

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4254623号
(P4254623)

(45) 発行日 平成21年4月15日(2009.4.15)

(24) 登録日 平成21年2月6日(2009.2.6)

(51) Int.Cl. F 1
G 0 6 F 17/30 (2006.01)
 G 0 6 F 17/30 2 2 0 Z
 G 0 6 F 17/30 1 7 0 A

請求項の数 9 (全 16 頁)

(21) 出願番号 特願2004-170612 (P2004-170612)
 (22) 出願日 平成16年6月9日(2004.6.9)
 (65) 公開番号 特開2005-352613 (P2005-352613A)
 (43) 公開日 平成17年12月22日(2005.12.22)
 審査請求日 平成17年7月20日(2005.7.20)

(73) 特許権者 000004237
 日本電気株式会社
 東京都港区芝五丁目7番1号
 (74) 代理人 100088812
 弁理士 ▲柳▼川 信
 (72) 発明者 森永 聡
 東京都港区芝五丁目7番1号 日本電気株
 式会社内
 (72) 発明者 山西 健司
 東京都港区芝五丁目7番1号 日本電気株
 式会社内
 審査官 紀田 馨

最終頁に続く

(54) 【発明の名称】 トピック分析方法及びその装置並びにプログラム

(57) 【特許請求の範囲】

【請求項1】

テキストデータが時間とともに追加されていくような状況のもとで、該データを順次読み込みつつトピックを検出するトピック分析装置であって、

トピックの生成モデルを混合分布モデルで表現し、データのタイムスタンプに応じて過去のデータほど激しく忘却しながら該トピックの生成モデルをオンラインで学習する学習手段と、

前記生成モデルを格納する記憶手段と、

前記記憶手段に格納された複数の候補となるトピックの生成モデルの中で、情報量基準に基づいて最適なトピックの生成モデルを選択して、その混合成分としてトピックを検出する手段と、

特定の時間のトピックの生成モデルの混合成分と、別の時間のトピック生成モデルの混合成分を比較して、新しいトピックの生成と既存のトピックの消滅を判定するトピック形成消滅判定手段と、

を含むことを特徴とするトピック分析装置。

【請求項2】

トピックの生成モデルの各混合成分に対応するトピックの特徴表現を、混合成分のパラメータに基づいて抽出して、各トピックを特徴付けるトピック特徴表現抽出手段を、更に含むことを特徴とする請求項1記載のトピック分析装置。

【請求項3】

テキストデータが時間とともに追加されていくような状況のもとで、該データを順次読み込みつつトピックを検出するトピック分析装置であって、

トピックの生成モデルを混合分布モデルで表現し、データのタイムスタンプに応じて過去のデータほど激しく忘却しながら該トピックの生成モデルをオンラインで学習する学習手段と、

前記生成モデルを格納する記憶手段と、

前記記憶手段に格納された複数の候補となるトピックの生成モデルの中で情報量基準に基づいて最適なトピックの生成モデルを選択してその混合成分としてトピックを検出する手段と、

トピックの生成モデルの各混合成分に対応するトピックの特徴表現を、混合成分のパラメータに基づいて抽出して、各トピックを特徴付けるトピック特徴抽出手段と、
を含むことを特徴とするトピック分析装置。

10

【請求項4】

テキストデータが時間とともに追加されていくような状況のもとで、該データを順次読み込みつつトピックを検出分析するコンピュータによるトピック分析方法であって、

前記コンピュータの学習機能により、トピックの生成モデルを混合分布モデルで表現し、データのタイムスタンプに応じて過去のデータほど激しく忘却しながら該トピックの生成モデルをオンラインで学習して、記憶手段に記憶するステップと、

前記コンピュータのモデル選択機能により、前記記憶手段に記憶された複数の候補となるトピックの前記生成モデルの中で、情報量基準に基づいて最適なトピックの生成モデルを選択して、その混合成分としてトピックを検出するステップと、

20

前記コンピュータのトピック消滅判定機能により、特定の時間のトピックの生成モデルの混合成分と、別の時間のトピック生成モデルの混合成分とを比較して、新しいトピックの生成と既存のトピックの消滅を判定するステップと、
を含むことを特徴とするトピック分析方法。

【請求項5】

前記コンピュータのトピック特徴表現抽出機能により、トピックの生成モデルの各混合成分に対応するトピックの特徴表現を、混合成分のパラメータに基づいて抽出して、各トピックを特徴付けるステップを、更に含むことを特徴とする請求項4記載のトピック分析方法。

30

【請求項6】

テキストデータが時間とともに追加されていくような状況のもとで、該データを順次読み込みつつトピックを検出分析するコンピュータによるトピック分析方法であって、

前記コンピュータの学習機能により、トピックの生成モデルを混合分布モデルで表現し、データのタイムスタンプに応じて過去のデータほど激しく忘却しながら該トピックの生成モデルをオンラインで学習して、記憶手段に記憶するステップと、

前記コンピュータのモデル選択機能により、前記記憶手段に記憶された複数の候補となるトピックの前記生成モデルの中で情報量基準に基づいて最適なトピックの生成モデルを選択して、その混合成分としてトピックを検出するステップと、

前記コンピュータのトピック特徴表現抽出機能により、トピックの生成モデルの各混合成分に対応するトピックの特徴表現を、混合成分のパラメータに基づいて抽出することにより、各トピックを特徴付けるステップと、
を含むことを特徴とするトピック分析方法。

40

【請求項7】

テキストデータが時間とともに追加されていくような状況のもとで、該データを順次読み込みつつトピックを検出する方法をコンピュータに実行させるためのプログラムであって、

前記コンピュータを、トピックの生成モデルを混合分布モデルで表現し、データのタイムスタンプに応じて過去のデータほど激しく忘却しながら該トピックの生成モデルをオンラインで学習して、記憶手段に記憶する機能として動作させる処理と、

50

前記コンピュータを、前記記憶手段に記憶された複数の候補となるトピックの前記生成モデルの中で、情報量基準に基づいて最適なトピックの生成モデルを選択して、その混合成分としてトピックを検出する機能として動作させる処理と、

