゴリの写真そのものを描写しているとは限らない。例えばタワーに関する写真はタワーそのもののほかにタワーからの眺め、電車に関する写真は、電車自体のほかに車窓や駅舎などが含まれることがある。

データソースから得られる実験データセットを、より客観的なものに近付けるべく、我々は Flickr 画像を用いて認識の精度検証を行うにあたり、適切な画像を手動で選定することによりデータセットを作成する。

正解データセットについては、表 1 に示すような 28 種類のカテゴリを準備した。これらのデータは、タグなどのメタデータである程度分類したものの中から、カテゴリごとに写真の外観、位置情報ともに適切なものを手動で日本国内で撮影された 200 枚を選定した。正解画像の一例として、図 4 を挙げることができる。

カテゴリの選定に当たっては、Flickr で多く投稿される写真の内容の偏りを考慮しつつ、主に次のような 8 つの観点 (ジャンル) から選択した。また、それによる認識精度に関する仮説を示す。

位置に特有なランドマーク ディズニーリゾート、東京タワー [2 カテゴリ]
必然的に位置情報が偏るので、航空写真の特徴量に類似の外観ものが集中するので、航空写真の特徴量に依存する形で認識精度に改善がみられると考えられる。

狭い範囲の地理構成物 橋、神社、建物、城、鉄道 [5 カテゴリ]
地理的要素のうち、特定の物体を対象に描写している。航空写真がある程度詳細なものであれば、これらは把握可能であり、それらの特徴量によって認識精度に若干の改善がみられると考えられる。

広い範囲の地理構成物 公園、庭園、風景 [3 カテゴリ]
地理的要素のうち、撮影対象が広範囲で、画像中の物体も多岐にわたる。公園や庭園は、余ほど詳細なものでない限りは航空写真との相関は考えにくい。風景については、航空写真と関連することも考えられるが、撮影位置と視界の位置が異なるので、航空写

真による精度はさほど変化しないと考えられる。

地形 湖畔、川、海岸 [3 カテゴリ]
位置情報が特定の箇所に偏在する。色特徴、SIFT 特徴共に航空写真との相関が考えられる。

屋外の人工物 像・モニュメント、自動車、自転車、落書き、自動販売機 [5 カテゴリ]
主に屋外の物体であるが、航空写真から自明なものを確認することはほぼ不可能である。位置情報・時間情報ともに位置情報による精度変化は考えにくい。

時期依存的要素 紅葉、桜・花見、夕日、コスプレ、祭 [5 カテゴリ]
周期的な季節・時刻を伴うものと、人為的なイベントに関係する。夕日、桜・花見、紅葉については時系列的な相関と色特徴が精度の変化に影響を与えると考えられる。位置情報による変化は考えにくいであろう。コスプレと祭については、そのイベントの性質上時間と共に位置情報に大いに依存する可能性が高い。

天然の物体 猫、鳥、花 [3 カテゴリ]
動植物。航空写真・時間特徴共に何ら関連性は持たない。画像特徴で十分な精度が期待できるが、位置情報・時間情報ともに位置情報による精度変化は考えにくい。

食べ物 ラーメン、寿司 [2 カテゴリ]
位置情報に関しては若干関連があるかもしれない。画像特徴で十分な精度が期待できるが、位置情報・時間情報ともに位置情報による精度変化は考えにくい。

不正解データセットは、Flickr から収集した画像データの中から、正解データに用いられていないものをランダムに選定することで 200 枚準備した。

いずれのデータセットも、一般性を保証するため写真の撮影者が偏向しないように選定した。画像は縦と横のうち長い方の画素が 500 ピクセルになるようにリサイズしてある。

カテゴリごとのデータセットの枚数を決定するにあたり、各カテゴリにおいて Flickr から得られる適切な画像の枚数の平均を考慮した。正確な実験結果を保証するためには、一般に大量のデータセットを構築する必要があるが、本研究においては身近なデータからデータセットを構築することを重視した上で、今後万データソースを見直す場合は認識精度の比

