

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5744228号
(P5744228)

(45) 発行日 平成27年7月8日(2015.7.8)

(24) 登録日 平成27年5月15日(2015.5.15)

(51) Int. Cl.		F I			
G06F 13/00	(2006.01)	G06F 13/00	540E		
G06F 17/27	(2006.01)	G06F 17/27			

請求項の数 16 (全 14 頁)

(21) 出願番号 特願2013-545039 (P2013-545039)
 (86) (22) 出願日 平成23年12月26日(2011.12.26)
 (65) 公表番号 特表2014-502754 (P2014-502754A)
 (43) 公表日 平成26年2月3日(2014.2.3)
 (86) 国際出願番号 PCT/CN2011/084699
 (87) 国際公開番号 W02012/083892
 (87) 国際公開日 平成24年6月28日(2012.6.28)
 審査請求日 平成26年3月28日(2014.3.28)
 (31) 優先権主張番号 201010621142.1
 (32) 優先日 平成22年12月24日(2010.12.24)
 (33) 優先権主張国 中国(CN)

早期審査対象出願

(73) 特許権者 507231932
 北大方正集▲団▼有限公司
 PEKING UNIVERSITY F
 OUNDER GROUP CO., L
 TD
 中華人民共和国北京市▲海▼淀区成府路2
 98号中▲関▼村方正大厦5▲層▼
 5 Floor, Zhongguanc
 un Founder Building
 , No. 298, Chengfu R
 oad, Haidian Distri
 ct, Beijing 100871,
 China

最終頁に続く

(54) 【発明の名称】 インターネットにおける有害情報の遮断方法と装置

(57) 【特許請求の範囲】

【請求項1】

インターネットにおける有害情報の遮断方法であって、
 遮断待ちテキスト情報、システムプレリサーチモデル情報及びユーザーフィードバック
 モデル情報を取得するステップと、

前記遮断待ちテキスト情報を前処理するステップと、

前記前処理された遮断待ちテキスト情報と前記システムプレリサーチモデル情報とに対
 して特徴情報マッチングし、第一マッチング結果を取得するステップと、

前記前処理された遮断待ちテキスト情報と前記ユーザーフィードバックモデル情報とに
 対して特徴情報マッチングし、前記第一マッチング結果から独立した第二マッチング結果
 を取得するステップと、

第一マッチング結果と第二マッチング結果とが一致しているか否かに基いて、前記遮断
 待ちテキスト情報を遮断するステップと、

を備えることを特徴とする方法。

【請求項2】

更に、

前記システムプレリサーチモデル情報のコーパス及び前記ユーザーフィードバックモデ
 ル情報のコーパスを取得するステップを備えることを特徴とする請求項1に記載の方法。

【請求項3】

前記ユーザーフィードバックモデル情報のコーパスには、ユーザーフィードバックコー

10

20

パス及び/または被遮断コーパスが含まれることを特徴とする請求項 2 に記載の方法。

【請求項 4】

更に、

前記ユーザーフィードバックモデル情報のコーパス量及びそれに対応する閾値を取得するステップと、

前記ユーザーフィードバックモデル情報のコーパス量及びそれに対応する閾値に基づいて、前記ユーザーフィードバックモデル情報を更新するステップと、

を備えることを特徴とする請求項 3 に記載の方法。

【請求項 5】

前記遮断待ちテキスト情報を前処理するステップにおいては、

前記遮断待ちテキスト情報に対してセグメント処理をし、

前記セグメント処理がされた特徴項候補量を統計することを特徴とする請求項 2、3 または 4 に記載の方法。

10

【請求項 6】

前記前処理された遮断待ちテキスト情報と前記システムプレリサーチモデル情報とに対して特徴情報マッチングし、第一マッチング結果を取得するステップにおいては、

前記前処理された遮断待ちテキスト情報及び前記システムプレリサーチモデル情報を取得し、

前記前処理された遮断待ちテキスト情報と前記システムプレリサーチモデル情報とをマッチングし、特徴項を取得し、

20

前記特徴項のコーパス情報の得点を計算し、

コーパス情報の得点に基づいて、前記特徴項に対応する遮断待ちテキスト情報が有害情報であるかどうかを判断し、

判断結果に基づいて、第一マッチング結果を取得することを特徴とする請求項 5 に記載の方法。

【請求項 7】

前記前処理された遮断待ちテキスト情報と前記ユーザーフィードバックモデル情報とに対して特徴情報マッチングし、第二マッチング結果を取得するステップにおいては、

前記前処理された遮断待ちテキスト情報及び前記ユーザーフィードバックモデル情報を取得し、

30

前記前処理された遮断待ちテキスト情報と前記ユーザーフィードバックモデル情報とをマッチングし、特徴項を取得し、

前記特徴項のコーパス情報の得点を計算し、

コーパス情報の得点に基づいて、前記特徴項に対応する遮断待ちテキスト情報が有害情報であるかどうかを判断し、

判断結果に基づいて、第二マッチング結果を取得することを特徴とする請求項 6 に記載の方法。

【請求項 8】

前記システムプレリサーチモデル情報は規則索引データベースとシステムリサーチモデルの特徴項情報とを含んでおり、

40

前記ユーザーフィードバックモデル情報は規則索引データベースとユーザーフィードバックモデルの特徴項情報とを含むことを特徴とする請求項 6 または 7 に記載の方法。

