

研究種目：基盤研究（C）
 研究期間：2007～2009
 課題番号：19500076
 研究課題名（和文） 圧縮空間を用いたマルチメディアデータマイニングとそのウェブマイニングへの応用
 研究課題名（英文） Compression feature space based data mining and its application to web mining
 研究代表者
 渡邊 俊典（WATANABE TOSHINORI）
 電気通信大学・大学院情報システム学研究所・教授
 研究者番号：10242348

研究成果の概要（和文）：

インターネットや携帯電話の発展の中で文章、音声、画像などのマルチメディアデータが爆発的に増大している。本研究では、人手介入なしに、計算機によってこれらを分類あるいは検索する方式を検討した。我々が過去検討してきた、圧縮率によるテキストの特徴表現方式を原理としつつ、より高性能な圧縮性特徴空間の構成可能性の検討、文書や画像分類への適用などを試みた。文法知識を事前準備せずに、文章や画像に適用できること、従来方式をしのぐ性能も発揮できることなどを確認できた。なお、EU での衛星画像利用地球環境管理国際プロジェクト（GEOSS）関係機関より招待され、衛星画像処理への応用可能性について講演も実施した。

研究成果の概要（英文）：

Rapid explosion of internet and mobile telephone require us highly automatic computerized systems that can classify and/or recognize multimedia data including documents, sounds, and images. In this project, based on the principle of text featuring by compressibility which we have been exploiting, new investigations are made on the new problem of how to construct efficient compressibility feature space and its applications. Compared with conventional approaches, its merits include: free of linguistic knowledge (grammar), wide applicability, simplicity, and high performance.

交付決定額

（金額単位：円）

	直接経費	間接経費	合計
2007年度	800,000	240,000	1,040,000
2008年度	1,100,000	330,000	1,430,000
2009年度	1,400,000	420,000	1,820,000
年度			
年度			
総計	3,300,000	990,000	4,290,000

研究分野：総合領域

科研費の分科・細目：情報学、メディア情報学・データベース

キーワード：マルチメディア情報処理、データマイニング、データ圧縮

1. 研究開始当初の背景

近年、情報化が急速に発展するにつれ、計算機の扱うデータの大容量化、種類の多様化

が進んでいる。データベースに保存されたデータと比べ、インターネットや携帯電話で利用されるテキスト、画像、音声などは構文木

などの構造化手法になじみにくく、これらを、文法規則によってモデル化しようとする検討もなされているが、期待されたほどには成功していない。これらの構造化しにくいデータへの柔軟なアクセスを実現することは未解決の課題である。現状では主要な取り組みは、下記のようなものである。

マルチメディアデータの、①意味ラベルを人が付与し、以降の計算機処理を容易化するアプローチと、②なるべく人手によらず計算機で構造や意味を自動把握することを目指すデータマイニングアプローチとに大別される。前者は、多様な意味を持つマルチメディアデータに特定のラベルを与えることが将来の検索を限定してしまうという問題や、情報付与のための人手コストの問題を含んでいるが、これといった代替手段がないため広く採用されている。後者は、近年多くの研究がなされているにもかかわらず、顕著な応用領域を切り開くまでにはいたっていない。例えば、ベイジアンネットワークや Markov ランダム場 (MRF) でデータの内容をモデル化する手法が近年盛んであるが、モデルの基本構造は人による設計を要し、またモデルによるデータの解析には高い計算コストが必要となる。

2. 研究の目的

本研究では、圧縮率特徴空間（圧縮空間、あるいは PRDC (Pattern Representation Scheme Using Data Compression)）における多次元データ表現方式を利用した新たなマルチメディアデータマイニングの可能性を探求する。その可能性を左右すると予想される下記の諸機能の検討を目的とする。

(1) テキスト解析の検討: PRDC 研究の中で、課題となっているテキストの内容分析に関し、文法を用いず、類似性解析・検索、断片の特性分析 (不要語、新語などの抽出)、要約を可能とする方式。および、時間変化するテキスト集合への本手法の対応能力を向上させるための自動適応方式、分析結果の人への視覚化表示方式を検討する。

