

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号  
特許第6092141号  
(P6092141)

(45) 発行日 平成29年3月8日(2017.3.8)

(24) 登録日 平成29年2月17日(2017.2.17)

(51) Int.Cl.

F I

G O 6 N 99/00 (2010.01)

G O 6 N 99/00 1 5 3

請求項の数 4 (全 14 頁)

(21) 出願番号	特願2014-46601 (P2014-46601)	(73) 特許権者	000004226
(22) 出願日	平成26年3月10日 (2014.3.10)		日本電信電話株式会社
(65) 公開番号	特開2015-170281 (P2015-170281A)		東京都千代田区大手町一丁目5番1号
(43) 公開日	平成27年9月28日 (2015.9.28)	(74) 代理人	110001519
審査請求日	平成28年2月15日 (2016.2.15)		特許業務法人太陽国際特許事務所
		(72) 発明者	貞光 九月
			東京都千代田区大手町一丁目5番1号 日
			本電信電話株式会社内
		(72) 発明者	松尾 義博
			東京都千代田区大手町一丁目5番1号 日
			本電信電話株式会社内
		(72) 発明者	浅野 久子
			東京都千代田区大手町一丁目5番1号 日
			本電信電話株式会社内

最終頁に続く

(54) 【発明の名称】 データ解析装置、方法、及びプログラム

(57) 【特許請求の範囲】

【請求項 1】

正例及び負例の何れか一方を示すラベルが付与されている学習データに基づいて、対象データが正例であるか又は負例であるかを識別するための識別モデルを学習し、学習された識別モデルに基づいて、前記ラベルが付与されていない未知データに前記ラベルを付与し、前記未知データに前記ラベルが付与された結果に基づいて、前記識別モデルを学習することを繰り返すブートストラップ法に従って、前記識別モデルを繰り返し学習したときに、繰り返し毎に前記未知データに前記ラベルが付与された結果から得られる、前記正例のラベルが付与された未知データの各々からなる正例集合の入力を受け付ける入力部と、

前記入力部によって受け付けた繰り返し毎の前記正例集合のうち、繰り返し開始から予め定められた第1の繰り返し回数までに得られた前記正例集合を擬似正例データとして選択し、予め定められた第2の繰り返し回数から繰り返し終了までに得られた前記正例集合を擬似負例データとして選択する擬似正例負例選択部と、

前記擬似正例負例選択部によって選択された前記擬似正例データ、及び前記擬似負例データから抽出された素性に基づいて、前記識別モデルを学習するモデル学習部と、

識別対象データの各々について、前記モデル学習部によって学習された前記識別モデルと、前記識別対象データから抽出された素性とに基づいて、前記識別対象データが正例である度合いを表すスコアを算出するスコアリング部と、

前記識別対象データの各々について、前記スコアリング部によって算出されたスコアが、予め定められた第1の閾値より大きい場合には、前記識別対象データが正例であると判

10

20

定し、前記スコアリング部によって算出されたスコアが、予め定められた第2の閾値より小さい場合には、前記識別対象データが負例であると判定する判定部と、  
を含むデータ解析装置。

【請求項2】

前記識別対象データを、前記正例集合とし、

前記判定部は、前記ブートストラップ法における各繰り返しにおいて、前記正例のラベルが付与された未知データの数、及び前記負例のラベルが付与された未知データの数の割合に基づいて、前記第2の閾値を変化させ、

前記識別対象データの各々について、前記スコアリング部によって算出されたスコアが、前記第1の閾値より大きい場合には、前記識別対象データが正例であると判定し、前記スコアリング部によって算出されたスコアが、前記変化した第2の閾値より小さい場合には、前記識別対象データが負例であると判定する

請求項1記載のデータ解析装置。

【請求項3】

入力部が、正例及び負例の何れか一方を示すラベルが付与されている学習データに基づいて、対象データが正例であるか又は負例であるかを識別するための識別モデルを学習し、学習された識別モデルに基づいて、前記ラベルが付与されていない未知データに前記ラベルを付与し、前記未知データに前記ラベルが付与された結果に基づいて、前記識別モデルを学習することを繰り返すブートストラップ法に従って、前記識別モデルを繰り返し学習したときに、繰り返し毎に前記未知データに前記ラベルが付与された結果から得られる、前記正例のラベルが付与された未知データの各々からなる正例集合の入力を受け付けるステップと、