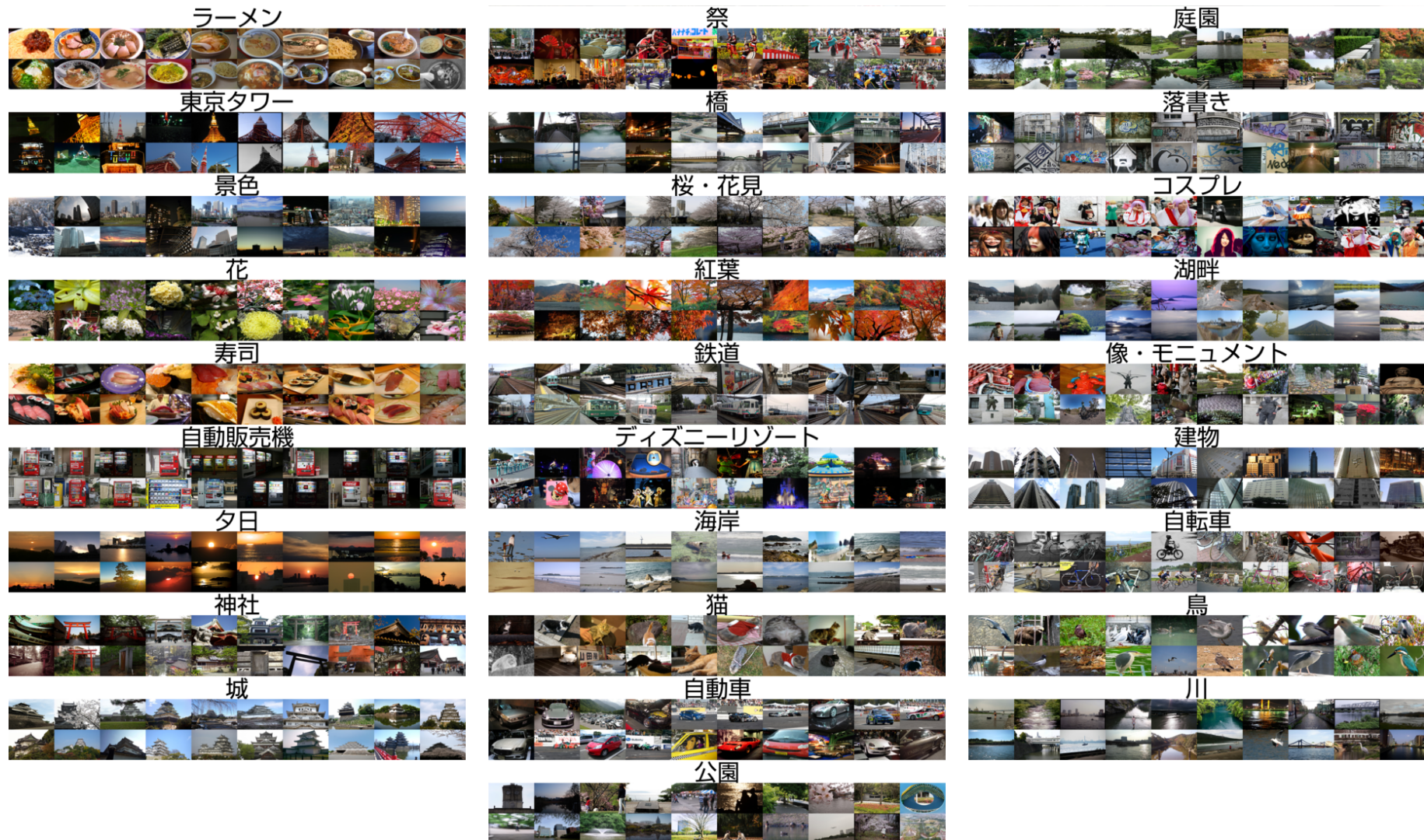


図 4 実験に用いる 28 種類の正解データセットの例.

表 2 MKL に与える特徴量の組み合わせ

組合せ	BoK	HSV	レベル 1	レベル 2	レベル 3	レベル 4	緯度経度	位置テキスト	時刻
2 種類 (H)	○	○	-	-	-	-	-	-	-
2 種類 (L1)	○	-	○	-	-	-	-	-	-
2 種類 (L2)	○	-	-	○	-	-	-	-	-
2 種類 (L3)	○	-	-	-	○	-	-	-	-
2 種類 (L4)	○	-	-	-	-	○	-	-	-
2 種類 (G)	○	-	-	-	-	-	○	-	-
2 種類 (T)	○	-	-	-	-	-	-	○	-
2 種類 (D)	○	-	-	-	-	-	-	-	○
5 種類 (A)	○	-	○	○	○	○	-	-	-
5 種類 (NA)	○	○	-	-	-	-	-	○	○
7 種類	○	-	○	○	○	○	○	○	-
9 種類	○	○	○	○	○	○	○	○	○

較検討と共に考慮することになる。

7.2 特徴量の種類と組み合わせ

画像の bag of keypoints 表現 (以下, 結果記述にあたり BoK とも記す), 航空写真の bag of keypoints 表現 (4 レベルでの 4 種類), 緯度経度, 周辺情報に関するヒストグラム, 撮影時刻に関するヒストグラム表現, HSV カラーヒストグラムの, 計 9 種類を扱う。これらを, MKL に与えるにあたり表 2 のように組み合わせる。