前記コンピュータを、特定の時間のトピックの生成モデルの混合成分と、別の時間のトピック生成モデルの混合成分とを比較して、新しいトピックの生成と既存のトピックの消滅を判定する機能として動作させる処理と、

を含むことを特徴とする、コンピュータにより読取可能なプログラム。

【請求項 8】

前記コンピュータを、トピックの生成モデルの各混合成分に対応するトピックの特徴表現を、混合成分のパラメータに基づいて抽出して、各トピックを特徴付ける機能として動作させる処理を、更に含むことを特徴とする請求項 7 記載のプログラム。

10

【請求項 9】

テキストデータが時間とともに追加されていくような状況のもとで、該データを順次読み込みつつトピックを検出する方法をコンピュータに実行させるためのプログラムであって、

前記コンピュータを、トピックの生成モデルを混合分布モデルで表現し、データのタイムスタンプに応じて過去のデータほど激しく忘却しながら該トピックの生成モデルをオンラインで学習して、記憶手段に記憶する機能として動作させる処理と、

前記コンピュータを、前記記憶手段に記憶された複数の候補となるトピックの前記生成モデルの中で情報量基準に基づいて最適なトピックの生成モデルを選択して、その混合成分としてトピックを検出する機能として動作させる処理と、

20

前記コンピュータを、トピックの生成モデルの各混合成分に対応するトピックの特徴表現を、混合成分のパラメータに基づいて抽出することにより、各トピックを特徴付ける機能として動作させる処理と、

を含むことを特徴とする、コンピュータにより読取可能なプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明はトピック分析方法及びその装置並びにプログラムに関し、特にテキストマイニングや自然言語処理の分野において、時系列で追加されていくテキスト集合に対して、各時刻の主要なトピックを同定して各トピックの内容および変化を分析するトピック分析方式に関するものである。

30

【背景技術】

【0002】

一括で与えられた時系列のテキストデータに対して、各時刻における主要な表現を抽出する方式としては、例えば、非特許文献 1 に示された方式が知られている。この方式では、テキストデータに現れる単語の中で、その出現頻度が特定の時間期間で上昇しているものが抽出され、その時間期間の開始時刻が主要トピックの出現時刻、期間の終了時刻がそのトピックの消滅時刻、その単語がトピックの内容を表現するものとされていた。

【0003】

40

また、トピックの時系列的変化を可視化する方式としては、非特許文献 2 に開示の方式が知られている。しかし、上記 2 つの方式はいずれもデータが逐次的に与えられる語毎にオンラインでリアルタイムに処理することはできなかった。

【0004】

ある特定の単語を含む文章の時系列の塊を検出する方式としては、非特許文献 3 に示された方式が知られているが、これは異なる単語を使っても同一内容のトピックを表すようなトピックの分析には不向きであり、また、リアルタイムに分析できないという問題があった。

【0005】

有限混合確率モデルを用いてトピックの同定や変化検出を行う方式としては、非特許文

50

献 4 に示された方式が知られているが、いずれもデータが逐次的に与えられる語毎にオンラインでリアルタイムに処理することはできなかった。

【 0 0 0 6 】

リアルタイムに有限混合確率モデルを学習する方式については、非特許文献 5 が知られているが、これはデータの時系列的順序を考慮するが、データの発生時間そのものを反映できないという問題があった。

【 0 0 0 7 】

【非特許文献 1】R. Swan, J. Allan, "Automatic generation of overview timelines," Proc. SIGIR Intl. Conf. Information Retrieval, 2000. S.Harve, B.Hetzler, and L.Norwell: ThemeRiver: Visualizing theme changes over time,

10

【非特許文献 2】Proceedings of IEEE Symposium on Information Visualization, 2000

【非特許文献 3】J.Kleinberg: Bursty and hierarchical structure in streams, Proceedings of KDD2002, pp:91-101, ACM Press, 2003

【非特許文献 4】X.Liu, Y.Gong, W.Xu, and S.Zhu: Document clustering with cluster refinement and model selection capabilities, Proceedings of SIGIR International Conference on Information Retrieval, 2002 や H.Li and K.Yamanishi: Topic analysis using a finite mixture model, Information Processing and Management, Vol.39/4, pp 521-541, 2003

【非特許文献 5】K.Yamanishi, J.Takeuchi, G.Williams, and P.Milne: On-line unsupervised outlier detection using finite mixtures with discounting learning algorithms," in Proceedings of KDD2000, ACM Press, pp:320--324 2000

20

【発明の開示】

【発明が解決しようとする課題】

【 0 0 0 8 】

テキストデータが時間を追って追加されていくような状況で、随時、主要なトピックの内容同定をしたい場合には、従来の多くの方式は、多大な記憶容量と処理時間とが必要になるという問題があった。しかしながら、CRM (Customer Relationship Management) やナレッジマネジメントおよびWEB監視などの目的で、時間的に追加されていくテキストデータのトピックを分析する際には、なるべく少ない記憶容量と処理時間でリアルタイムに分析を行う必要がある。

30

【 0 0 0 9 】

さらに上記の各方式においては、単一のトピックの内容が時間と共に微妙に変化していく場合に、「同じトピックだが内容が微妙に変化している」ことを知ることが出来ない。しかしながら、CRMやWEB監視目的のトピック分析などにおいては、「特定の商品に対する苦情内容の変化」の抽出などのように、単一トピックの内容変化を同定することによって得られる知見は大きい。

【 0 0 1 0 】

本発明の目的は、時間的に追加されていくテキストデータに対して、なるべく少ない記憶容量と処理時間で、随時、リアルタイムに主要トピックの個数および生成と消滅を同定すること、および主要トピックの特徴を抽出すること、それによって、単一トピックの内容が変化した場合にも、それを分析者が知ることが出来るようにしたトピック分析方法およびその装置並びにプログラムを提供することである。

40

【課題を解決するための手段】

【 0 0 1 1 】

本発明によるトピック分析装置は、テキストデータが時間とともに追加されていくような状況のもとで、該データを順次読み込みつつトピックを検出するトピック分析装置であって、トピックの生成モデルを混合分布モデルで表現し、データのタイムスタンプに応じて過去のデータほど激しく忘却しながら該トピックの生成モデルをオンラインで学習する学習手段と、前記生成モデルを格納する記憶手段と、前記記憶手段に格納された複数の候補となるトピックの生成モデルの中で、情報量基準に基づいて最適なトピックの生成モデ