擬似正例負例選択部が、前記入力部によって受け付けた繰り返し毎の前記正例集合のうち、繰り返し開始から予め定められた第1の繰り返し回数までに得られた前記正例集合を擬似正例データとして選択し、予め定められた第2の繰り返し回数から繰り返し終了までに得られた前記正例集合を擬似負例データとして選択するステップと、

モデル学習部が、前記擬似正例負例選択部によって選択された前記擬似正例データ、及び前記擬似負例データから抽出された素性に基づいて、前記識別モデルを学習するステップと、

スコアリング部が、識別対象データの各々について、前記モデル学習部によって学習された前記識別モデルと、前記識別対象データから抽出された素性とに基づいて、前記識別対象データが正例である度合いを表すスコアを算出するステップと、

判定部が、前記識別対象データの各々について、前記スコアリング部によって算出されたスコアが、予め定められた第1の閾値より大きい場合には、前記識別対象データが正例であると判定し、前記スコアリング部によって算出されたスコアが、予め定められた第2の閾値より小さい場合には、前記識別対象データが負例であると判定するステップと、

を含むデータ解析方法。

【請求項4】

コンピュータに、請求項1又は請求項2記載のデータ解析装置の各部として機能させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データ解析装置、方法、及びプログラムに係り、特に、識別対象データが正例であるか負例であるかを判定するデータ解析装置、方法、及びプログラムに関する。

【背景技術】

【0002】

一度学習したモデルに基づいて未知事例の識別を行い、識別結果の正例・負例を新たな学習データとして用いていく繰り返し学習の枠組みをブートストラップ法と呼ぶ。

【0003】

例えば、少量の教師信号となる単語または文字列のセットを入力とし、未知の単語または文字列が正例か負例なのかを逐次的に識別していくブートストラップを行う手法がこれまでに数多く提案されている（例えば、非特許文献１）。

【０００４】

また、ブートストラップの過程に着目し、初期のイテレーションで正例と判定された事例と終盤のイテレーションで正例と判定された事例を分け、識別対象となる任意の事例が、これら２つのグループのうち、いずれに近いかを、単純な情報量（分布類似度）に基づいて判定することが提案されている（例えば、非特許文献２）。

【先行技術文献】

【非特許文献】

10

【０００５】

【非特許文献１】Patrick Pantel and Marco Pennacchiotti, “Espresso: Leveraging Generic Patterns for Automatically Harvesting Semantic Relations.” , COLING-ACL , 2006.

【非特許文献２】Tara Mcintosh, “Unsupervised discovery of negative categories in lexicon bootstrapping” , ACL2010 , p.356-365

【発明の概要】

【発明が解決しようとする課題】

【０００６】

単語や文字列に対する自動ラベリング等のタスクにおいて、その判別の基準となる正例のみが与えられる場合は多い。

20

【０００７】

特に、上記非特許文献１に記載の技術のようにブートストラップ法を用いる手法においては、イテレーションが進むほど、正例として識別したい対象とは異なる事例を誤って正例を識別するセマンティックドリフトと呼ばれる現象が頻発する。

【０００８】

また、正例のみが与えられる場合には、負例が存在しないために、どこまでを正例識別対象とするのかといった基準が、システムには自明でなく、そのため自動的な負例判別法が必要とされる。

【０００９】

30

また、上記非特許文献２に記載の技術では、単純な情報量しか用いないため、より高度な素性（例えば単語の連鎖情報や、文書全体のトピック情報等）を考慮できないので、負例検出性能が達成されない。

【００１０】

本発明は、上記の事情を鑑みてなされたもので、負例データを精度よく抽出することができるデータ解析装置を提供することを目的とする。

【課題を解決するための手段】

【００１１】

上記の目的を達成するために本発明に係るデータ解析装置は、正例及び負例の何れか一方を示すラベルが付与されている学習データに基づいて、対象データが正例であるか又は負例であるかを識別するための識別モデルを学習し、学習された識別モデルに基づいて、前記ラベルが付与されていない未知データに前記ラベルを付与し、前記未知データに前記ラベルが付与された結果に基づいて、前記識別モデルを学習することを繰り返すブートストラップ法に従って、前記識別モデルを繰り返し学習したときに、繰り返し毎に前記未知データに前記ラベルが付与された結果から得られる、前記正例のラベルが付与された未知データの各々からなる正例集合の入力を受け付ける入力部と、前記入力部によって受け付けた繰り返し毎の前記正例集合のうち、繰り返し開始から予め定められた第１の繰り返し回数までに得られた前記正例集合を擬似正例データとして選択し、予め定められた第２の繰り返し回数から繰り返し終了までに得られた前記正例集合を擬似負例データとして選択する擬似正例負例選択部と、前記擬似正例負例選択部によって選択された前記擬似正例デ