7.3 評価

認識実験に際しては, 後述のように写真の内容を予め各カテゴリに分類し, カテゴリごとに, 画像と航空写真をペアにしたものを MKL-SVM を用いて学習・分類・重みの推定を行う。

学習と分類はクロスバリデーションで行う。すなわち, データセットを 5 分割し, 1 つを分類用に, 残り 4 つを学習用に充てることを, 分類用に相当する各 fold について 5 通り繰り返す。

認識精度は各分類結果における SVM の出力値に対して, 以下のように平均適合率を計算することで評価する。SVM による出力値に基づいて, テストデータをソートする。ソートしたデータを最初から順番に読み込み, Positive のデータが出現した時点で, それまでの読み込んだデータの数 m_i と, Positive データの出現頻度 i を記録する。ここで $p_i = \frac{i}{m_i}$ とおく。最後までテストデータを読み込んだときのすべての Positive データの出現数を n とすると, 平均適合率 (Average Precision) は,

$$P = \frac{1}{n} \sum_{i=1}^n p_i \quad (2)$$

で計算される。これは precision-recall グラフの面積に相当する (図 3)。

7.4 実験結果

画像と航空写真を MKL-SVM で分類した結果について示す。ここでは, 画像と航空写真の 2 組について MKL-SVM で分類, 重み推定を行った。結果については, 5 回のクロスバリデーションの平均値を示す。

平均適合率の結果を表 3 に示す。表 3 において, 見出し冒頭の数字は, 融合した特徴の

表 3 MKL-SVM における平均適合率の計算結果。各カテゴリごと, 各ジャンルごとの結果をそれぞれ示す。見出しにおける数値は統合した特徴の数を表し, 2 種類の場合は画像 BoK と共に航空写真 BoK(L1~4), 緯度経度 (G), 位置テキスト (T), 時間 (D), 画像 HSV(H) と統合した。5 種類については, 航空写真 BoK 以外の特徴量 (NA) と, 航空写真 BoK の組み合わせ (A) でそれぞれ統合した。