【請求項 9】

前記システムプレリサーチモデル情報の規則索引データベースは、システムプリセット規則を含んでおり、

前記ユーザーフィードバックモデル情報の規則索引データベースは、ユーザー設定規則を含むことを特徴とする請求項 8 に記載のインターネットにおける有害情報の遮断方法。

【請求項 10】

インターネットにおける有害情報の遮断装置であって

遮断待ちテキスト情報、システムプレリサーチモデル情報及びユーザーフィードバック

50

モデル情報を取得する情報取得モジュールと、

前記遮断待ちテキスト情報を前処理する前処理モジュールと、

前記前処理された遮断待ちテキスト情報と前記システムプレリサーチモデル情報とに対して特徴情報マッチングし、第一マッチング結果を取得する第一マッチングモジュールと、

前記前処理された遮断待ちテキスト情報と前記ユーザーフィードバックモデル情報とに対して特徴情報マッチングし、前記第一マッチング結果から独立した第二マッチング結果を取得する第二マッチングモジュールと、

第一マッチング結果と第二マッチング結果とが一致しているか否かに基いて、前記遮断待ちテキスト情報を遮断する遮断モジュールと、

を備えることを特徴とする装置。

10

【請求項 1 1】

前記情報取得モジュールは、更に、前記ユーザーフィードバックモデル情報のコーパスを取得することを特徴とする請求項 1 0 に記載の装置。

【請求項 1 2】

前記ユーザーフィードバックモデル情報のコーパスには、ユーザーフィードバックコーパス及び/または被遮断コーパスが含まれることを特徴とする請求項 1 1 に記載の装置。

【請求項 1 3】

更に、

前記ユーザーフィードバックモデル情報のコーパス量及びそれに対応する閾値を取得する閾値取得モジュールと、

前記ユーザーフィードバックモデル情報のコーパス量及びそれに対応する閾値に基づいて、前記ユーザーフィードバックモデル情報を更新する更新モジュールと、

を備えることを特徴とする請求項 1 2 に記載の装置。

20

【請求項 1 4】

前記前処理モジュールは、

前記遮断待ちテキスト情報に対してセグメント処理をするセグメントサブモジュールと、

前記セグメント処理された特徴項候補量を統計する統計サブモジュールと、

を備えることを特徴とする請求項 1 1、1 2 または 1 3 に記載の装置。

30

【請求項 1 5】

前記第一マッチングモジュールは、

前記前処理された遮断待ちテキスト情報及び前記システムプレリサーチモデル情報を取得する情報取得サブモジュールと、

前記前処理された遮断待ちテキスト情報と前記システムプレリサーチモデル情報とをマッチングし、特徴項を取得するマッチングサブモジュールと、

前記特徴項のコーパス情報の得点を計算する統計サブモジュールと、

コーパス情報の得点に基づいて、前記特徴項に対応する遮断待ちテキスト情報が有害情報であるかどうかを判断する判断サブモジュールと、

判断結果に基いて、第一マッチング結果を取得する結果出力サブモジュールと、

を備えることを特徴とする請求項 1 4 に記載の装置。

40

【請求項 1 6】

前記第二マッチングモジュールは、

前記前処理された遮断待ちテキスト情報及び前記システムプレリサーチモデル情報を取得する情報取得サブモジュールと、

前記前処理された遮断待ちテキスト情報と前記ユーザーフィードバックモデル情報とをマッチングし、特徴項を取得するマッチングサブモジュールと、

前記特徴項のコーパス情報の得点を統計する統計サブモジュールと、

コーパス情報の得点に基づいて、前記特徴項に対応する遮断待ちテキスト情報が有害情報であるかどうかを判断する判断サブモジュールと、

50

判断結果に基づいて、第二マッチング結果を取得する結果出力サブモジュールと、
を備えることを特徴とする請求項15に記載の装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、コンピューター情報処理及び情報遮断の技術に関し、特に、統計と規則に基づくインターネットにおける有害情報の遮断方法と装置に関する。

【背景技術】

【0002】

インターネットが迅速に発展するにつれて、情報を伝播するスピードも速くなる。インターネットにおいて、様々なコンテンツが混在しているため、広告、色情、暴力などの有害情報を禁止することが難しい。そして、このような有害情報はますますもっと隠蔽の形で拡散されているため、有害情報の拡散を抑制し、インターネット空間を浄化することが非常に重要である。インターネットにおける膨大なデータ情報は、人工的にはインターネットにおける有害情報を遮断する場合、極めて大量の労力と財力を必要とする。そのため、近年、インターネットにおける有害情報を自動的に遮断する技術の研究が注目されている。

10

【0003】

現在、一般的には、インターネットにおける有害情報を自動的に遮断する技術としては、下記の二つの方法が取り上げられる。

20

(1) キーワードマッチングに基づく遮断方法

判定プロセスにおいて、この方法は精確なマッチング法でキーワードがあるテキストを遮断する。当該方法が採用される場合に、インターネットにおける有害情報を速く遮断でき、簡単で使いやすい。

(2) 統計のテキスト分類モデルに基づく遮断方法

この方法において、本質的には、統計に基づく有害テキストの遮断モデルはテキストを二種類に分類する。テキスト分類は自然言語の処理領域における重要な研究方向であり、大量のモデルが参考にできる。理論上、統計のテキスト分類モデルは、効果的であるはずであるが、実際の適用時には性能が望ましくない。誤判断の場合がよくあり、主な原因が下記で示され、