(2) 画像解析の検討: 画像に対する、類似性解析・検索、特徴的領域の同定、要約方式を検討する。

(3) マルチメディアデータ解析: 圧縮性特徴量とデータの共起性を利用した、異タイプメディアデータ間の自動関連付けの可能性を検討する。

(4) Web マイニングシステムのプロトタイプ: 上記の諸手法を統合的に用いて、Web 情報の自動分類、検索、内容推測、結果のビジュアル表示などを行なうプロトタイプシステムを Web サーバ上に実現し、現在の先端検索エンジンと比較し、問題点を発見して本研究を更に改良する。

3. 研究の方法

(1) テキスト解析の検討: テキスト全体の圧縮性による類似度解析とともに、テキスト断片の圧縮性に基づくトピック抽出を試みた。提案手法を検証するために、人工のモデル文書を生成し、原理レベルでの検証を行なった後、実テキストによる検証を行なった。流行語など、テキストの時間変化については新トピックの出現問題として扱った。また、解析の基礎となる圧縮空間の構成方式に関し、分類能力を保持したままの次元低下法などを検討した。

(2) 画像解析の検討: 画像をランレングス方向 (行方向) に走査してコーディングして得られるテキストを用い、上記 (1) を適用することを試みた。さらに、テクスチャ画像について、向きの変化やズームに対するロバスト性も併せて検討した。

(3) マルチメディアデータ解析: 上記 (2) の画像解析研究の結果を踏まえ、Web データによく見られるような、画像とテキストの隣接配置に対する本方式の適用法を考察することとした。

(4) Web マイニングシステムプロトタイプ: Web クローラで収集した実文書 (日本語、英語、中国語、およびこれらの混在) について、上記 (1) の応用可能性を検証し、既存手法と比較した。

(5) その他: 上記研究の中で、新たに発掘される問題については、重要度に応じた検討を実施することとした。

4. 研究成果

以下に、研究目的にリンクさせつつ、研究成果を説明する。

(1) テキスト解析の検討については、主として独立成分分析を利用した圧縮空間の独立化、文書間の関係分析方式、対話形式による多次元空間の次元削減方式について下記の成果を得た。

まず、独立圧縮空間の構成方式について検討した。特徴表現方法として有名な bag-of-words 法と N-gram モデルが、自然言語処理、テキストマイニングと文書解析によく利用されてきた。一方、PRDC はテキストデータのみではなく、マルチメディアデータ解析において、良い性能を示したが、圧縮辞書の選択や特徴空間の構築などの問題が残されている。本研究では、独立な圧縮率空間の構築法を検討した。

提案法では、独立成分分析法を利用して、圧縮率空間の次元数を削減し、より性能がよい特徴表現空間を構築に成功した。実験では、日本語のニュースサイトから収集したデータと、URCS、CLASSIC3、Reuters-21578 という

表 2. 次元縮小後の分類結果

人手による分類		次元縮小後の分類クラス							
トピック	記事数	クラス 1	クラス 2	クラス 3	クラス 4	クラス 5	クラス 6	クラス 7	クラス 8
履修	10	10							
松坂	10		10						
日興	10			10					
中越	10				7				
アジア杯	10					9			
参院選	20						15		
石川	10							9	
年金	20						2		17
分類再現率		100	100	100	70	90	75	90	85
分類精度		100	100	100	100	100	88	100	100

高次特徴空間の構成について、更に以下の検討も行った。すなわち、様々な文書を分類するための高次元特徴空間を定義したとき、その特徴空間の良さを計測したい。例えば物理学書は数学的な記述を多く含むので、そのクラス(集まり)は特徴空間上で歴史書や芸術書のクラスよりは数学書のクラスの近くに配されるはずであるが、実際にそうなっているか。他にもクラス間の距離が大きすぎたり小さすぎたりしないかなどを計測したい。こういった目的を達する最も直感的な方法は視覚的に特徴空間を観察することである。特徴空間は一般に非常に高次元である一方で、4次元以上の空間を平面的な紙の上やモニターに描くことは困難である。このような場合、従来では高次元特徴空間を主成分分析やFastMapなどの次元縮約法を用いて2次元、あるいは3次元空間に変換することが試みられている。