ジャンル	カテゴリ	2(L1)	2(L2)	2(L3)	2(L4)	2(G)	2(T)	2(D)	2(H)	5(NA)	7	5(A)	9	BoK	HSV
狭い地理構成物	橋	73.81	72.62	73.46	72.97	69.73	74.19	69.63	69.78	74.65	75.00	74.32	75.12	69.78	63.87
	神社	71.75	70.96	73.04	75.90	70.02	75.96	71.07	71.34	76.08	75.73	75.01	76.20	70.32	70.60
	建物	79.92	79.53	80.20	80.29	79.15	80.63	78.28	79.92	81.35	80.73	79.80	81.04	78.58	74.35
	坂	82.56	82.58	82.81	83.12	81.83	83.31	80.32	82.09	83.74	83.65	82.99	83.73	80.13	81.21
	鉄道	75.29	75.74	74.90	76.10	74.39	77.48	74.76	77.71	78.83	77.90	77.03	78.98	74.43	72.59
	平均	76.67	76.29	76.88	77.68	75.02	78.31	74.81	76.17	78.93	78.60	77.83	79.02	72.52	74.65
屋外の人工物	像	69.21	68.47	70.34	70.17	66.46	71.97	65.65	67.78	72.34	72.30	70.78	72.49	66.45	65.94
	自動車	74.59	73.13	73.24	73.69	71.56	76.55	72.18	75.72	77.01	76.46	74.84	76.70	70.80	70.74
	自転車	77.82	77.37	76.75	77.63	76.20	78.74	76.51	76.72	78.87	78.61	77.90	78.80	70.71	81.85
	落書き	78.33	77.85	79.45	79.31	74.83	80.75	77.05	76.68	81.81	80.70	79.53	81.42	72.59	75.60
	自動販売機	81.34	81.04	80.86	80.67	79.91	81.96	80.88	83.26	83.39	81.93	81.59	83.28	73.88	82.81
	平均	76.26	75.57	76.13	76.29	73.79	77.99	74.45	76.03	78.69	78.00	76.93	78.54	73.39	73.15
時期依存的要素	紅葉	82.24	82.08	82.13	81.12	81.77	82.39	83.05	83.27	83.96	82.44	82.08	83.97	80.77	82.53
	桜・花見	80.20	80.28	80.62	80.41	80.19	80.83	82.28	82.13	82.80	80.78	80.56	82.55	80.22	79.70
	夕日	82.90	82.92	82.89	82.88	82.90	82.89	83.30	83.68	83.86	82.82	83.84	83.40	82.90	83.41
	コスプレ	82.81	82.99	83.21	83.66	78.37	83.63	79.12	80.29	83.73	83.69	83.68	83.73	75.67	79.42
	祭	73.52	73.31	74.09	75.95	73.40	76.97	75.37	76.90	79.23	77.36	76.20	79.58	73.31	74.89
	平均	80.33	80.31	80.59	80.81	79.32	81.34	80.62	81.25	82.72	81.42	81.07	82.73	79.99	78.57
広い地理構成物	公園	71.86	71.58	72.77	74.01	70.09	78.64	70.49	71.24	77.76	78.05	74.99	77.33	70.42	69.05
	庭園	80.00	80.55	80.30	80.07	79.16	81.71	78.31	79.80	82.20	81.91	81.06	82.27	77.57	77.61
	風景	74.75	74.77	75.61	74.57	74.29	76.95	74.83	75.75	78.00	77.04	75.16	78.06	74.23	74.28
	平均	75.54	75.63	76.22	76.22	74.51	79.10	74.54	75.60	79.32	79.00	77.07	79.22	73.65	74.08
地形	湖畔	80.88	81.12	81.06	80.48	79.01	82.13	79.34	79.41	82.67	82.06	81.42	82.39	78.27	76.25
	川	76.91	78.36	79.68	79.41	75.05	79.65	75.99	75.81	80.06	80.44	80.11	80.78	74.56	74.22
	海岸	83.22	83.29	83.22	83.46	81.02	83.72	81.80	81.70	83.72	83.53	83.50	83.55	81.02	80.06
	平均	80.34	80.92	81.32	81.12	78.36	81.83	79.04	78.97	82.15	82.01	81.68	82.24	76.84	77.95
天然の物体	猫	70.04	70.16	69.77	69.58	69.01	71.61	68.11	71.72	74.23	70.92	69.62	74.14	67.96	67.24
	鳥	77.44	77.79	78.62	79.03	69.71	79.75	74.30	73.00	80.36	79.79	78.89	80.26	69.75	71.32
	花	77.84	77.79	78.02	78.54	78.22	78.67	79.65	80.71	82.23	78.93	78.30	82.22	77.98	77.70
	平均	75.11	75.25	75.47	75.72	72.31	76.68	74.02	75.14	78.94	76.55	75.60	78.87	72.09	71.89
ランドマーク	ディズニー	84.09	84.10	83.89	84.01	75.32	84.10	68.97	73.20	84.10	84.09	84.14	84.09	68.37	70.98
	東京タワー	83.46	83.56	83.68	83.52	79.93	83.78	79.25	81.54	83.78	83.78	83.67	83.78	79.34	80.59
	平均	83.78	83.83	83.79	83.76	77.62	83.94	74.11	77.37	83.94	83.93	83.90	83.93	75.79	73.86
	食べ物	ラーメン	82.58	82.62	82.73	82.77	82.58	82.59	82.63	83.03	83.00	82.76	82.77	82.97	82.59
寿司		82.42	82.05	81.88	82.02	79.93	81.99	80.96	81.76	83.20	81.97	82.09	83.19	73.55	73.51
平均		82.50	82.33	82.31	82.40	81.25	82.29	81.79	82.40	83.10	82.36	82.43	83.08	80.18	81.22
全体平均		78.28	78.16	78.54	78.76	76.21	79.77	76.57	77.71	80.61	79.83	79.10	80.59	75.33	75.49

数を表わす。ベースラインである画像 BoK 単体の精度を BoK に, 同様に HSV 単体の精度を右端に示す。

特徴統合による精度の変化について, まず精度のみに着目して考察する。平均適合率におけるベースライン (BoK) との差分 (gain) を表 4 に示す。平均のみに着目すると, 単純に緯度経度を使う場合 gain は 0.72 程度でほとんどないが, 位置テキストや航空写真に変換することで 2.88 に上昇することを示唆しており, これにより提案手法の有効性が示される。

ジャンルごとの平均に着目すると, 位置に固有なランドマークが+10.05 で位置情報が有効で, 緯度経度でも+3.77 の gain である。緯度経度は固有ランドマーク以外では 1%未満で殆どが 0.5%以下であるから, 固有ランドマーク以外では基本的には有効ではない。一方, 天然の物体, 広い範囲の地理構成物, 地形は, +3.2~3.7 で位置情報に由来する特徴が有効であるが, 緯度経度は 0.4 程度でほとんど有効ではない。このことから緯度経度を位置テキストや航空写真に変換する本手法の有効性が確認できる。

位置情報が単なる付加的情報ではなく, 適切な重みにより分類されていることを確認するために, 表 6 にいくつかの組み合わせで MKL を行った時の重みの比較を示す。また, いかなる位置テキストが実験に影響を及ぼしたかを示すため, 位置テキストの重みが高かった

表 4 平均適合率におけるベースライン (BoK あるいは BoK+HSV) との差分 (gain).