50

ルを選択して、その混合成分としてトピックを検出する手段と、特定の時間のトピックの生成モデルの混合成分と、別の時間のトピック生成モデルの混合成分を比較して、新しいトピックの生成と既存のトピックの消滅を判定するトピック形成消滅判定手段とを含むことを特徴とする。

【0012】

本発明による他のトピック分析装置は、テキストデータが時間とともに追加されていくような状況のもとで、該データを順次読み込みつつトピックを検出するトピック分析装置であって、トピックの生成モデルを混合分布モデルで表現し、データのタイムスタンプに応じて過去のデータほど激しく忘却しながら該トピックの生成モデルをオンラインで学習する学習手段と、前記生成モデルを格納する記憶手段と、前記記憶手段に格納された複数の候補となるトピックの生成モデルの中で情報量基準に基づいて最適なトピックの生成モデルを選択してその混合成分としてトピックを検出する手段と、トピックの生成モデルの各混合成分に対応するトピックの特徴表現を、混合成分のパラメータに基づいて抽出して、各トピックを特徴付けるトピック特徴抽出手段とを含むことを特徴とする。

10

【0016】

本発明によるトピック分析方法は、テキストデータが時間とともに追加されていくような状況のもとで、該データを順次読み込みつつトピックを検出分析するコンピュータによるトピック分析方法であって、前記コンピュータの学習機能により、トピックの生成モデルを混合分布モデルで表現し、データのタイムスタンプに応じて過去のデータほど激しく忘却しながら該トピックの生成モデルをオンラインで学習して、記憶手段に記憶するステップと、前記コンピュータのモデル選択機能により、前記記憶手段に記憶された複数の候補となるトピックの前記生成モデルの中で、情報量基準に基づいて最適なトピックの生成モデルを選択して、その混合成分としてトピックを検出するステップと、特定の時間のトピックの生成モデルの混合成分と、別の時間のトピック生成モデルの混合成分を比較して、新しいトピックの生成と既存のトピックの消滅を判定するトピック形成消滅判定手段とを含むことを特徴とする。

20

【0017】

本発明による他のトピック分析装置は、テキストデータが時間とともに追加されていくような状況のもとで、該データを順次読み込みつつトピックを検出するトピック分析装置であって、トピックの生成モデルを混合分布モデルで表現し、データのタイムスタンプに応じて過去のデータほど激しく忘却しながら該トピックの生成モデルをオンラインで学習する学習手段と、前記生成モデルを格納する記憶手段と、前記記憶手段に格納された複数の候補となるトピックの生成モデルの中で情報量基準に基づいて最適なトピックの生成モデルを選択してその混合成分としてトピックを検出する手段と、トピックの生成モデルの各混合成分に対応するトピックの特徴表現を、混合成分のパラメータに基づいて抽出して、各トピックを特徴付けるトピック特徴抽出手段とを含むことを特徴とする。

30

【0021】

本発明によるプログラムは、テキストデータが時間とともに追加されていくような状況のもとで、該データを順次読み込みつつトピックを検出する方法をコンピュータに実行させるためのプログラムであって、前記コンピュータを、トピックの生成モデルを混合分布モデルで表現し、データのタイムスタンプに応じて過去のデータほど激しく忘却しながら該トピックの生成モデルをオンラインで学習して、記憶手段に記憶する機能として動作させる処理と、前記コンピュータを、前記記憶手段に記憶された複数の候補となるトピックの前記生成モデルの中で、情報量基準に基づいて最適なトピックの生成モデルを選択して、その混合成分としてトピックを検出する機能として動作させる処理と、前記コンピュータを、特定の時間のトピックの生成モデルの混合成分と、別の時間のトピック生成モデルの混合成分とを比較して、新しいトピックの生成と既存のトピックの消滅を判定する機能として動作させる処理とを含むことを特徴とする。

40

【0022】

本発明による他のプログラムは、テキストデータが時間とともに追加されていくような

50

状況のもとで、該データを順次読み込みつつトピックを検出する方法をコンピュータに実行させるためのプログラムであって、前記コンピュータを、トピックの生成モデルを混合分布モデルで表現し、データのタイムスタンプに応じて過去のデータほど激しく忘却しながら該トピックの生成モデルをオンラインで学習して、記憶手段に記憶する機能として動作させる処理と、前記コンピュータを、前記記憶手段に記憶された複数の候補となるトピックの前記生成モデルの中で情報量基準に基づいて最適なトピックの生成モデルを選択して、その混合成分としてトピックを検出する機能として動作させる処理と、前記コンピュータを、トピックの生成モデルの各混合成分に対応するトピックの特徴表現を、混合成分のパラメータに基づいて抽出することにより、各トピックを特徴付ける機能として動作させる処理とを含むことを特徴とする。

10

【0026】

本発明の作用を述べる。本発明では、各テキストを文書ベクトルとして表現し、その生成モデルとして混合分布モデルを用いる。混合分布の一つのコンポーネントが一つのトピックに対応するとする。混合分布モデルはコンポーネントの個数等が異なる複数のものが保持される。新規テキストデータが追加されるたびに、学習手段によって各モデルのパラメータが追加学習され、モデル選択手段によって情報量基準に基づいて最も適切なモデルが選択される。選択されたモデルの各コンポーネントが主要なトピックを表している。また、モデル選択手段によってどのモデルが選択されるかが変化した場合には、トピック形成消滅判定手段により以前に選択されていたモデルと今回選択されたモデルの比較が行われ、どれが新たに形成されたトピックであるか、どのトピックが消滅したのかが判定される。

20

【0027】

さらに本発明では、モデル選択手段によって選択されたモデルの各トピック、トピック生成/消滅判定手段によって判定された新たに生成されたトピック/消滅したトピックに関して、トピック特徴表現抽出手段が該当する混合分布のパラメータから、そのトピックの特徴表現を抽出し出力される。