40

50

ータ、及び前記擬似負例データから抽出された素性に基づいて、前記識別モデルを学習するモデル学習部と、識別対象データの各々について、前記モデル学習部によって学習された前記識別モデルと、前記識別対象データから抽出された素性とに基づいて、前記識別対象データが正例である度合いを表すスコアを算出するスコアリング部と、前記識別対象データの各々について、前記スコアリング部によって算出されたスコアが、予め定められた第1の閾値より大きい場合には、前記識別対象データが正例であると判定し、前記スコアリング部によって算出されたスコアが、予め定められた第2の閾値より小さい場合には、前記識別対象データが負例であると判定する判定部と、を含んで構成されている。

【0012】

本発明に係るデータ解析方法は、入力部が、正例及び負例の何れか一方を示すラベルが付与されている学習データに基づいて、対象データが正例であるか又は負例であるかを識別するための識別モデルを学習し、学習された識別モデルに基づいて、前記ラベルが付与されていない未知データに前記ラベルを付与し、前記未知データに前記ラベルが付与された結果に基づいて、前記識別モデルを学習することを繰り返すブートストラップ法に従って、前記識別モデルを繰り返し学習したときに、繰り返し毎に前記未知データに前記ラベルが付与された結果から得られる、前記正例のラベルが付与された未知データの各々からなる正例集合の入力を受け付けるステップと、擬似正例負例選択部が、前記入力部によって受け付けた繰り返し毎の前記正例集合のうち、繰り返し開始から予め定められた第1の繰り返し回数までに得られた前記正例集合を擬似正例データとして選択し、予め定められた第2の繰り返し回数から繰り返し終了までに得られた前記正例集合を擬似負例データとして選択するステップと、モデル学習部が、前記擬似正例負例選択部によって選択された前記擬似正例データ、及び前記擬似負例データから抽出された素性に基づいて、前記識別モデルを学習するステップと、スコアリング部が、識別対象データの各々について、前記モデル学習部によって学習された前記識別モデルと、前記識別対象データから抽出された素性とに基づいて、前記識別対象データが正例である度合いを表すスコアを算出するステップと、判定部が、前記識別対象データの各々について、前記スコアリング部によって算出されたスコアが、予め定められた第1の閾値より大きい場合には、前記識別対象データが正例であると判定し、前記スコアリング部によって算出されたスコアが、予め定められた第2の閾値より小さい場合には、前記識別対象データが負例であると判定するステップと、を含んで構成されている。

【0013】

本発明は、前記識別対象データを、前記正例集合とし、前記判定部は、前記ブートストラップ法における各繰り返しにおいて、前記正例のラベルが付与された未知データの数、及び前記負例のラベルが付与された未知データの数の割合に基づいて、前記第2の閾値を変化させ、前記識別対象データの各々について、前記スコアリング部によって算出されたスコアが、前記第1の閾値より大きい場合には、前記識別対象データが正例であると判定し、前記スコアリング部によって算出されたスコアが、前記変化した第2の閾値より小さい場合には、前記識別対象データが負例であると判定するようにすることができる。

【0014】

本発明のプログラムは、コンピュータに、上記のデータ解析装置の各部として機能させるためのプログラムである。

【発明の効果】

【0015】

以上説明したように、本発明のデータ解析装置、方法、及びプログラムによれば、繰り返し毎の正例集合のうち、繰り返し開始から予め定められた第1の繰り返し回数までに得られた正例集合を擬似正例データとして選択し、予め定められた第2の繰り返し回数から繰り返し終了までに得られた正例集合を擬似負例データとして選択し、選択された擬似正例データ、及び擬似負例データから抽出された素性に基づいて、識別モデルを学習し、識別対象データの各々について、学習された識別モデルと、識別対象データから抽出された素性とに基づいて、識別対象データが正例である度合いを表すスコアを算出し、スコアが