記号	算出方法	概要
Text	2 種類 (T)-(BoK)	位置テキストによる性能向上
Date	2 種類 (D)-(BoK)	時間による性能向上
Geo	2 種類 (G)-(BoK)	緯度経度による性能向上
Aerial	5 種類 (A)-(BoK)	【参考】前回提案の航空写真による性能向上
NA	5 種類 (NA)-2 種類 (H)	追加した 3 種類 (緯度経度, 時間, テキスト) による性能向上
All	9 種類 -2 種類 (H)	すべての位置情報に由来する特徴による性能向上

ジャンル	Text	Date	Geo	Aerial	NA	All
位置に特有なランドマーク	10.08	0.25	3.77	10.05	6.57	6.57
屋外の人工物	4.84	1.30	0.64	3.77	2.65	2.51
狭い範囲の地理構成物	3.66	0.17	0.37	3.18	2.76	2.85
広い範囲の地理構成物	5.02	0.47	0.44	2.99	3.72	3.62
時期依存的要素	2.77	2.05	0.75	2.50	1.46	1.48
食べ物	1.07	0.57	0.03	1.21	0.71	0.68
地形	3.88	1.09	0.41	3.73	3.18	3.27
天然の物体	4.78	2.12	0.42	3.71	3.80	3.73
全体平均	4.28	1.08	0.72	3.61	2.89	2.88

6 カテゴリと、そうでない 6 カテゴリについて、出現頻度の高い上位 20 語の単語を表 5 に示す。

位置テキストにおいて頻繁に出現する単語は、概ね地理的には均一に分散していると思われる構成物の名称が多い。ただし、位置テキストの重みが大きいカテゴリに関しては、特定の地名を表すいくつかの固有名詞 (有明, 舞浜, 姫路, 六本木など) において高い頻度を記録している。このテキストの偏在が学習において大きく影響していると考えられる。

航空写真については、レベル 4 の重みが最も高いが、位置テキストを導入した場合、この重みが最も高くなる。このことから位置情報付き画像の認識にはテキストが精度向上に影響していると思える。

一方、緯度経度の重みは著しく少ないことから、緯度経度情報のみではその位置のコンテキストを記述するには十分でないと言える。したがって、位置情報を、対応するコンテキストに変換することの有効性を見て取ることができる。ただし、表 4 において航空写真による gain が +3.61 であり、位置テキストの gain に次いで多いことから、航空写真特徴のみの場合でも十分な精度向上を見込める可能性はあると考えられる。

時間情報の特徴は、1 日と 1 年の特定の周期に集中する場合に、認識に対して有効に作用する。本実験においては、カテゴリ「桜・花見」に対して高い重みを示したが、「夕日」や「紅葉」など、期間が長い場合やイベントの発生する時期に (同じ周期内であっても) ばらつきが目立つ場合、時間情報は有効性ではない。したがって、時間情報を画像認識に十分に適用するには、時間情報を補強するに足るメタデータの利用が必要になると思われる。

以上のことから、周辺テキストの特徴量は画像認識に有効であり、緯度経度・時間情報は有効ではなく、これらについては追加のメタデータが必要であると考えられるが、緯度経度については追加の位置情報特徴により認識への有効性が示された。

なお、Flickr から得られる写真は個人により撮影されているものである以上、位置情報の偏向については常に考慮しなければならない問題であるが、これについてはデータの大量収集と効率的なデータ・カテゴリの選定に対して包括的に配慮することで、引き続き対処して

いきたいと考える。

8. おわりに

8.1 まとめ

本研究では、位置情報付き写真の一般画像認識を拡張するにあたり、写真の撮影位置に対応する航空写真と周辺テキストの情報を付加的な画像特徴量として利用する手法を導入した。認識実験においては各種特徴量の組み合わせによる認識精度の変化を検証するとともに、マルチカーネル学習 (MKL, Multiple Kernel Learning) を導入することで、特徴量の種類ごとの認識への関与を定量的に分析した。28 種類のカテゴリによる実験では、画像単体の場合と比較して 75% から 80% へ向上した。

緯度経度のみでは、特定のランドマークを除いて十分な精度向上は得られなかった。一方、位置情報に由来する特徴量は、緯度経度の拡張・補強し、認識精度の向上に有効に作用した。位置テキストは、その有効性について航空写真特徴を凌駕した。時間情報は、特定の周期を捉えられない場合は有効に作用せず、画像認識への適切な導入にあたりこれを補完するメタデータが必要とされる。