【0028】

複数の混合分布モデルを全て独立に学習し選択するのではなく、一つもしくは複数の上位モデルを学習し、学習された上位モデルからサブモデル生成手段によって複数のサブモデルを生成し、モデル選択手段によって、その中から適切なモデルを選択するのも良い。さらに、サブモデルを独立に生成して保持するのではなく、サブモデル生成選択手段によって、特定のサブモデルの情報量基準を上位モデルから直接に計算し、最も適切なサブモデルを選択するのもよい。

30

【0029】

学習手段による各モデルのパラメータの追加学習においては、到着順の早いテキストデータに比べて到着順が後ろのテキストデータの内容を重視するようにしてもよい。さらに、テキストデータにタイムスタンプが付随している場合には、到着順のみならずタイムスタンプの内容を利用して、古いテキストデータに比べて最近のテキストデータほど内容を重視するようにしてもよい。

【0030】

モデル選択手段およびサブモデル生成選択手段において適切なモデルを選択する際に、新たに入力されたテキストデータを用いて追加学習した前後の分布間の距離や、追加学習する前の分布において前記入力テキストデータが発生するのはどれくらい稀なのか、を各モデルに対して計算し、それを最小にするモデルを選択するのもよい。さらに、これらを計算した結果をモデルの次元数で割ったものや、特定の時刻からの値の累積値、最近の値を重視するように重み付けした平均値などを計算するのもよい。

40

【0031】

トピック形成/消滅判定手段において、以前に選択されていたモデル(旧モデル)と今回選択されたモデル(新モデル)を比較する際に、旧モデルに含まれるコンポーネントと新モデルに含まれるコンポーネントの全ての組み合わせのペアについて類似度を計算し、

50

どの旧モデルのコンポーネントとも類似度が低い新モデルのコンポーネントを新たに形成されたトピックと判定、どの新モデルのコンポーネントとも類似度が低い旧モデルのコンポーネントを消滅したトピックと判定してもよい。コンポーネント間の類似度は、平均値間の距離や、分布の同一性検定におけるP値を用いてもよい。モデルが上位モデルから生成されたサブモデルである場合は、コンポーネント間の類似度として上位モデルにおける同一のコンポーネントから生成されているかどうかを用いてもよい。

【0032】

トピック特徴表現抽出手段においては、各トピックを表すコンポーネントの確率分布に従ってテキストデータを発生させ、テキストデータを入力とする公知の特徴抽出技術を用いて各トピックの特徴表現を抽出してもよい。前記公知の特徴抽出技術で必要となるテキストデータの各種統計量が、コンポーネントのパラメータから計算できる場合は、その値を使って特徴抽出してもよい。サブ分布生成手段においては、上位モデルの幾つかのコンポーネントをコンポーネントとする混合分布をサブ分布としてもよい。

10

【発明の効果】

【0033】

本発明の第一の効果は、時系列のテキストデータを複数の混合分布でモデル化し、忘却型逐次学習アルゴリズムによるパラメータ学習とモデル選択によって、主要トピックおよびその生成/消滅を、少ない記憶容量と処理時間で随時同定することができるということである。この際、データのタイムスタンプを利用して、古いものほど、その効果を失いながらトピック構造を同定することができる。また、テキストデータが追加されるごとに新しい単語が出現して、その表現ベクトルの次元が上がっても、これに対応して、最適な主要トピックを同定することができる。

20

【0034】

また、本発明の第二の効果は、学習された混合分布のパラメータから各トピックの特徴表現を同定することによって、トピックの内容を随時抽出できること、それによって、単一トピックの内容が変化した場合にも、それを分析者が知ることができるということである。

【発明を実施するための最良の形態】

【0035】

以下に、図面を参照して本発明の実施の形態について詳細に説明する。図1は本発明の第一の実施の形態にかかるトピック分析装置の構成を示すブロック図である。本トピック分析装置は、全体としてコンピュータからなり、テキストデータ入力手段1、学習手段2 1, ..., 2 n、混合分布モデル(モデル記憶手段) 3 1, ..., 3 n、モデル選択手段4、トピック形成/消滅判定手段5、トピック特徴表現抽出手段6、出力手段8を含んでいる。

30

【0036】

テキストデータ入力手段1は、コールセンターのコンタクト内容や、Webから収集した監視対象ページの内容、新聞記事の内容などテキスト(文字情報)を入力する手段であり、対象とするデータを一括して入力するだけでなく、データが発生したり収集されたりする毎に、追加的に入力することも可能となっているものである。また、入力されたテキストは公知の形態素解析技術や構文解析技術によって分解され、さらに公知の属性選択技術や重み付け技術を用いることで、後記モデル3 1 ~ 3 nが対象とするデータ形式に変換される。

40

【0037】

例えば、全ての単語のうち、名詞だけを取り出し、それらを w_1, \dots, w_N として、それらのテキスト中の頻度を $t_f(w_1), \dots, t_f(w_N)$ としてベクトル $(t_f(w_1), \dots, t_f(w_N))$ をテキストデータの表現としたり、全体のテキストの数を M 、単語 w_i を含むテキストの数を $d_f(w_i)$ として、 $t_f - i d_f$ 値である、

$$t_f - i d_f(w_i) = t_f(w_i) \times \log(M / d_f(w_i))$$

を各要素とするベクトル、

50

(t f - i d f (w i) , ... , t f - i d f (w N))

をテキストデータの表現としたりする。これらを構成する際に、予め頻度がしきい値に達しないものは最初から要素に入れないなどの前処理を行うこともあり得る。

【 0 0 3 8 】

本テキストデータ入力手段 1 は、テキストデータを操作入力するためのキーボードや、コールセンターデータベースの内容を逐次転送するプログラム、Web 上のテキストデータをダウンロードするアプリケーションなどの一般的な情報入力手段により構成される。

【 0 0 3 9 】

学習手段 2 1 ~ 2 n は、テキストデータ入力手段 1 によって入力されたテキストデータに基づき、混合分布 3 1 ~ 3 n を更新する手段である。混合分布 3 1 ~ 3 n は、入力されるテキストデータの従う確率分布の候補として、テキストデータ入力手段 1 によって入力されたテキストデータに基づき推定されたものである。