予め定められた第2の閾値より小さい場合には、識別対象データが負例であると判定することにより、負例データを精度よく抽出することができる、という効果が得られる。

【図面の簡単な説明】

【0016】

【図1】本発明の実施の形態に係るデータ解析装置の構成を示す概略図である。

【図2】本発明の第1の実施の形態に係るデータ解析装置におけるデータ解析処理ルーチンの内容を示すフローチャートである。

【図3】本発明の第2の実施の形態に係るデータ解析装置におけるデータ解析処理ルーチンの内容を示すフローチャートである。

【発明を実施するための形態】

10

【0017】

本発明の実施の形態は、単語や文字列等に対し、ある基準を用いて識別を行う際のモデル学習に必要な負例を、自動的に抽出する技術である。以下、図面を参照して本発明の実施の形態を詳細に説明する。

【0018】

<第1の実施の形態>

<システム構成>

本発明の第1の実施の形態に係るデータ解析装置100は、入力された、ブートストラップ法で得られた正例集合から、新負例データを抽出する。このデータ解析装置100は、CPUと、RAMと、後述するデータ解析処理ルーチンを実行するためのプログラムを記憶したROMとを備えたコンピュータで構成され、機能的には次に示すように構成されている。図1に示すように、データ解析装置100は、入力部10と、演算部20と、出力部30とを備えている。

20

【0019】

本実施の形態におけるブートストラップ法は、まず、正例及び負例の何れか一方を示すラベルが付与されている学習データに基づいて、識別モデルを学習する。そして、学習された識別モデルに基づいて、ラベルが付与されていない未知データにラベルを付与し、未知データにラベルが付与された結果に基づいて、識別モデルを学習する。そして、この未知データに対するラベル付与と、識別モデルの学習とを繰り返す。なお、識別モデルは、対象データが正例であるか又は負例であるかを識別するためのモデルである。

30

【0020】

本実施の形態に係るデータ解析装置100には、ブートストラップ法に従って識別モデルを繰り返し学習したときに、繰り返し毎に、未知データにラベルが付与された結果から得られる、正例のラベルが付与された未知データの各々を含む正例集合が入力される。

【0021】

入力部10は、ブートストラップ法における繰り返し毎に得られた正例集合の入力を受け付ける。また、入力部10は、後述する素性化参照データを受け付ける。

【0022】

例えば、対象データが単語であって、単語が車名であるかどうかを識別したいとする。この場合、識別モデルは、単語が車名である場合には当該対象データに正例を付与し、単語が車名でない場合には当該対象データに負例を付与する。この場合、ブートストラップ法における各繰り返しにおいて得られた正例集合の一例を以下に示す。

40

【0023】

iteration 1: シビック, ヴィッツ

iteration 2: プリウス, スカイライン

...

iteration 9: カブ, ヘリ, フィット

iteration10: ブルートレイン, エアバス

【0024】

演算部20は、正例集合データベース200と、擬似正例負例選択部202と、擬似正

50

例データベース204と、擬似負例データベース206と、識別対象データベース208と、素性化参照用データベース210と、素性化部212と、素性化済み訓練データベース214と、素性化済み識別対象データベース216と、モデル学習部218と、識別モデルデータベース220と、スコアリング部222と、判定部224と、新正例データベース226と、新負例データベース228とを備えている。

【0025】

正例集合データベース200には、入力部10によって受け付けた、ブートストラップ法における繰り返し毎に得られた正例集合が格納される。

【0026】

擬似正例負例選択部202は、正例集合データベース200に格納された繰り返し毎の正例集合のうち、ブートストラップ法における繰り返し開始から予め定められた初期の第1の繰り返し回数Bまでに得られた正例集合を擬似正例データとして選択する。

10

【0027】

また、擬似正例負例選択部202は、正例集合データベース200に格納された繰り返し毎の正例集合のうち、予め定められた終盤の第2の繰り返し回数Eから繰り返し終了までに得られた正例集合を擬似負例データとして選択する。

【0028】

そして、擬似正例負例選択部202は、それ以外の繰り返し（繰り返し回数B+1～繰り返し回数E-1）において得られた正例集合を識別対象データとして選択する。

【0029】

20

擬似正例データ、及び擬似負例データは、ブートストラップ法における各繰り返しの正例集合に対して、擬似的に正例・負例のラベルを付与したものである。例えば、対象データが単語であって、単語が車名であるかどうかを識別する場合の、擬似正例データ、擬似負例データ、及び識別対象データの一例を以下に示す。