8.2 今後の課題

本実験は実際のところ、データセットを作成する際のデータの入手可能性に関する不可抗力的な制約に対し、ある程度の妥協を容認した。位置情報の認識精度をさらに検証し、役立てるには、大量かつ良質なデータ収集の手段について改めて模索していく必要がある。また、より一般的な画像認識の枠組みを実現するにあたり、海外の位置情報付き写真も利用することを視野に入れることが今後不可欠である。

参考文献

- 1) G. Csurka, C. Bray, C. Dance, and L. Fan. Visual categorization with bags of keypoints. *Proc. of ECCV Workshop on Statistical Learning in Computer Vision*, pp. 59-74, 2004.
- 2) D. Lowe. Distinctive Image Features from Scale-Invariant Keypoints. *International Journal of Computer Vision*, Vol. 60, No. 2, pp. 91-110, 2004.
- 3) 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.
- 4) 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.
- 5) Shogun - A Large Scale Machine Learning Toolbox. <http://www.shogun-toolbox.org/>
- 6) D. Joshi and J. Luo. Inferring generic activities and events from image content and bags of geo-tags. In *Proc. of ACM International Conference on Image and Video Retrieval*, 2008.
- 7) J. Luo, J. Yu, D. Joshi, and W. Hao. Event recognition: Viewing the world with a third eye. In *Proc. of ACM International Conference Multimedia*, 2008.
- 8) G. Qi, X. Hua, Y. Rui, J. Tang, T. Mei, and H. Zhang. Correlative multi-label video annotation. In *Proc. of ACM International Conference Multimedia*, pp. 17-26, 2007.
- 9) K. Yaegashi and K. Yanai. Can Geotags Help Image Recognition?. *Proc. of the Pacific-Rim Symposium on Image and Video Technology*, pp. 361-373, 2009.
- 10) 八重樫恵太, 柳井啓司. マルチカーネル学習を用いた画像特徴と航空写真特徴の重要度の推定, 電

表 5 周辺テキストの抽出結果. ヒストグラムビンをカテゴリごとに平均し, 対応する単語を降順で上位 20 語まで示す. 表の上段に 9 種類の分類において位置テキストの重みが特に高かった 6 カテゴリ, 下段には特に低かった 6 カテゴリを示す.

コスプレ		鳥		ディズニーリゾート		城		神社		東京タワー	
有明	0.0391	市立	0.0202	場	0.1170	姫路	0.0324	ビル	0.0340	ビル	0.1089
ビル	0.0388	寺	0.0197	駐車	0.1074	市立	0.0252	局	0.0187	大使館	0.0395
東京	0.0250	局	0.0197	東京	0.0649	局	0.0201	郵便	0.0182	会館	0.0255
場	0.0223	郵便	0.0191	舞浜	0.0424	ビル	0.0199	院	0.0174	六本木	0.0211
原宿	0.0222	小学校	0.0189	リゾート	0.0416	学校	0.0173	寺	0.0169	神谷町	0.0199
センター	0.0181	ビル	0.0171	パーキング	0.0347	橋	0.0161	小学校	0.0136	日本	0.0186
公園	0.0144	橋	0.0150	ホテル	0.0321	郵便	0.0158	神社	0.0132	麻布	0.0183
棟	0.0130	線	0.0136	南葛西	0.0240	寺	0.0149	市立	0.0131	聖	0.0183
幕張	0.0128	センター	0.0131	ベイ	0.0197	公園	0.0143	線	0.0122	虎ノ門	0.0181
号	0.0124	院	0.0121	入口	0.0139	小学校	0.0133	幼稚園	0.0113	飯倉	0.0180
明治	0.0119	公園	0.0117	東急	0.0120	センター	0.0131	センター	0.0107	寺	0.0170
神宮	0.0099	幼稚園	0.0114	区立	0.0115	高等	0.0126	館	0.0104	タワー	0.0164
橋	0.0088	保育園	0.0111	ンター	0.0108	神社	0.0124	会館	0.0099	館	0.0140
ホテル	0.0083	場	0.0094	公園	0.0105	幼稚園	0.0123	保育園	0.0094	泉	0.0111
ホール	0.0081	学校	0.0088	ランド	0.0100	会館	0.0114	ホテル	0.0092	局	0.0109
東郷	0.0079	神社	0.0088	野毛	0.0100	私立	0.0109	公園	0.0086	郵便	0.0109
会館	0.0078	東京	0.0087	立体	0.0092	城	0.0107	学校	0.0085	芝	0.0101
通り	0.0077	中学校	0.0087	クリスタル	0.0089	大阪	0.0101	東京	0.0082	殿	0.0097
記念	0.0076	旭川	0.0078	保育園	0.0089	中学校	0.0100	場	0.0079	教会	0.0093
代々木	0.0071	館	0.0070	小学校	0.0087	小学校	0.0092	中学校	0.0078	復興	0.0092