【 0 0 4 0 】

一般に、確率モデルの考え方では、与えられたデータ x は、ある確率変数の実現値とみなされる。特に、この確率変数の確率密度関数が有限次元のパラメータ a を持つ固定された関数形 $f(x; a)$ を持つと仮定すると、その確率密度関数族、

$$F = \{ f(x; a) \mid a \in A \}$$

をパラメトリック確率モデルという。なお、 A は a のとり得る値の集合である。また、データ x に基づきパラメータ a の値を推測することを推定という。例えば、 $\log f(x; a)$ を a の関数(対数尤度関数)とみなし、これを最大にする a を推定値とする最尤推定法などが一般的である。

【 0 0 4 1 】

また、複数の確率モデルの線形結合、

$$\begin{aligned} M &= \{ f(x; C_1, \dots, C_n, a_1, \dots, a_n) \\ &= C_1 * f_1(x; a_1) + \dots + C_n * f_n(x; a_n) \mid a_i \in A_i, \\ &C_1 + \dots + C_n = 1, \quad C_i > 0 \quad (i = 1, \dots, k) \} \end{aligned}$$

によって与えられる確率モデル M を混合モデル、その確率分布を混合分布、線形結合の対象となった元の各分布をコンポーネント、 C_i を i 番目のコンポーネントの混合比率とよぶ。これは、 y を 1 から n までの整数を値域とする確率変数とし、

$$Pr\{y = i\} = C_i, \quad f(x \mid y = i) = f_i(x; a_i)$$

を満たす確率変数 $z = (y, x)$ に対して、 y を隠れ変数として x のみをモデル化したものと同一である。

【 0 0 4 2 】

ただし、 $y = i$ という条件の下での x の条件付密度関数を $f(x \mid y = i)$ としている。また、後の記述の簡単化のために、 $z = (y, x)$ の確率密度関数を、

$$g(z; C_1, \dots, C_n, a_1, \dots, a_n)$$

とする。

【 0 0 4 3 】

本発明においては、モデル 3 1 ~ 3 n は、コンポーネント数やコンポーネントのパラメータが異なる混合モデルであるとし、各コンポーネントは特定の主要なトピックについて述べているテキストデータの従う確率分布であるとする。すなわち、与えられたモデルにおいて、コンポーネントの個数はテキストデータ集合の中の主要トピックの数を表し、各コンポーネントが各主要トピックに相当することになる。

【 0 0 4 4 】

混合モデルに対して、与えられたデータに基づいて最尤推定を行うことは非常に大きな計算量を必要とするが、計算量を節約して近似解を求める方法として、EM (Expectation Maximization) アルゴリズムがよく知られている。この EM アルゴリズムにおいては、対数尤度を直接に最大化するのではなく、隠れ変数 y の値の事後分布の計算と、 y の値で条件付けした x の対数尤度の前記事後分布による平均値 $E_y[\log g(x \mid y)]$ の最大化を繰り返すことで、混合分布の各パラメータの推定が行われる。ただし、 y の前記

10

20

30

40

50

事後分布による平均値を $E y [*]$ としている。

【 0 0 4 5 】

さらに、データが一括で与えられるのではなく、逐次的に追加到着する状況で、混合分布のパラメータ推定結果をデータ追加時に更新していく逐次型のEMアルゴリズムも公知となっている。特に、非特許文献5では、データの到着順序が考慮され、最近到着したものが重要視され、昔に到着したデータの影響は徐々に軽くなっていく手法が記されている。これは、到着したデータの総数を L 個とし l 番目のデータを x_l 、そのときの隠れ変数を y_l とした場合に、 y_l の事後分布の計算と、最近到着したものの重みを大きくした対数尤度、

$$E y_l [(1 - r)^{(L-l)} r \log g(y_l | x_l)]$$

10

の最大化を逐次的に行うものである。

【 0 0 4 6 】

ただし、 α は $l = 1 \sim L$ の和を表すとし、 $E y_l [*]$ は y_l の事後分布による平均とする。上記の特別な場合として $r = 0$ としたものが、データの到着順序による重み付けをしない逐次型のEMアルゴリズムである。

【 0 0 4 7 】

本発明の学習手段 2 1 ~ 2 n は前記の逐次型EMアルゴリズムによって、テキストデータ入力手段 1 からデータが与えられるたびに、モデル 3 1 ~ 3 n における混合分布の推定結果を更新する。さらに、テキストデータにタイムスタンプが付随している場合は、

$$E y_l [(1 - r)^{(L-t_l)} r \log g(x_l, y_l | y_l)]$$

20

を最大化するように逐次学習をおこなってもよい。ただし、 l 番目のデータのタイムスタンプを t_l としている。こうすることによって、データの到着間隔が不ぞろいである場合にも、時間的に最近のデータを重要視し古いデータの影響を軽くするようにコンシスタントに推定が行われる。

【 0 0 4 8 】

例えば、混合モデルとして、各コンポーネントがガウス分布であるような場合を考えると、 i 番目のコンポーネントは平均 μ_i 、分散共分散行列 Σ_i をパラメータとするガウス密度関数として、

$$\left(\frac{1}{(2\pi)^{d/2} |\Sigma_i|} \right) \exp \left[- \frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) \right]$$

30

で表される。コンポーネントの数は k 個あるとし、 i 番目のコンポーネントの混合比率を α_i とする。

【 0 0 4 9 】

この場合、 t_{old} 時刻のデータを x_n とし、 t_{new} 時刻に新しいデータ x_{n+i} を入力としたとき、 i 番目のコンポーネントの更新前の平均パラメータ、分散共分散行列パラメータ、混合比率をそれぞれ μ_i^{old} 、 Σ_i^{old} 、 α_i^{old} とし、更新後のそれらを μ_i^{new} 、 Σ_i^{new} 、 α_i^{new} は、例えば以下のように計算することができる。ここで、 d 、 W_{ij} 、 S_i は助変数である。

【 0 0 5 0 】

【 数 1 】

40

$$P = \frac{1}{\sum_{l=1}^k \exp \left\{ \log \xi_l^{old} + \log \phi \left(x_{n+1} \mid \mu_l^{old}, \Sigma_l^{old} \right) - \log \xi_l^{old} - \log \phi \left(x_{n+1} \mid \mu_l^{old}, \Sigma_l^{old} \right) \right\}}$$