30

(1) 順方向コーパス(corpus)と逆方向コーパスはバランスが取れていない。その中、順方向コーパスに少量の種類しか含まれていなく、例えば、広告、色情、暴力など、ユーザーが関心を持っている有害情報はメインである。一方、逆方向コーパスには、大量の種類が含まれており、テキスト内容によって分類すると、経済、体育、政治、医薬、アート、歴史、政治、文化、環境、交通、コンピューター、教育、軍事などが分けられている。

(2) 有害情報の内容の表現は非常に隠蔽で変わりやすい。伝播者は常に通常の言語をわざと避け、代わりに、同音字、漢字分解、略字、造語などが使用されている。

(3) ユーザー辞書にキーワードを精確にマッチングする方法しか提供されないため、判定方法は機械的で融通性がなくなる。しかも、単一のキーワードの単語感情極性は代表的なものではなく、誤判断率が高い。例えば、「免費(無料)」と「發票(インボイス)」が同時にコンテキストに現れる場合は、単一の「發票(インボイス)」より説得力がある。

40

(4) 従来中国語情報処理方法はテキスト分類に基づく有害情報の遮断には適用できない。例えば、一定規模の禁止用語の使用や、特徴項に二文字以上の語彙しか含まれないなど。

(5) 広告、色情、暴力などの有害情報を総合的に遮断するための統一モデルがない。

【0004】

上記したインターネットにおける有害情報の自動的遮断を実現するプロセスにおいて、従来の技術では、現在のインターネットからの要請を満足できなく、そして、自動的な更

50

新も実現できない。

【発明の概要】

【0005】

本発明は、インターネットにおける有害情報の遮断方法と装置を提供することを目的とする。

本発明は、このような目的を達成するために、インターネットにおける有害情報の遮断方法であって、遮断待ちテキスト情報、システムプレリサーチ（pre-research）モデル情報及びユーザーフィードバックモデル情報を取得するステップと、前記遮断待ちテキスト情報を前処理するステップと、前記前処理された遮断待ちテキスト情報と前記システムプレリサーチモデル情報とに対して特徴情報マッチングし、第一マッチング結果を取得するステップと、前記前処理された遮断待ちテキスト情報と前記ユーザーフィードバックモデル情報とに対して特徴情報マッチングし、第二マッチング結果を取得するステップと、第一マッチング結果と第二マッチング結果に基いて、前記遮断待ちテキスト情報を遮断するステップと、を備えることを特徴とする。

10

【0006】

また、本発明は、このような目的を達成するために、インターネットにおける有害情報の遮断装置であって、遮断待ちテキスト情報、システムプレリサーチモデル情報及びユーザーフィードバックモデル情報を取得する情報取得モジュールと、前記遮断待ちテキスト情報を前処理する前処理モジュールと、前記前処理された遮断待ちテキスト情報と前記システムプレリサーチモデル情報とに対して特徴情報マッチングし、第一マッチング結果を取得する第一マッチングモジュールと、前記前処理された遮断待ちテキスト情報と前記ユーザーフィードバックモデル情報とに対して特徴情報マッチングし、第二マッチング結果を取得する第二マッチングモジュールと、第一マッチング結果と第二マッチング結果に基いて、前記遮断待ちテキスト情報を遮断する遮断モジュールと、を備えることを特徴とする。

20

【0007】

以上のように、本発明は、遮断待ちテキスト情報、システムプレリサーチモデル情報及びユーザーフィードバックモデル情報を取得し、前記遮断待ちテキスト情報を前処理し、前記前処理された遮断待ちテキスト情報と前記システムプレリサーチモデル情報とに対して特徴情報マッチングし、第一マッチング結果を取得し、前記前処理された遮断待ちテキスト情報と前記ユーザーフィードバックモデル情報とに対して特徴情報マッチングし、第二マッチング結果を取得し、第一マッチング結果と第二マッチング結果に基いて、前記遮断待ちテキスト情報を遮断する。二回のマッチングによって遮断処理が行われるため、システムにおいて有害情報を自動的に遮断する正確性が高く、システムの性能を向上することができる。また、本発明はユーザーフィードバックモデル情報を利用して有害情報を遮断するため、ユーザーフィードバック情報を適時に有害情報の自動的遮断プロセスに適用することが出来、システムモデル情報の自動的更新機能を実現できる。

30

【図面の簡単な説明】

【0008】

【図1】本発明の実施例で、インターネットにおける有害情報の遮断方法を示すフローチャートである。

40

【図2】本発明の更なる実施例で、インターネットにおける有害情報の遮断方法を示すフローチャートである。

【図3】本発明の実施例で、インターネットにおける有害情報の遮断装置の構成を示す模式図である。

【図4】本発明の更なる実施例で、インターネットにおける有害情報の遮断装置の構成を示す模式図である。

【発明を実施するための形態】

【0009】

以下、図面を参照しながら、実施例を使って本発明に係るインターネットにおける有害

50

情報の遮断方法と装置を詳細に説明する。

【0010】

図1で示されるように、本発明の一実施例はインターネットにおける有害情報の遮断方法であって、遮断待ちテキスト情報、システムプレリサーチモデル情報及びユーザーフィードバックモデル情報を取得するステップ101と、前記遮断待ちテキスト情報を前処理するステップ102と、前記前処理された遮断待ちテキスト情報と前記システムプレリサーチモデル情報とに対して特徴情報マッチングし、第一マッチング結果を取得するステップ103と、前記前処理された遮断待ちテキスト情報と前記ユーザーフィードバックモデル情報とに対して特徴情報マッチングし、第二マッチング結果を取得するステップ104と、第一マッチング結果と第二マッチング結果に基づいて、前記遮断待ちテキスト情報を遮断するステップ105と、を備える。