本研究では高次元特徴空間の形状を測り、無用な次元を削減するために、高次元空間内の多数の点の散布状況を調べるための外郭点の高速計算アルゴリズム(球分割法)を検討した。考案した球分割法は単純ガウス分布だけではなく様々な混合ガウス分布から外郭を取り出すことができることを確認した。また、単純ガウス分布に対する球分割法の実行時間は確率アルゴリズム(LSH)を用いることで10倍から20倍程度高速化できた。しかし、高速化の度合いは理論値よりも小さかった。この原因はデータの分布がアルゴリズムで仮定した球対称ではないためと考えられる。(学会発表⑤)

(2) 画像解析の検討については、ランレングスコーディングによって得たテキストを用いることで原画像の分類を試みた。また、テキスト画像にも適用し、ある程度のサイズや角度変化の下でも良好な分類性能を得ることを確認した。

最初に、低周波成分主体の画像、高周波成分主体の画像を意図的に混合した人工画像や自然画像を用いて、画像をテキスト化する際に必要となるアルファベット定義用の単

位ランレングス長Lの適正值を設定した。100種類1,000枚の画像データセットに対してテキスト化画像の分類を試みたところ、平均F値尺度が37%であり高精度とは言えない結果となった。原因として主に画像内の低周波背景が問題であると考えられた。画像が周波数成分の異なる領域を含む場合、単一の単位ランレングス長Lを用いると、情報損失が発生する。これに対しては、同一画像を複数のLでテキスト化し、得られる圧縮率ベクトルを結合することで、分類能力が若干高められることを確かめた(表3)。背景を除き、前景を取り出した場合についても同様の実験を行った。この場合、分類精度はより大きく向上した(表4)(学会発表①と⑦)。

表 3. L=2, 8(画素)の結合空間を用いた10種類データの分類結果

人手による分類		k-means法により得られたクラス									
データセット	種類	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10
パラ	10	9		1							
ヒマワリ	10	1	5	1			1	2			
象	10	3	3	3							1
日の出	10	8			2						
地球	10	1	1	3		5					
パンダ	10					1	6		1		2
トラ	10					1	3	6			
岩	10					1			4	2	3
ヨット	10			2					1	6	
家	10					1			2	1	6
分類再現率		0.9	0.5	0.3	0.2	0.5	0.6	0.6	0.4	0.6	0.6
分類適合率		0.41	0.56	0.38	0.5	0.56	0.60	0.75	0.50	0.66	0.5
F-値		0.56	0.53	0.34	0.28	0.52	0.60	0.66	0.44	0.62	0.54

表 4. 前景を切り出した画像の分類結果

人手による分類		k-means法により得られたクラス		
テストデータ		クラス1	クラス2	クラス3
ジャンル	種類	花	象	トラ
植物	花	26	4	
動物	象		30	
人工画像(イラスト)	トラ		17	13
分類適合率		1.00	0.59	1.00
分類再現率		0.86	1.00	0.44
F-値		0.93	0.74	0.62

次に、画像処理で重要なテキスト解析について検討した。従来法の場合、テキスト分類と学習の2ステップで構成される方式がよくみられる。本研究では、テキスト圧縮法を利用して、画像中に現れる本質的な頻出特徴パターンを発見することにより、学習ステップを不要とした。

実験では、Brodatz データベースと物体のテキスト表面を用いて、画像の回転、ズームなどの変化の下で認識能力の評価を行った。また、テキスト画像の検索についても実験を行った。提案法がかなり広範な変動に対してロバストであることを確認した。

(3) マルチメディアデータ解析: テキストや画像が同時出現することが多いウェブページの分類では、例えば、ページを決定木アルゴリズムを用いて画像、テキストの2つ

のカテゴリに分け、テキスト部位では単語間の繋がりを考慮する CNN-like word net などの特徴抽出法を使用し、画像部位に関しては領域輪郭線から領域の意味を推定する前段処理のあと、ベイズ推定法を用いてテキスト及び画像を結合させた分類を行っている。このような現在の主流アプローチでは、情報の種類毎に別々の手法や表現形式を用いたあと、面倒な重みパラメータの調整を要する統計的情報統合の段階が必要となる。この問題に対して、上記(1)のテキスト解析、(2)の画像解析研究の結果から本方式では双方を圧縮率特徴ベクトルで表現し、双方を結合したベクトルによって、画像とテキストの共起事象を表現する方式が有効であると考えられる。これによって情報統合の段階が単なるベクトル結合によって実現できる。関連データ収集の困難さから、実証については将来課題となった。