【0030】

（例）

擬似正例データ(iteration1を取得)：シビック，ヴィッツ

擬似負例データ(iteration10を取得)：ブルートレイン，エアバス

識別対象データ(iteration 2~9) プリウス，スカイライン，…，カブ，ヘリ

【0031】

30

上記の例では、ブートストラップ法における繰り返し開始から1回目の繰り返しまでに得られた正例集合を擬似正例データとして選択し、10回目の繰り返しから繰り返し終了までに得られた正例集合を擬似負例データとして選択している。そして、それ以外の繰り返し（2回目の繰り返しから9回目の繰り返しまで）において得られた正例集合を識別対象データとして選択している。

【0032】

擬似正例データベース204には、擬似正例負例選択部202によって選択された擬似正例データが格納される。また、擬似負例データベース206には、擬似正例負例選択部202によって選択された擬似負例データが格納される。識別対象データベース208には、擬似正例負例選択部202によって選択された識別対象データが格納される。

40

【0033】

素性化参照用データベース210には、入力部10によって受け付けた素性化参照用データが格納される。素性化参照用データは、正例集合データベース200に格納された正例集合に含まれる正例データの各々を、素性化するために必要となるデータである。素性化参照用データベース210には、例えば、形態素済みテキストや、文書に付与されたトピック情報等が、素性化参照用データとして格納されている。

【0034】

なお、上述したように、例えば、対象データが単語であって、単語が車名であるかどうかを識別したい場合には、上記の正例集合の例で示した「シビック」「ヴィッツ」等の単語が含まれる文書についての形態素済み文書や、当該文書に付与されたトピック情報等が

50

、素性化参照用データとなる。

【 0 0 3 5 】

素性化部 2 1 2 は、擬似正例データベース 2 0 4 に格納された擬似正例データの各々と、擬似負例データベース 2 0 6 に格納された擬似負例データの各々について、素性化参照用データベース 2 1 0 に格納された素性化参照用データに基づいて、擬似正例データの素性、及び擬似負例データの素性を抽出し、素性化済み訓練データとする。素性化済み訓練データは、擬似正例データと擬似負例データとを、識別モデルを学習するために、素性関数を用いて変換したデータである。

【 0 0 3 6 】

また、素性化部 2 1 2 は、識別対象データベース 2 0 8 に格納された識別対象データの各々について、素性化参照用データベース 2 1 0 に格納された素性化参照用データに基づいて、識別対象データの素性を抽出し、素性化済み識別対象データとする。素性化済み識別対象データは、識別対象データを、識別モデルを適用するために、素性関数を用いて変換したデータである。

【 0 0 3 7 】

擬似正例データと擬似負例データとを、素性関数により変換したデータの一例を以下に示す。

【 0 0 3 8 】

例：（最初の + 1 / - 1 は擬似正例 / 擬似負例を表す）

+ 1 シビック

素性 I D 1：（1 つ後ろから「に / 乗る」が連鎖する回数 = ） 5

素性 I D 2：（「車」との共起回数 = ） 1 0 0

- 1 ブルートレイン

素性 I D 1：（1 つ後ろから「に / 乗る」が連鎖する回数 = ） 3

素性 I D 2：（「車」との共起回数 = ） 0

【 0 0 3 9 】

なお、具体的な素性関数は従来提案されたものでよく、例えば対象データが単語である場合には、単語と共起する任意の単語の出現数や、対象単語の直後に連鎖して出現した任意の単語の出現数などがあげられる。

【 0 0 4 0 】

素性化済み訓練データベース 2 1 4 には、素性化部 2 1 2 によって変換された素性化済み訓練データが格納される。

【 0 0 4 1 】

素性化済み識別対象データベース 2 1 6 には、素性化部 2 1 2 によって変換された素性化済み識別対象データが格納される。

【 0 0 4 2 】

モデル学習部 2 1 8 は、素性化済み訓練データベース 2 1 4 に格納された素性化済み訓練データに基づいて、識別モデルを学習する。なお識別モデルとしては、例えば、SVM や l o g i s t i c 回帰モデルなど、各素性に対し正負それぞれに対する重み（信頼度）を記憶するモデルを用いる。また、識別モデルの学習は、従来の機械学習法によって、識別モデルを学習する。学習方法についても従来の手法を用いてよい。

【 0 0 4 3 】

識別モデルデータベース 2 2 0 には、モデル学習部 2 1 8 によって学習された識別モデルが格納される。

【 0 0 4 4 】

スコアリング部 2 2 2 は、識別モデルデータベース 2 2 0 に格納された識別モデルと、素性化済み識別対象データベース 2 1 6 に格納されている素性化済み識別対象データとに基づいて、識別対象データの各々について、当該識別対象データが正例である度合いを表すスコアを算出する。