10

【0011】

図2で示されるように、本発明の更なる実施例はインターネットにおける有害情報の遮断方法であって、下記の各ステップを備える。すなわち、

ステップ201：前記システムプレリサーチモデル情報のコーパス及びユーザーフィードバックモデル情報のコーパスを取得する。ここで、前記ユーザーフィードバックモデル情報のコーパスには、ユーザーフィードバックコーパス及び/または被遮断コーパスが含まれる。通常、前記システムプレリサーチモデル情報及びユーザーフィードバックモデル情報の学習コーパスには、順方向コーパスと逆方向コーパスとが含まれる。順方向コーパスとしては、例えば、広告、色情、暴力などの有害情報を含むテキストが10000件用意される。一方、逆方向コーパスとしては、例えば、経済、政治、体育、文化、医薬、交通、環境、軍事、アート、歴史、コンピューター、教育、法律、不動産、科学技術、自動車、人材、娯楽などの非有害情報を含むテキストが30000件用意される。

20

ここで、前記学習コーパスの収集において、順方向コーパスと逆方向コーパスはバランスが取れていない場合がよくあり、一方は範囲が広すぎるが、もう一方は範囲が狭すぎる。本発明において、このようなバランスが取れていないコーパスの分布が許容される。コーパス範囲が広い場合は、量ではなく、できるだけ多くの種類を確保しながら用意する。

【0012】

ステップ202：遮断待ちテキスト情報、システムプレリサーチモデル情報及びユーザーフィードバックモデル情報を取得する。

30

【0013】

ステップ203：前記遮断待ちテキスト情報を前処理する。

このステップにおいては、前記遮断待ちテキスト情報に対してセグメント処理をする。例えば、句読点と常用語に基づいて、コーパスを区切る。常用語とは、よく使用され判定には無意味な語彙であり、例えば「的」、「了」など。しかし、「

您

(貴方)」はよく順方向コーパスに、「我(私)」はよく逆方向コーパスに使用されるが、いずれも常用語に使用されない。

ここで、自然言語処理においては、よく用いられる禁止用語リストが常用語リストとして適用されない。通常、「方正智思分詞4.0(ペキンファンダー社が開発ソフトウェア)」によって、コーパスに対してセグメントや品詞分類をすることができる。前記セグメント処理されたセグメントユニットは後工程における最小の処理単位である。

40

前記セグメント処理された特徴項候補量を統計する。例えば、前記セグメント処理されたセグメントユニットにおける非漢字部分を統計し、前記セグメントユニットの合計をN1、非漢字部分の合計をN2とする場合、 $N2/N1$ が閾値より大きければ、当該特徴項候補に対応する遮断待ちテキスト情報は有害情報と判断される。判断の理由としては、大量のノイズ文字がこの情報に含まれ、広告などのスパムテキストであるかもしれない。もしくは、前記セグメントユニットにおける、広告によく用いられるURL、電話番号、電子メールアドレス、QQ等の連絡方法の数量num(ad)を統計し、デフォルトウェイトsco

50

readを与える。

【0014】

ステップ204：前記前処理された遮断待ちテキスト情報と前記システムプレリサーチモデル情報とに対して特徴情報マッチングし、第一マッチング結果を取得する。このステップにおいては、

ステップ2041：前記前処理された遮断待ちテキスト情報及び前記システムプレリサーチモデル情報を取得する。前記システムプレリサーチモデル情報に規則索引データベースと前記システムリサーチモデルの特徴項情報とが含まれる。具体的には、前記規則索引データベースにおけるユーザー規則索引データベースとユーザーキーワード索引データベースが生成されるプロセスは以下のものである。すなわち、

ステップS1：キーワード解析。まず、常用漢字のピンインの索引を作成し、キーワードにおける各字のピンインの索引に基づいてキーワード全体の索引を生成する。それから、キーワードにおける各字に対して構造的に分解し、分解された結果に基づいて、キーワードを再帰し再組合せする。最後、キーワードの索引と、分解の集合によってキーバリュペア (key value pair) を形成させ、全ての解析結果を保存し、ユーザーキーワードの索引データベースを生成する。例えば「法輪功」は、キーワード解析後に、一つの索引値が生成され、しかも幾つかの分解結果がある。具体的には、「三去車侖工力」、「法車侖功」などが含まれる。

【0015】

ステップS2：文法解析。コンピューターによって規則文法を、処理できる形に解析する。前記規則文法には、ANDと、ORと、NEARと、NOTとが含まれる。例えば、「A AND B」の場合、AとBは解析待ちのキーワードであり、AND文法とはAとBが同時にコンテキストに現れる場合に、当該規則はマッチングに成功である。キーワードと規則文法に対してキーバリュペアを形成し、全ての解析結果を保存しユーザー規則索引データベースを生成する。