【数 2】

$$W_{in+1}^{new} = WA\left(P, \frac{1}{k} \mid 1, \alpha\right)$$

ここに、 α はユーザ指定の定数である。

【0051】

【数 3】

$$\mu_i^{new} = WA\left(\mu_i^{old}, x_{n+1} \mid \xi_i^{old} d^{old}, \lambda^{-(t_{new}-t_{old})} W_{in+1}^{new}\right)$$

10

ここに、 λ はユーザ指定の定数（忘却率）である。

【数 4】

$$S_i^{new} = WA\left(S_i^{old}, x_{n+1} x_{n+1} \mid \xi_i^{old} d^{old}, \lambda^{-(t_{new}-t_{old})} W_{in+1}^{new}\right)$$

20

【数 5】

$$\Sigma_i^{new} = S_i^{new} - \mu_i^{new} \mu_i^{new}$$

【数 6】

$$\xi_i^{new} = WA\left(\xi_i^{old}, W_{in+1}^{new} \mid d^{old}, \lambda^{-(t_{new}-t_{old})}\right)$$

30

【数 7】

$$d^{new} = \lambda^{-(t_{new}-t_{old})} d^{old} + 1$$

40

【0052】

ただし、上記では、表記の簡単化のために、

(式1 * 式3 + 式2 * 式4) / (式3 + 式4)

と書くところを、

WA(式1, 式2 | 式3, 式4)

として表している。

【0053】

モデル選択手段4では、入力されるテキストデータの従う確率分布の候補であるモデル31～3nのそれぞれに対し、テキストデータ入力手段1によって入力されたテキストに基づいて情報量基準の値が計算され、最も適切なモデルが選択される。例えば、Wをウイ

50

ンドウの大きさとし、 t 番目のデータのベクトル表現の次元を d_t とし、 $P^{(t)}(x|k)$ を k 個のコンポーネントからなる混合分布で、 t 番目のデータが入力されてから逐次的にパラメータを更新したものであるとすると、 n 番目のデータを受け取ったときの情報量基準の値 $I(k)$ は、例えば、以下のように計算できる。

$$I(k) = (1/W) \sum_{t=n-W}^n (-\log P^{(t)}(x_t|k)) / d_t$$

【0054】

この値を最小化するようなコンポーネント個数 k が最適なコンポーネント数であり、それを構成するコンポーネントが主要トピックを表すコンポーネントであると同定することができる。この基準の値は、入力テキストデータが追加されるごとに新しい単語が出現して、その表現ベクトルの次元が上がっても、これに対応して計算できるものである。 $P^{(t)}(x_t|k)$ を構成するコンポーネントは、独立なコンポーネントであっても、上位の混合モデルのサブコンポーネントであってもよいものとする。

10

【0055】

トピック形成/消滅判定手段5では、モデル選択手段4によって選択されモデルが変化した場合、新たに選択されたモデルのコンポーネントの中で、以前に選択されていたモデルには近いコンポーネントが存在しないものを、「新たに形成されたトピック」、逆に新しいモデルにおいて近いコンポーネントが存在しない古いモデルのコンポーネントを、「消滅したトピック」と判定し、出力手段7に出力する。コンポーネント間の近さの尺度としては、分布の同一性検定における P 値や、二つの確率分布の近さを計る量として公知の KL (Kullback Leibler) ダイバージェンス等を用いればよい。あるいは、さらに簡単に二つの確率分布の平均値の差などを用いても良い。

20

【0056】

トピック特徴抽出手段6は、モデル選択手段4によって選択されたモデルに対して、各コンポーネントの特徴を抽出し、該当トピックの特徴表現として出力手段7に出力する。特徴表現を抽出するには、単語の情報利得を計算して、その大きいものを抽出する方法を用いることができる。情報利得は、例えば、以下のように計算する。

【0057】

t 番目のデータが与えられたときに、全体のデータの数を t とし、全データの中で指定された単語 w を含むデータの数を m_w 、これに含まないデータの数を m'_w 、ある指定したコンポーネント(かりに i 番目とする)から発生したテキストの数を t_i 、単語 w を含むデータの中で i 番目のコンポーネントから発生したデータ数を m_w^+ 、単語 w を含まないデータの中で i 番目のコンポーネントから発生したデータ数を m'^w_+ とするとき、 $I(A, B)$ を情報量尺度として、 w の情報利得を、

30

$$IG(w) = I(t, t_i) - (I(m_w, m_w^+) + I(m'_w, m'^w_+))$$

のように計算する。

【0058】

ここで、 $I(A, B)$ の計算式としては、エントロピー、確率的コンプレキシティ、拡張型確率的コンプレキシティなどを用いることができる。エントロピーは、

$$I(A, B) = A H(B/A) = A (B \log(B/A) + (A - B) \log((A - B)/A))$$

40

で表されるものであり、確率的コンプレキシティは、

$$I(A, B) = A H(B/A) + (1/2) \log(A/2)$$

で表されるものであり、拡張型確率的コンプレキシティは、

$$I(A, B) = \min\{B, A - B\} + c (A \log A)^{1/2}$$

で表されるものである。

【0059】

また、 $IG(w)$ の代わりに情報利得として 自乗検定量、

$$(m_w + m'_w) \times (m_w^+ (m'_w - m'^w_+) - (m_w - m_w^+) m'_w) \times ((m_w^+ + m'^w_+) \times (m_w - m_w^+ + m'_w - m'^w_+) m_w m'_w)$$

50

)⁻¹

を用いることもできる。

【0060】

各 i について、 i 番目のコンポーネントに対し、各 w について、上記情報利得を計算し、大きい順に指定された数の言葉を抽出することにより、特徴語を抽出することができる。また、しきい値を予め与えて、そのしきい値以上の情報利得を与える言葉を抽出することにより、特徴語を抽出することができる。上記情報利得を計算するのに必要な統計量は、 t 番目のデータが与えられたときには、各 i と w に対し、 t 、 t_i 、 m_w 、 m'_w 、 m_w^+ 、 m'_w^+ であるが、これらはデータが与えられる毎にインクリメンタルに計算できる。

10

【0061】

本学習手段およびモデルは、CPUなどのマイクロプロセッサおよびその周辺回路と、モデル31～3nを記憶している記憶装置、およびこれらの動作を統括するプログラムとが協働することにより構成されている。