【 0 0 4 5 】

10

20

30

40

50

判定部 224 は、識別対象データの各々について、スコアリング部 222 によって算出されたスコアが、予め定められた第 1 の閾値 P より大きい場合には、当該識別対象データが正例であると判定する。また、判定部 224 は、識別対象データの各々について、スコアリング部 222 によって算出されたスコアが、予め定められた第 2 の閾値 N より小さい場合には、当該識別対象データが負例であると判定する。

【0046】

新正例データベース 226 には、判定部 224 によって正例であると判定された識別対象データが、新正例データとして格納される。

【0047】

新負例データベース 228 には、判定部 224 によって負例であると判定された識別対象データが、新負例データとして格納される。

10

【0048】

出力部 30 は、新正例データベース 226 に格納された新正例データの各々と、新負例データベース 228 に格納された新負例データの各々とを出力する。例えば、対象データが車の名称である場合に、出力部 30 が出力するデータの一例を以下に示す。

【0049】

例：

新正例データ：プリウス，スカイライン，フィット

新負例データ：カブ，ヘリ

【0050】

20

< データ解析装置の作用 >

次に、第 1 の実施の形態に係るデータ解析装置 100 の作用について説明する。まず、ブートストラップ法における繰り返し毎に得られた正例集合と、正例集合に含まれる正例データ毎の素性化参照用データとが、データ解析装置 100 に入力されると、データ解析装置 100 によって、図 2 に示すデータ解析処理ルーチンが実行される。

【0051】

まず、ステップ S100 において、入力部 10 によって、ブートストラップ法における繰り返し毎に得られた正例集合の入力を受け付ける。そして、入力部 10 によって、ブートストラップ法における繰り返し毎に得られた正例集合を正例集合データベース 200 に格納する。また、入力部 10 によって、素性化参照用データを受け付け、素性化参照用データベース 210 に格納する。

30

【0052】

ステップ S102 において、擬似正例負例選択部 202 によって、上記ステップ S100 で正例集合データベース 200 に格納された繰り返し毎の正例集合のうち、ブートストラップ法における繰り返し開始から予め定められた初期の第 1 の繰り返し回数 B までに得られた正例集合を擬似正例データとして選択し、擬似正例データベース 204 に格納する。

【0053】

ステップ S104 において、擬似正例負例選択部 202 によって、上記ステップ S100 で正例集合データベース 200 に格納された繰り返し毎の正例集合のうち、予め定められた終盤の第 2 の繰り返し回数 E から繰り返し終了までに得られた正例集合を擬似負例データとして選択し、擬似負例データベース 206 に格納する。

40

【0054】

ステップ S106 において、擬似正例負例選択部 202 によって、それ以外の繰り返し（繰り返し回数 B + 1 ~ 繰り返し回数 E - 1）において得られた正例集合を識別対象データとし、識別対象データベース 208 に格納する。

【0055】

ステップ S108 において、素性化部 212 によって、上記ステップ S102 で擬似正例データベース 204 に格納された擬似正例データの各々と、上記ステップ S104 で擬似負例データベース 206 に格納された擬似負例データの各々について、上記ステップ S

50



100で素性化参照用データベース210に格納された素性化参照用データに基づいて、擬似正例データの素性、及び擬似負例データの素性を抽出し、素性化済み訓練データとする。そして、素性化部212によって、得られた素性化済み訓練データを、素性化済み訓練データベース214に格納する。

【0056】

ステップS110において、素性化部212によって、上記ステップS106で識別対象データベース208に格納された識別対象データの各々について、上記ステップS100で素性化参照用データベース210に格納された素性化参照用データに基づいて、識別対象データの素性を抽出し、素性化済み識別対象データとする。そして、素性化部212によって、得られた素性化済み識別対象データを、素性化済み識別対象データベース216に格納する。

10

【0057】

ステップS112において、モデル学習部218によって、上記ステップS108で素性化済み訓練データベース214に格納された素性化済み訓練データに基づいて、識別モデルを学習し、識別モデルデータベース220に格納する。

【0058】

ステップS114において、スコアリング部222によって、上記ステップS112で識別モデルデータベース220に格納された識別モデルと、上記ステップS110で素性化済み識別対象データベース216に格納された素性化済み識別対象データとに基づいて、識別対象データの各々について、当該識別対象データが正例である度合いを表すスコアを算出する。