【0016】

ここで、上記した規則索引データベースにおいては、ユーザーが設定した規則でも良いし、システムのプリセット規則でもよい。上記したステップはユーザー設定規則を解析し相応する索引データベースを生成するプロセスであり、当該索引データベースは以下のマッチングプロセスを最適化できる。

【0017】

ステップ2042：前記前処理された遮断待ちテキスト情報と前記システムプレリサーチモデル情報とをマッチングし、特徴項を取得する。ここで、前記システムプレリサーチモデル情報には、規則索引データベースと前記システムプレリサーチモデル特徴項の情報とが含まれ、具体的には、システムプレリサーチモデル特徴項の情報を取得するプロセスは、以下のものである。すなわち、

ステップS1：前記セグメントユニットを文字列に組合せし、特徴項候補とする。

(例1)：連続的なセグメントユニットを文字列に組合せする場合。

各文のセグメントユニットに対して、一番目のセグメントユニットから、組合せウィンドウの最大値をNとして組合せする。順序があるセグメントユニット「A B C D」を例として挙げると、組合せウィンドウの最大値が3である場合に、文字列の組合せはA B C、B C D、A B、B C、C D、A、B、C、Dとの九つがある。

(例2)：

非連続のセグメントユニットを文字列に組合せする場合。

例1で組合せの文字列に対してピンインの索引を計算し、前記ステップ2041におけるステップS1で生成されたユーザーキーワードの索引データベースに基づいてマッチングする。マッチング成功の集合があれば、マッチング成功の数量num (user) を統計する。それから、前記ステップ2041におけるステップS2で生成されたユーザー規則索引データベースに基づいてマッチングし、マッチング成功すれば、非連続のセグメントユニットに対して一つの文字列が生成される。例えば、例1における九つの文字列。ユーザー

10

20

30

40

50

キーワードの索引データベースにおいて、二つの文字列 A、D がマッチング成功する。ユーザー規則索引データベースに規則「A NEAR 2 D」がある場合に、特徴項 AD が新たに生成される。ここに、2 は A と D の距離は 2 以下の意味とする。マッチング成功の数量 $\text{num}(\text{user})$ を累計し、デフォルトウェイト $\text{score}_{\text{user}}$ を与える。

【0018】

ステップ S 2 : 前記特徴項候補を頻度によって遮断する。具体的には、学習コーパスに特徴項候補が現れる回数を統計し、頻度に従って遮断し、頻度が閾値以上の特徴項候補を残しておき、頻度が閾値未満の特徴項候補を削除し、閾値を調整することによって、残す範囲を制御する。

【0019】

ステップ S 3 : 前記特徴項候補を頻度によって再遮断する。具体的な遮断プロセスは、まず、改めて不適切の頻度を評価し、例えば、全ての B が現れる時に、A も同時に現れ、AB になる場合であれば、B の頻度が 0 になる。頻度再評価式は：

$$\begin{cases} \log_2 |a| * f(a) & a \text{ が含まれていない場合} \\ \log_2 |a| * \left(f(a) - \frac{\sum_{b \in T_a} f(b)}{P(T_a)} \right) & \text{その他} \end{cases}$$

ここで、a は特徴項であり、 $f(a)$ は a のワード頻度であり、b は a が含まれる長い文字列の特徴項であり、 T_a は b の集合であり、 $P(T_a)$ は集合のサイズである。

【0020】

それから、再評価された頻度に従って再遮断を行い、頻度が閾値未満の特徴項候補を削除し、閾値を調整することによって、残す範囲を制御する。

【0021】

ステップ S 4 : 前記特徴項候補が自動的に選択されて、特徴項が抽出される。具体的には、当該ステップにおいて、前記ステップ 3 で順方向コーパスから取得される特徴項候補と前記ステップ 3 で逆方向コーパスから取得される特徴項候補とを組合せ、組合せによる特徴項候補は二つのワード頻度があり、それぞれ順方向頻度と逆方向頻度に対応する。統計学上のカイ 2 乗統計量によって特徴項を自動的に選択し、カイ 2 乗値が最大である前からの N 個の特徴項候補を残して最終の特徴項情報として、 χ^2 統計量の式は：

$$\chi^2(\omega_i, C_k) = \frac{N(AD - BC)^2}{(A + C)(A + B)(B + D)(C + D)}$$

【0022】

その中、A、B、C、D の意味はそれぞれ下記で示され、

	C_k に属する テキスト	C_k に属しない テキスト	合計
w 特徴項 が含まれる	A	B	A + B
w 特徴項 が含まれない	C	D	C + D
合計	A + C	B + D	N

10

20

30

40

50

表における k は「0」または「1」で、順方向タイプと逆方向タイプの二タイプを代表する。

【0023】

ここで、前記特徴項は一文字単語（単一の文字からなる単語）と複数文字単語（複数の文字からなる単語）とを含む。一文字単語は逆方向テキストの判定に影響が大きい。特に、フォーラムテキスト情報の内容において、一文字単語に基づくセグメントユニットがよく用いられ、一文字単語を考えなければ、逆方向テキストが誤判断しやすくなる。