【0062】

図2は本発明の動作を示すフローチャートである。まず、ステップ101では、テキストデータ入力手段によってテキストデータが入力され、以降のステップでの処理の対象とするデータ形式に変換される。続いて、ステップ102では、前記変換されたテキストデータに基づき、学習手段によってモデルのパラメータ推定の更新を行う。これによって、各モデルにおいては今回入力されたデータの値を反映した新しいパラメータの値を保持することになる。

20

【0063】

次に、ステップ103においては、保持されている複数のモデルの中から、これまでに入力されたテキストデータを鑑みて最も適切なモデルがモデル選択手段により選択される。選択されたモデルにおける混合分布の各コンポーネントが主要なトピックに対応している。

【0064】

ステップ104においては、どのモデルが選択されたかが今回のデータ入力の結果、前回のそれと変化したかが判定される。今回と前回で選択されるモデルが変わらなかった場合は、前回までのテキストデータにおける主要トピックに対して、今回のデータを入力することで新たに主要トピックの形成や消滅がおきなかったことを意味する。逆に、選択されるモデルが変化した場合は、一般に混合分布を構成するコンポーネントの数が増加しており、何らかの新規トピックの形成もしくは消滅が起きていることを意味する。

30

【0065】

そこで、ステップ105においては、今回選択されたモデルのコンポーネントの中で、前回選択されていたモデルのコンポーネントのどれとも近いものがないものをトピック形成/消滅判定手段により同定し、新規に形成された主要トピックを表すコンポーネントであるとする。同様に、ステップ106においては、前回選択されていたモデルのコンポーネントの中で、今回選択されたモデルのコンポーネントのどれとも近いものがないものを同定し、主要でなくなったトピックを表すコンポーネントであるとする。

40

【0066】

ステップ107では、今回選択されたモデルの各コンポーネントおよび新規形成/消滅したとされたコンポーネントの特徴がトピック特徴抽出手段により抽出され、該当するトピックの特徴表現とされる。新たにテキストデータが入力された場合は、ステップ101に戻り、一連の処理がなされる。また、ステップ103から107の処理は、入力される各テキストデータに対して毎回行う必要は無く、主要トピックの同定や新規形成/消滅トピックの同定を行うように、ユーザーなどから指示された場合やタイマーなどで指定された時刻にだけ行うようにしてもよい。

【0067】

図3は本発明の第二の実施形態にかかるトピック分析装置の構成を示すブロック図であ

50

り、図1と同等部分は同一符号により示している。第一の実施形態との違いは、モデル選択手段でモデル選択する際の候補となるモデルが、上位モデルの複数のサブモデルである場合になっていることである。サブモデル生成手段9によって生成されたサブモデルに対して、第一の実施の形態と同様のモデル選択を行う。例えば、上位モデルとしては比較的多数のコンポーネントをもつ混合モデルを想定し、サブモデルとしてはそのコンポーネントを幾つか取り出して混合モデルを作った場合が相当する。

【0068】

このような構成にすることで、並列に複数のモデルを保持する必要と、それぞれを学習手段によって更新する必要が無くなり、処理に必要な記憶容量や計算量を縮減することができる。また、トピック形成/消滅判定手段においても、二つのコンポーネントの間の近さの尺度として、「上位モデルで同一コンポーネントから生成されたかどうか」を採用することにより、確率分布間の距離等を尺度とする場合に比べて必要な計算量を縮減することができる。

10

【0069】

図4は本発明の第三の実施形態にかかるトピック分析装置の構成を示すブロック図であり、図1と同等部分は同一符号にて示している。ここでも、モデル選択手段でモデル選択する際の候補となるモデルが、上位モデルの複数のサブモデルとして与えられるが、第二の実施形態との違いは、複数のサブモデルを並列に計算するのではなく、サブモデル生成選択手段41によって、順番に情報量基準を計算し、最も適切なサブモデルを選択することにある。このような構成にすることで、サブモデル全てを保持しておく必要も無くなり、必要な記憶容量をさらに縮減することができる。

20

【0070】

図5に本発明への入力データの例を示す。特定のタイプの電気製品に関して議論を行うWEB上の掲示板に対する監視データで、掲示板への書き込みが行われた日付時刻を付加された書き込み内容(テキストデータ)が1レコードを構成している。WEB掲示板自体は投稿が随時行われるので、時間的にデータが随時追加されていくことになる。スケジュールに従って動くプログラムあるいは掲示板サーバー自体等により、新規に追加されたデータが本発明のトピック分析装置に入力され、各処理が行われるとする。

【0071】

図6は、ある特定の時刻までデータが入力された場合の、本発明によるトピック分析の出力例である。各列が各主要トピックに相当し、モデル選択手段によって選択されたモデルにおける各コンポーネントに対して、トピック特徴表現抽出手段の出力を記載したものである。この分析例では、選択されたモデルには二つのコンポーネントがあり、一つ目のコンポーネントは、「商品XX」、「遅い」、「メール」などを特徴表現とする主要トピック、二つ目のコンポーネントは、「音」、「ZZ」、「よい」などを特徴表現とする主要トピックとなっている。

30

【0072】

図7は、さらに特定の時刻までデータ入力が進んだ場合の、本発明によるトピック分析の出力例である。ただし、本出力例はこの時刻でモデル選択手段によってどのモデルが選択されたかが変化した場合を記載している。本出力例で、トピック形成/消滅判定手段により新規形成と判定されたトピックには「主要トピック：新規」、消滅したと判定されたトピックには「消滅トピック」、新しく選択されたモデルのコンポーネントで、以前のモデルのコンポーネントに近いものが存在するトピックには「主要トピック：継続」と列名がついている。

40

【0073】

「商品XX」を特徴語とするトピックは、「主要トピック：継続」の列名を持つので、以前から主要であったトピックである。しかしながら、図6の「商品XX」のトピックと比較すると、「メール」の代わりに「ウイルス」が特徴語となっており、同じトピックでも内容が変化してきていることを分析者が見て取ることが可能となっている。