20

【0059】

ステップS116において、判定部224によって、識別対象データの各々について、上記ステップS114で算出されたスコアが、予め定められた第1の閾値Pより大きい場合には、当該識別対象データが正例であると判定する。また、判定部224によって、識別対象データの各々について、上記ステップS114で算出されたスコアが、予め定められた第2の閾値Nより小さい場合には、当該識別対象データが負例であると判定する。そして、判定部224によって、正例であると判定された識別対象データを新正例データとして新正例データベース226に格納する。また、判定部224によって、負例であると判定された識別対象データを新負例データとして新負例データベース228に格納する。

30

【0060】

ステップS118において、上記ステップS116で新正例データベース226に格納された新正例データの各々と、新負例データベース228に格納された新負例データの各々とを、結果として出力して、データ解析処理ルーチンを終了する。

【0061】

以上説明したように、第1の実施の形態に係るデータ解析装置によれば、ブートストラップ法における繰り返し毎の正例集合のうち、ブートストラップ法における繰り返し開始から予め定められた初期の第1の繰り返し回数Bまでに得られた正例集合を擬似正例データとして選択し、予め定められた終盤の第2の繰り返し回数Eから繰り返し終了までに得られた正例集合を擬似負例データとして選択し、選択された擬似正例データ、及び擬似負例データから抽出された素性に基づいて、識別モデルを学習し、識別対象データの各々について、学習された識別モデルと、識別対象データから抽出された素性とに基づいて、識別対象データが正例である度合いを表すスコアを算出し、スコアが予め定められた第2の閾値Nより小さい場合には、識別対象データが負例であると判定することにより、負例データを精度よく抽出することができる。

40

【0062】

また、識別モデルとして、一般の識別モデル（サポートベクタマシンや最大エントロピー法）を使うことが可能であるため、自由度の高い情報を用いることが可能となる。

【0063】

また、ブートストラップ法における繰り返し初期と繰り返し終盤とにおける正例データ

50

を、擬似正例データ、及び擬似負例データとみなし、自由度の高い素性表現と識別モデルを併用することで、semantic driftを抑制し、高い精度で新負例データを抽出することができる。

【0064】

<第2の実施の形態>

次に、第2の実施の形態について説明する。なお、第1の実施の形態と同様の構成となる部分については、同一符号を付して説明を省略する。

【0065】

第2の実施の形態では、ブートストラップ法の各繰り返しにおいて、正例・負例と判定された未知データの割合を用いて、semantic driftの度合いを推測し、それに応じて第2の閾値を変化させる点が、第1の実施の形態と異なっている。

10

【0066】

<システム構成>

第2の実施の形態における入力部12は、ブートストラップ法における繰り返し毎に、正例のラベルが付与された未知データの数、及び負例のラベルが付与された未知データの数の割合を更に受け付ける。

【0067】

判定部2224は、入力部12によって受け付けた、繰り返し毎の上記割合に基づいて、第2の閾値Nを変化させる。具体的には、判定部2224は、ブートストラップ法における繰り返し毎の上記割合に基づいて、第2の閾値Nを変化させる。

20

【0068】

例えば、iteration1において51:49の割合で正例、負例が判定され、iteration9において48:52の割合で正例、負例が判定された場合、iterationが進行しても負例の量が増えていない。これはsemantic driftがそれほど生じていないことを表していると考えられるため、第2の閾値Nを引き下げるように変化させる(負例になりづらくする)。

【0069】

一方、例えば、iteration1において90:10の割合で正例、負例が判定され、iteration9において10:90の割合で正例、負例が判定された場合、iterationが進行するに従い負例の量が増えている。これはsemantic driftが生じていることを表していると考えられるため、第2の閾値Nを引き下げないこととする。

30

【0070】

このときの第2の閾値Nの変動量は、固定値でも良いし、semantic driftの度合いを示す数値に比例した値としてもよい。第2の閾値Nの変動量が、semantic driftの度合いを示す数値に比例した値とする場合には、例えば、semantic driftの度合いを、最終繰り返しでの負例の数と初期繰り返しでの負例の数との比とし(49/52)、当該比と一定量Kとの積を算出し、第2の閾値Nの変動量とする。

【0071】

また、判定部2224は、識別対象データの各々について、スコアリング部222によって算出されたスコアが、第1の閾値Pより大きい場合には、識別対象データが正例であると判定し、スコアリング部222によって算出されたスコアが、変化した第2の閾値Nより小さい場合には、識別対象データが負例であると判定する。