【0024】

ステップ2043：前記特徴項のコーパス情報の得点を計算する。ステップS4で前記特徴項の頻度が既に保存され、特徴項ごとに順方向頻度と逆方向頻度をそれぞれ代表する二つの頻度を有する。例えば、「発票（インボイス）」の順方向頻度は逆方向頻度よりずっと大きく、「発票（インボイス）」は広告の有害情報によく用いられるからである。各特徴項の順方向頻度を特徴項の順方向ウェイトとして、各特徴項の逆方向頻度を特徴項の逆方向ウェイトとする。全ての特徴項の順方向/逆方向ウェイトに対して正規化を行い、これによってこそ、ウェイト値は比較する意味がある。正規化の式は：

$$score(\omega_i) = \frac{freq(\omega_i)}{\sum freq(\omega_i)}$$

【0025】

生成された特徴項とそのウェイトがシステムより準備しておく標準的二種類のコーパスに基づいて学習することによって取得されるため、生成された結果を保存しシステムプレリサーチモデル特徴項情報とする。

【0026】

前記前処理された遮断待ちテキスト情報と前記システムプレリサーチモデル特徴項情報とに対して特徴情報マッチングし、遮断待ちテキスト特徴項情報を取得し、前記特徴項情報の順方向得点を計算し、その計算式は：

$$score_{pos}(doc) = \sum \log(score(\omega_i)_{pos})$$

前記特徴項情報の逆方向得点を計算し、その計算式は：

$$score_{neg}(doc) = \sum \log(score(\omega_i)_{neg})$$

なお、 $num(ad)$ と $num(user)$ も考慮すると、上記計算式の右側が下記のようになる：

$$\sum \log(score(\omega_i)_{neg}) + num(ad) * score_{ad} + num(user) * score_{user}$$

【0027】

ステップ2044：コーパス情報の得点に基づいて、前記特徴項に対応する遮断待ちテキスト情報が有害情報であるかどうかを判断する。 $score_{pos}(doc) > score_{neg}(doc)$ の場合に、システムプレリサーチモデルはこの遮断待ちテキスト情報が有害テキストと判断する。また、 $score_{pos}(doc) == score_{neg}(doc)$ の場合に、このモデルが無効となり、判定が無効される。また、 $score_{pos}(doc) < score_{neg}(doc)$ の場合に、システムプレリサーチモデルはこの遮断待ちテキスト情報が通常テキストと判断する。

【0028】

ステップ2045：判断結果に基づいて、第一マッチング結果を取得する。

【0029】

ステップ205：前記前処理された遮断待ちテキスト情報と前記ユーザーフィードバックモデル情報とに対して特徴情報マッチングし、第二マッチング結果を取得する。当該ス

10

20

30

40

50

トップに含まれるプロセスとステップ204におけるプロセスが大体同じである。

【0030】

ここで、前記ユーザーフィードバックモデル情報を取得するプロセスとシステムプレリサーチモデル情報を取得するプロセスについて、主な相違点はステップ201で学習コーパスの選択である。前記ユーザーフィードバックモデル情報の学習コーパスが下記の二つの方面から取得できる：

(1) ユーザーフィードバックメカニズム

実際の使用するプロセスにおいて、判定には問題があると発見され、主に有害情報が通常情報と誤判断される場合に、ユーザーはシステムにエラーを報告し、システムはユーザーからの標準回答を受けフィードバックコーパスとする。

(2) 判断モデルメカニズム

処理待ちのテキストがステップ206での有害情報判定を受け、当該テキストの判定結果が出力される。結果は有害情報テキストが通常テキストである。信頼性を判定する状況に基づいて、処理待ちのテキストはフィードバック学習に用いられるかどうかを判断する。

【0031】

ステップ206：前記第一マッチング結果と第二マッチング結果に基づいて、前記遮断待ちテキスト情報に対して遮断処理を行う。具体的には、前記第一マッチング結果と第二マッチング結果、つまり、システムプレリサーチモデル情報とユーザーフィードバックモデル情報が一致するかどうかを判断する。一致と判定され、いずれも有害情報テキストまたは通常情報テキストであれば、判定結果の信頼性が高く、フィードバック学習に使用できる。一方、不一致と判定されれば、判定結果の信頼性がより低くなる。比較的厳しい遮断方針が採用される場合に、このテキストが遮断されるが、フィードバック学習に使用できない。その中には一つのモデルが無効である場合に、結果は残りのモデルの判定結果次第であり、ある程度の信頼性があり、フィードバック学習に使用できる。二つのモデルが両方とも無効である場合に、無効標識をリターンさせ、フィードバック学習に使用できない。

【0032】

ここで、前記方法において、遮断待ちテキスト情報の判定プロセスの完成後に、

前記ユーザーフィードバックモデル情報のコーパス量及びそれに対応する閾値を取得する。具体的には、フィードバック学習に使用できるコーパス量を統計し、前記コーパス量はそれに対応する閾値を越えているかどうかを判断する。

【0033】

前記ユーザーフィードバックモデル情報のコーパス量及びそれに対応する閾値に基づいて、前記ユーザーフィードバックモデル情報を更新する。コーパス量は閾値より大きい場合に、改めてフィードバックコーパスを学習し、ユーザーフィードバックモデル情報を更新する。閾値を調整することによって、更新の周期が調整される。