花		桜・花見		自動販売機		紅葉		ラーメン		夕日	
ビル	0.0326	ビル	0.0359	ビル	0.0286	院	0.0274	ビル	0.0522	線	0.0270
局	0.0186	局	0.0177	局	0.02	寺	0.0241	局	0.0173	公園	0.0187
郵便	0.0178	郵便	0.0171	郵便	0.0187	市立	0.0202	郵便	0.0169	市立	0.0175
寺	0.0163	寺	0.0151	保育園	0.0168	局	0.0166	市立	0.0117	小学校	0.0166
小学校	0.0161	新宿	0.0136	富士山	0.0154	小学校	0.0163	センター	0.0115	局	0.0166
市立	0.0158	小学校	0.0133	小学校	0.0141	郵便	0.0161	ホテル	0.0109	ビル	0.0165
新宿	0.0136	東京	0.0124	区立	0.0132	線	0.0159	小学校	0.0109	郵便	0.0158
幼稚園	0.0121	市立	0.0115	寺	0.0123	ビル	0.0142	寺	0.0106	センター	0.0141
東京	0.0117	センター	0.0113	センター	0.0122	嵯峨	0.0125	東京	0.0102	橋	0.0139
センター	0.0116	院	0.0111	幼稚園	0.0122	保育園	0.0116	保育園	0.0099	寺	0.0127
線	0.0112	会館	0.0108	ホテル	0.0114	京都	0.0103	幼稚園	0.0099	保育園	0.0124
保育園	0.0111	学校	0.0108	会館	0.0108	堂	0.0101	公園	0.0095	横浜	0.0107
学校	0.0103	幼稚園	0.0107	東京	0.0104	橋	0.0101	会館	0.0093	東京	0.0101
公園	0.0098	保育園	0.0104	公園	0.0100	神社	0.0100	渋谷	0.0083	場	0.0094
院	0.0094	館	0.0084	神社	0.0099	幼稚園	0.0097	橋	0.0083	幼稚園	0.0092
セブン	0.0092	公園	0.0082	橋	0.0093	館	0.0095	セブン	0.0081	セブン	0.0087
イレブン	0.0091	中学校	0.0079	橋	0.0090	場	0.0089	イレブン	0.008	イレブン	0.0087
中学校	0.0088	区立	0.0078	セブン	0.0086	病院	0.0088	京都	0.008	ホテル	0.0084
橋	0.0085	橋	0.0076	イレブン	0.0084	公園	0.0084	ローソン	0.0079	中学校	0.0077
会館	0.0082	神社	0.0075	病院	0.0081	学校	0.0084	銀行	0.0078	国道	0.0074

表 6 MKL による 9 種類の特徴の重み推定結果.