【0074】

50

「音」や「ZZ」を特徴語としていたトピックは図6では主要トピックであったが、図7では「消滅トピック」として出力されている。図7の分析を行った時点で、このトピックが消滅したことが見て取れる。逆に、「新WW」などを特徴表現とするトピックは「主要トピック：新規」と同定されており、この時点であらたに主要トピックとなったことを分析者が見て取ることが出来る。

【図面の簡単な説明】

【0075】

【図1】本発明の第一の実施形態に係るトピック分析装置の構成を表すブロック図である。

【図2】本発明の第一の実施形態に係るトピック分析装置の動作を示すフロー図である。

10

【図3】本発明の第二の実施形態に係るトピック分析装置の構成を表すブロック図である。

【図4】本発明の第三の実施形態に係るトピック分析装置の構成を表すブロック図である。

【図5】本発明への入力データ例である。

【図6】本発明の分析結果出力例（その1）である。

【図7】本発明の分析結果出力例（その2）である。

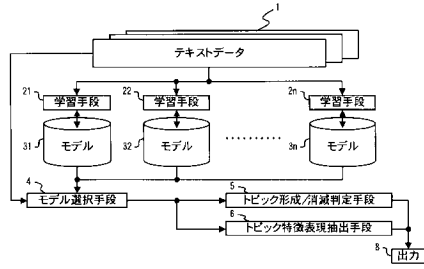
【符号の説明】

【0076】

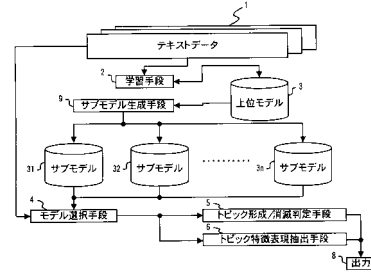
- 1 テキストデータ入力手段
- 2 1 ~ 2 n 学習手段
- 3 1 ~ 3 n モデル（または上位モデル、サブモデル）
- 4 モデル選択手段
- 5 トピック形成 / 消滅判定手段
- 6 トピック特徴表現抽出手段
- 8 出力手段
- 9 サブモデル生成手段
- 4 1 サブモデル生成選択手段

20

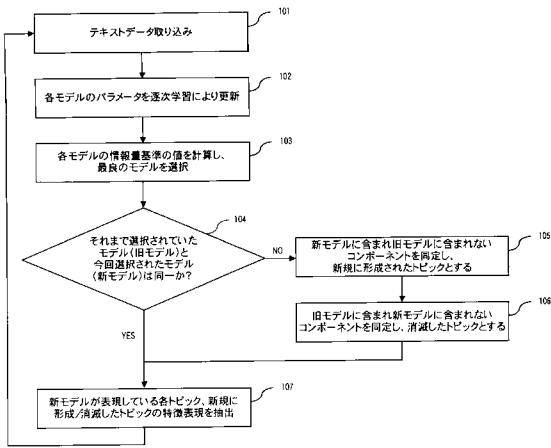
【図1】



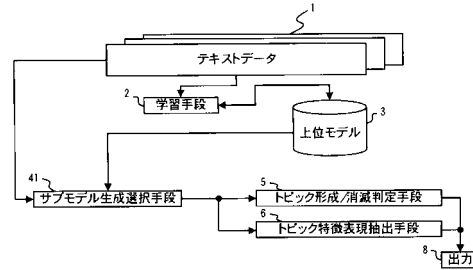
【図3】



【図2】



【図4】



【図5】

日付時刻	テキスト
○月○日○時○分	XXはかえって遅くなってしまった。
○月○日○時○分	テキストだけならモノクロのXXもよい。
○月○日○時○分	商品YYに役に立つようになっています。
○月○日○時○分	商品ZZ音がいいですよ。
○月○日○時○分	○○さんが思っている商品YYあたりいかがでしょうか。
○月○日○時○分	今度のYYは画面の色がいいですよわ。
○月○日○時○分	入力YYがいいのは分かっていますが高すぎます。
○月○日○時○分	XXは遅好みますますね。
○月○日○時○分	XXは遅いのはどうしようもない。
○月○日○時○分	本当だ、XXは遅すぎで使えない。
○月○日○時○分	本当にメールだけならXXだけど、他は死ぬほど遅いよ。
○月○日○時○分	だれかZZの話をしてくれ。
.....

【図7】

	主要トピック:継続	消滅トピック	主要トピック:新規
特徴表現	商品XX ウイルス 遅い	音 ZZ よい	新WW 軽い キーボード

【図6】

	主要トピック:継続	主要トピック:継続
特徴表現	商品XX 遅い メール	音 ZZ よい

フロントページの続き

- (56)参考文献 特開2004-054370(JP,A)
特開2003-330922(JP,A)
特開2001-101154(JP,A)
石川 佳治, 忘却の概念に基づくインクリメンタルな文書クラスタリング手法, 情報処理学会研究報告, 日本, 社団法人情報処理学会, 2001年 7月18日, Vol.2001 No.70, 第313頁乃至第320頁
H.Li, Topic Analysis Using a Finite Mixture Model, Information Processing and Management, 2003年, Vol.39/4, pp.521-543, 以下のWebページより入手 http://www.sciencedirect.com/science?_ob=ArticleURL&_udi=B6VC8-4607XFJ-1&_user=1583244&_rdoc=1&_fmt=&_orig=search&_sort=d&view=c&_acct=C000053883&_version=1&_urlVersion=0&_userid=1583244&md5=dd4785d37b652bdb96b2c5a56f28ad1f
K.Yamanishi, On-line Unsupervised Outlier Detection Using Finite Mixtures with Discounting Learning Algorithms, in Proceedings of the Sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining(KDD2000, ACM Press, 2000年, pp.320-324, 著者Webページより入手 (http://www.nec.co.jp/rd/DTmining/members/yamanishi/index_e.html))
山西 健司, 情報論的学習理論に基づくマイニング技術 - 外れ値検出とテキストマイニングを例に -, 第44回 人工知能基礎論研究会資料, 日本, 社団法人人工知能学会, 2001年 3月 8日, (SIG-FAI-A004), 第35頁乃至第40頁

(58)調査した分野(Int.Cl., DB名)

G06F 17/30

JSTPlus(JDreamII)