【0034】

図3で示されるように、本発明の一実施例のインターネットにおける有害情報の遮断装置は、

遮断待ちテキスト情報、システムプレリサーチモデル情報及びユーザーフィードバックモデル情報を取得する情報取得モジュール301と、前記遮断待ちテキスト情報を前処理する前処理モジュール302と、前記前処理された遮断待ちテキスト情報と前記システムプレリサーチモデル情報とに対して特徴情報マッチングし、第一マッチング結果を取得する第一マッチングモジュール303と、前記前処理された遮断待ちテキスト情報と前記ユーザーフィードバックモデル情報とに対して特徴情報マッチングし、第二マッチング結果を取得する第二マッチングモジュール304と、第一マッチング結果と第二マッチング結果に基づいて、前記遮断待ちテキスト情報を遮断する遮断モジュール305と、を備える。

【0035】

図4で示されるように、本発明の一実施例のインターネットにおける有害情報の遮断装

10

20

30

40

50

置は、下記のモジュールを備える。すなわち、

情報取得モジュール401：遮断待ちテキスト情報、システムプレリサーチモデル情報及びユーザーフィードバックモデル情報を取得し、ユーザーフィードバックモデル情報の学習コーパスを取得する。ここで、前記ユーザーフィードバックモデル情報のコーパスには、ユーザーフィードバックコーパス及び/または被遮断コーパスが含まれる。

【0036】

前処理モジュール402：前記遮断待ちテキスト情報を前処理する。このモジュールは、前記遮断待ちテキスト情報に対してセグメント化処理をするセグメントサブモジュール4021と、前記セグメント処理された特徴項候補量を統計する統計サブモジュール4022と、を備える。

10

【0037】

第一マッチングモジュール403：前記前処理された遮断待ちテキスト情報と前記システムプレリサーチモデル情報とに対して特徴情報マッチングし、第一マッチング結果を取得する。このモジュールは、前記前処理された遮断待ちテキスト情報及び前記システムプレリサーチモデル情報を取得する情報取得サブモジュール4031と、前記前処理された遮断待ちテキスト情報と前記システムプレリサーチモデル情報とをマッチングし、特徴項を取得するマッチングサブモジュール4032と、前記特徴項のコーパス情報の得点を計算する統計サブモジュール4033と、コーパス情報の得点に基づいて、前記特徴項に対応する遮断待ちテキスト情報が有害情報であるかどうかを判断する判断サブモジュール4034と、判断結果に基づいて、第一マッチング結果を取得する結果出力サブモジュール4035と、を備える。

20

【0038】

第二マッチングモジュール404：前記前処理された遮断待ちテキスト情報と前記ユーザーフィードバックモデル情報とに対して特徴情報マッチングし、第二マッチング結果を取得する。このモジュールは、前記前処理された遮断待ちテキスト情報及び前記システムプレリサーチモデル情報を取得する情報取得サブモジュール4041と、前記前処理された遮断待ちテキスト情報と前記ユーザーフィードバックモデル情報とをマッチングし、特徴項を取得するマッチングサブモジュール4042と、前記特徴項のコーパス情報の得点を統計する統計サブモジュール4043と、コーパス情報の得点に基づいて、前記特徴項に対応する遮断待ちテキスト情報が有害情報であるかどうかを判断する判断サブモジュール4044と、判断結果に基づいて、第二マッチング結果を取得する結果出力サブモジュール4045と、を備える。

30

【0039】

遮断モジュール405：前記第一マッチング結果と第二マッチング結果に基づいて、前記遮断待ちテキスト情報に対して遮断処理を行う。

【0040】

閾値取得モジュール406：前記ユーザーフィードバックモデル情報のコーパス量及びそれに対応する閾値を取得する。

【0041】

更新モジュール407：前記ユーザーフィードバックモデル情報のコーパス量及びそれに対応する閾値に基づいて、前記ユーザーフィードバックモデル情報を更新する。前記閾値取得モジュールが取得したユーザーフィードバックモデル情報のコーパス量はそれに対応する閾値に達する場合に、前記更新モジュールが前記ユーザーフィードバックモデル情報のコーパス量及びそれに対応する閾値に基づいて、前記ユーザーフィードバックモデル情報を更新する。

40

【0042】

以上のように、本発明の実施例に提供されるインターネットにおける有害情報の遮断方法と装置は、遮断待ちテキスト情報、システムプレリサーチモデル情報及びユーザーフィードバックモデル情報を取得し、前記遮断待ちテキスト情報を前処理し、前記前処理された遮断待ちテキスト情報と前記システムプレリサーチモデル情報とに対して特徴情報マッ

50

チングし、第一マッチング結果を取得し、前記前処理された遮断待ちテキスト情報と前記ユーザーフィードバックモデル情報とに対して特徴情報マッチングし、第二マッチング結果を取得し、第一マッチング結果と第二マッチング結果に基づいて、前記遮断待ちテキスト情報を遮断する。二回のマッチングによって遮断処理が行われるため、システムにおいて有害情報を自動的に遮断する正確性が高く、システムの性能を向上することができる。また、本発明はユーザーフィードバックモデル情報を利用して有害情報を遮断するため、ユーザーフィードバック情報を適時に有害情報の自動的遮断プロセスに適用することが出来、システムモデル情報の自動的更新機能を実現できる。