カテゴリ	画像情報 (BoK,HSV)	航空写真 (1,2,3,4)	緯度経度	位置テキスト	時刻
狭い地理構成物	橋	0.3843(0.2690,0.1153)	0.3408(0.0529,0.0557,0.0295,0.2028)	0.0089	0.2289 0.0371
	神社	0.4504(0.0999,0.3505)	0.0690(0.0012,0.0020,0.0001,0.0657)	0.0043	0.4307 0.0456
	建物	0.6147(0.3897,0.2251)	0.1207(0.0500,0.0425,0.0004,0.0278)	0.0455	0.2123 0.0068
	城	0.2441(0.0654,0.1787)	0.1094(0.0877,0.0041,0.0124,0.0053)	0.0521	0.5941 0.0002
	鉄道	0.6918(0.3016,0.3903)	0.1004(0.0212,0.0151,0.0004,0.0637)	0.0001	0.1869 0.0208
	平均	0.4771(0.2251,0.2519)	0.1480(0.0426,0.0239,0.0085,0.0731)	0.0222	0.3306 0.0221
屋外の人工物	像	0.2223(0.2156,0.0067)	0.1381(0.0587,0.0041,0.0464,0.0288)	0.4492	0.0336 0.1568
	自動車	0.6772(0.3290,0.3481)	0.0285(0.0065,0.0000,0.0001,0.0219)	0.0146	0.2678 0.0119
	自転車	0.5520(0.4410,0.1111)	0.0642(0.0076,0.0143,0.0003,0.0420)	0.0098	0.3391 0.0349
	落書き	0.2304(0.0682,0.1621)	0.2846(0.0141,0.1064,0.1043,0.0598)	0.0493	0.3728 0.0630
	自動販売機	0.8755(0.2851,0.5904)	0.0533(0.0064,0.0127,0.0094,0.0248)	0.0032	0.0085 0.0594
	平均	0.5115(0.2678,0.2437)	0.1137(0.0187,0.0275,0.0321,0.0355)	0.1052	0.2044 0.0652
時期依存的要素	紅葉	0.6158(0.2298,0.3861)	0.1490(0.0836,0.0297,0.0037,0.0318)	0.0066	0.0025 0.2261
	桜・花見	0.4437(0.3291,0.1147)	0.0402(0.0002,0.0019,0.0381,0.0001)	0.0016	0.0093 0.5051
	夕日	0.8096(0.4022,0.4074)	0.0635(0.0001,0.0004,0.0000,0.0631)	0.0000	0.0002 0.1267
	コスプレ	0.0485(0.0352,0.0133)	0.0036(0.0001,0.0001,0.0001,0.0033)	0.0126	0.9291 0.0062
	祭	0.5667(0.2347,0.3320)	0.1248(0.0012,0.0004,0.0005,0.1226)	0.0015	0.2154 0.0916
	平均	0.4969(0.2462,0.2507)	0.0762(0.0170,0.0065,0.0085,0.0442)	0.0045	0.2313 0.1911
広い地理構成物	公園	0.3971(0.1688,0.2283)	0.0881(0.0034,0.0101,0.0205,0.0541)	0.0169	0.4178 0.0800
	庭園	0.4740(0.1338,0.3402)	0.1227(0.0080,0.0033,0.0170,0.0944)	0.0269	0.3512 0.0252
	風景	0.5634(0.2895,0.2739)	0.0452(0.0000,0.0001,0.0078,0.0374)	0.0043	0.3587 0.0283
	平均	0.4782(0.1974,0.2808)	0.0853(0.0038,0.0045,0.0151,0.0620)	0.0161	0.3759 0.0445
地形	湖畔	0.3539(0.2184,0.1355)	0.2694(0.0749,0.0877,0.0150,0.0918)	0.0040	0.3435 0.0292
	川	0.3465(0.2492,0.0973)	0.3181(0.0063,0.0523,0.0614,0.1981)	0.0108	0.2447 0.0800
	海岸	0.1888(0.1816,0.0073)	0.5003(0.0108,0.1632,0.0013,0.3250)	0.0015	0.2934 0.0159
	平均	0.2964(0.2164,0.0800)	0.3626(0.0307,0.1011,0.0259,0.2049)	0.0054	0.2938 0.0417
天然の物体	猫	0.6929(0.2690,0.4239)	0.0718(0.0002,0.0032,0.0007,0.0677)	0.0100	0.2217 0.0036
	鳥	0.1869(0.0888,0.0980)	0.0296(0.0004,0.0039,0.0002,0.0252)	0.0000	0.7126 0.0708
	花	0.8196(0.3356,0.4840)	0.0063(0.0001,0.0007,0.0015,0.0041)	0.0179	0.0269 0.1293
	平均	0.5665(0.2312,0.3353)	0.0359(0.0002,0.0026,0.0008,0.0323)	0.0093	0.3204 0.0679
ランドマーク	ディズニー	0.0001(0.0001,0.0001)	0.3288(0.1460,0.1827,0.0001,0.0001)	0.0004	0.6706 0.0000
	東京タワー	0.1479(0.1474,0.0005)	0.4230(0.0000,0.0011,0.4217,0.0001)	0.0007	0.4280 0.0003
	平均	0.0740(0.0737,0.0003)	0.3759(0.0730,0.0919,0.2109,0.0001)	0.0006	0.5493 0.0002
食べ物	ラーメン	0.9317(0.4588,0.4729)	0.0563(0.0000,0.0000,0.0092,0.0470)	0.0018	0.0005 0.0097
	寿司	0.6789(0.2435,0.4354)	0.0021(0.0001,0.0005,0.0000,0.0016)	0.0043	0.2616 0.0532
	平均	0.8053(0.3512,0.4541)	0.0292(0.0001,0.0002,0.0046,0.243)	0.0030	0.1310 0.0314
全体平均	0.4717(0.2314,0.2403)	0.1411(0.0229,0.0285,0.0286,0.0611)	0.0271	0.2915 0.0685	

子情報通信学会パターン認識・メディア理解研究会, 2009.

11) 茶釜. <http://chasen.naist.jp/hiki/ChaSen/>.