(4) ウェブプロトタイピング

インターネットの上クローラ(情報収集ソフトウェアロボット)を用いて、ウェブドキュメントを収集し、上記の(1)を用いたウェブページ分類実験を行った。

英語、日本語、中国語及びそれらの混在文書の分類が従来方式(Bag-of-words, N-gram)よりも良好な分類性能となることを確認した。

(5) その他

本研究の原理としている圧縮性に基づくデータ分類法(PRDC)に関し、Romania 政府機関主催の GEOSS ワークショップに招待され、画像データマイニングにおける課題解決のための圧縮性特徴量を利用したアプローチについて講演した。(その他①)

5. 主な発表論文等

(研究代表者、研究分担者及び連携研究者には下線)

[学会発表] (計9件)

- ① Nuo Zhang and Toshinori Watanabe, Image Representation and Classification based on Data Compression, Proceedings of the 25th Annual ACM, pp.981-982, Mar. 22-26, 2010, Sierre, Switzerland.
- ② Nuo Zhang and Toshinori Watanabe, Documents Representation Based on Independent Compressibility Feature Space, Proceedings of ICAART'10, INSTICC, pp.217-222, Jan. 22-24, 2010, Valencia, Spain.
- ③ Nuo Zhang and Toshinori Watanabe,

Documents Clustering Based on Optimized Compressibility Vector Space, Proceedings of International Conference on Computational Intelligence and Software Engineering (CiSE), IEEE, Dec. 11-13, 2009, Wuhan, China.

- ④ Toshinori Watanabe, On the Possibility of Highly Automated Image Information Mining: Problems and Possible Solutions, Workshop on Innovative Data Mining Techniques in Support of GEOSS, Aug. 31 - Sep. 2, 2009, Sinaia, Romania.
- ⑤ 小林 郁弥, 渡辺 俊典, 古賀 久志, ユークリッド空間内の点分布の外郭を求めるアルゴリズム, 信学技報 PRMU2009-171, pp.115-120, 2010/1/21, 京都.
- ⑥ Nuo Zhang, Daisuke Matsuzaki, Toshinori Watanabe and Hisashi Koga, Document Relation Analysis Based on Compressibility Vector, Proceedings of ICAART09, INSTICC, pp.255-260, Jan. 19-21, 2009, Porto, Portugal.
- ⑦ 平井 敦之, 張 諾, 渡辺 俊典, 古賀 久志, テキスト化を介した画像分類手法の提案, AI2008-42, pp.7-12, 2009/1/16, 東京.
- ⑧ Nuo Zhang, Toshinori Watanabe, Daisuke Matsuzaki and Hisashi Koga, A Novel Document Analysis Method Using Compressibility Vector, Proceedings of the 1st International Symposium on Data, Privacy, and E-Commerce (ISDPE'07), IEEE-CS, pp.38-40, Nov. 1-3, 2007, Chengdu, China.
- ⑨ 山崎 啓介, 張 諾, 渡辺 俊典, 古賀 久志, 高次元圧縮空間の対話的手法による次元縮小, 情報処理学会研究報告 2007-NL-181, pp.35-40, 2007/09/25, 東京.

[その他]

ホームページ等

- ① http://events.rosa-rc.ro/docs/presentations/sinaia%2031_08/ROSA_WS_T.Watanabe.pdf

6. 研究組織

(1) 研究代表者

渡邊 俊典 (WATANABE TOSHINORI)
電気通信大学・大学院情報システム学
研究科・教授
研究者番号:10242348

(2)研究分担者

古賀 久志 (KOGA HISASHI)
電気通信大学・大学院情報システム学
研究科・准教授
研究者番号：40361836

張 諾 (ZHANG NUO)
電気通信大学・大学院情報システム学
研究科・助教
研究者番号：20436736

横山 貴紀 (YOKOYAMA TAKANORI)
電気通信大学・大学院情報システム学
研究科・助教
研究者番号：10401621

(3)連携研究者