【 0 0 4 3 】

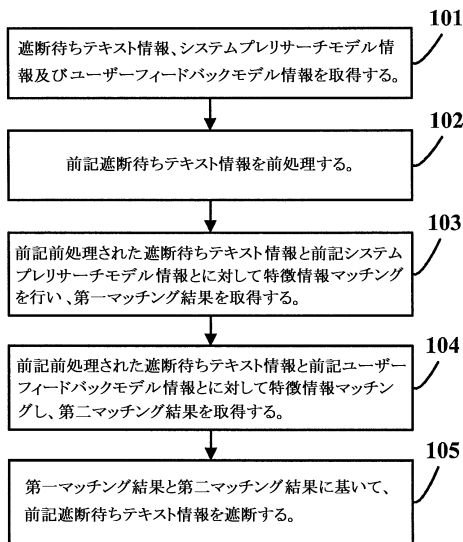
上記した説明から、当業者であれば分かるように、前記した実施例における全部または一部のステップは、プログラムによって関連するハードウェアで実行することができる。前記プログラムは例えば、ROM/RAM、磁気ディスク、光ディスクなどの記憶装置に記憶されてもよい。

【 0 0 4 4 】

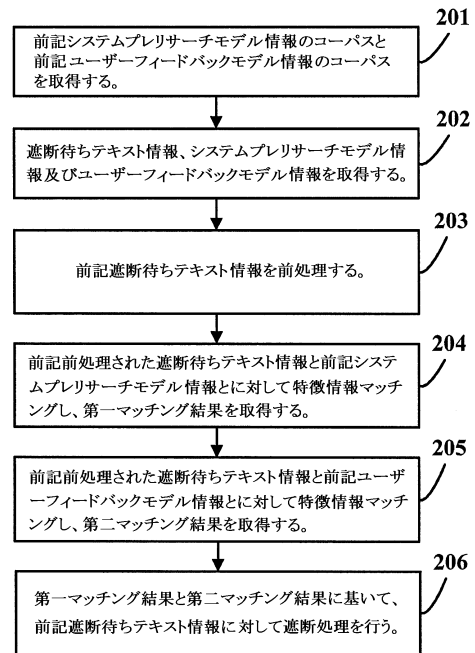
以上で説明した内容はただ本発明の各実施形態であり、本発明が保護しようとする範囲はここに限定されるものではなく、いかなる当業者は本発明より開示された技術範囲で容易に想到できる適宜組み合わせで得られる実施形態についても本発明の技術的範囲に含まれる。

10

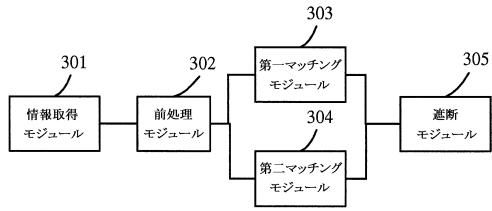
【 図 1 】



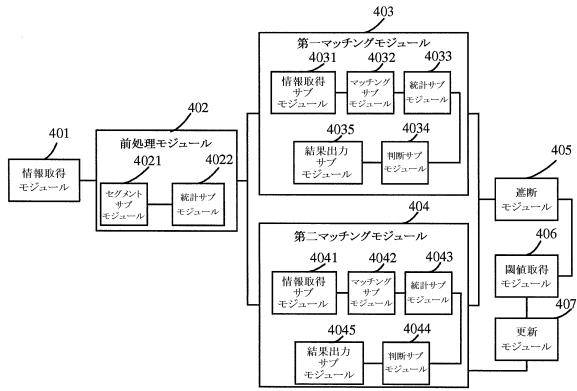
【 図 2 】



【図3】



【図4】



フロントページの続き

(73)特許権者 507232478

北京大学

PEKING UNIVERSITY

中華人民共和国北京市 海 淀区 頤 和 園 路5号

No.5, Yiheyuan Road, Haidian District, Beijing 100871, China

(73)特許権者 507232456

北京北大方正 電 子有限公司

BEIJING FOUNDER ELECTRONICS CO., LTD.

中華人民共和国北京市 海 淀区上地五街9号方正大厦

Founder Building, No.9, Shangdiwu Street, Haidian District, Beijing 100085, China

(73)特許権者 513157039

北京北大方正技 術 研究院有限公司

PEKING UNIVERSITY FOUNDER R & D CENTER

中華人民共和国北京市 海 淀区成府路298号中 関 村方正大厦4 層

4 Floor, Zhongguancun Founder Building, No.298, Chengfu Road, Haidian District, Beijing 100871, China

(74)代理人 110000729

特許業務法人 ユニアス国際特許事務所

(72)発明者 チェン、イェン

中華人民共和国 100085 ベイジン、北京市 海 淀区上地五街九号方正大厦

(72)発明者 ユー、シャオミン

中華人民共和国 100085 ベイジン、北京市 海 淀区上地五街九号方正大厦

(72)発明者 ヤン、チエンウー

中華人民共和国 100085 ベイジン、北京市 海 淀区上地五街九号方正大厦

審査官 木村 雅也

(56)参考文献 特開2005-235206(JP,A)

特開2007-179540(JP,A)

特表2010-507153(JP,A)

特開2009-140437(JP,A)

米国特許第05987457(US,A)

中国特許出願公開第101477544(CN,A)

中国特許出願公開第101639824(CN,A)

中国特許出願公開第101702167(CN,A)

中国特許出願公開第101877704(CN,A)

米国特許第05867799(US,A)

(58)調査した分野(Int.Cl., DB名)

G06F 13/00

G06F 17/27