40

【0072】

<データ解析装置の作用>

次に、第2の実施の形態に係るデータ解析装置の作用について説明する。まず、ブートストラップ法における繰り返し毎に得られた正例集合と、正例集合に含まれる正例データ毎の素性化参照用データと、ブートストラップ法における繰り返し毎の、正例のラベルが付与された未知データの数、及び負例のラベルが付与された未知データの数の割合とが、データ解析装置に入力されると、データ解析装置によって、図3に示すデータ解析処理ルーチンが実行される。なお、第1の実施の形態と同様の処理については、同一符号を付し

50

て説明を省略する。

【0073】

まず、ステップS200において、入力部12によって、ブートストラップ法における繰り返し毎に得られた正例集合の入力を受け付ける。そして、入力部12によって、ブートストラップ法における繰り返し毎に得られた正例集合を正例集合データベース200に格納する。また、入力部12によって、素性化参照用データを受け付け、素性化参照用データベース210に格納する。また、入力部12によって、正例のラベルが付与された未知データの数、及び負例のラベルが付与された未知データの数の割合を、ブートストラップ法における繰り返し毎に受け付ける。

【0074】

そして、ステップS215において、判定部2224によって、上記ステップS200で受け付けた、ブートストラップ法における繰り返し毎の、正例のラベルが付与された未知データの数、及び負例のラベルが付与された未知データの数の割合に応じて、第2の閾値Nを変化させる。

【0075】

そして、ステップS216において、判定部2224によって、識別対象データの各々について、上記ステップS114で算出されたスコアが、予め定められた第1の閾値Pより大きい場合には、当該識別対象データが正例であると判定し、上記ステップS114で算出されたスコアが、上記ステップS215で変化させた第2の閾値Nより小さい場合には、当該識別対象データが負例であると判定する。そして、判定部2224によって、正例であると判定された識別対象データを新正例データとして新正例データベース226に格納する。また、判定部2224によって、負例であると判定された識別対象データを新負例データとして新負例データベース228に格納する。

【0076】

以上説明したように、第2の実施の形態に係るデータ解析装置によれば、ブートストラップ法における各繰り返しにおいて、正例のラベルが付与された未知データの数、及び負例のラベルが付与された未知データの数の割合に基づいて、第2の閾値Nを変化させ、識別対象データの各々について、算出されたスコアが、変化した第2の閾値Nより小さい場合には、当該識別対象データが負例であると判定することにより、semantic driftの度合いを考慮して、負例データを精度よく抽出することができる。

【0077】

なお、本発明は、上述した実施形態に限定されるものではなく、この発明の要旨を逸脱しない範囲内で様々な変形や応用が可能である。

【0078】

例えば、本実施の形態では、識別対象データを、擬似正例データ及び擬似負例データとして選択された正例集合と異なる正例集合とする場合を例に説明したが、これに限定されるものではない。例えば、識別対象データに、擬似正例データ及び擬似負例データを含めてもよい。また、識別対象データに、負例集合を含めてもよい。

【0079】

また、本実施の形態のデータ解析装置は、正例集合データベース200、擬似正例データベース204、擬似負例データベース206、識別対象データベース208、素性化参照用データベース210、素性化済み訓練データベース214、素性化済み識別対象データベース216、識別モデルデータベース220、新正例データベース226、及び新負例データベース228を備えている場合について説明したが、例えば正例集合データベース200、擬似正例データベース204、擬似負例データベース206、識別対象データベース208、素性化参照用データベース210、素性化済み訓練データベース214、素性化済み識別対象データベース216、識別モデルデータベース220、新正例データベース226、及び新負例データベース228の少なくとも1つがデータ解析装置の外部装置に設けられ、データ解析装置は、外部装置と通信手段を用いて通信することにより、正例集合データベース200、擬似正例データベース204、擬似負例データベース20

10

20

30

40

50

6、識別対象データベース208、素性化参照用データベース210、素性化済み訓練データベース214、素性化済み識別対象データベース216、識別モデルデータベース220、新正例データベース226、及び新負例データベース228の少なくとも1つを参照するようにしてもよい。

【0080】

また、上述のデータ解析装置は、内部にコンピュータシステムを有しているが、「コンピュータシステム」は、WWWシステムを利用している場合であれば、ホームページ提供環境（あるいは表示環境）も含むものとする。

【0081】

また、本願明細書中において、プログラムが予めインストールされている実施形態として説明したが、当該プログラムを、コンピュータ読み取り可能な記録媒体に格納して提供することも可能である。

10

【符号の説明】

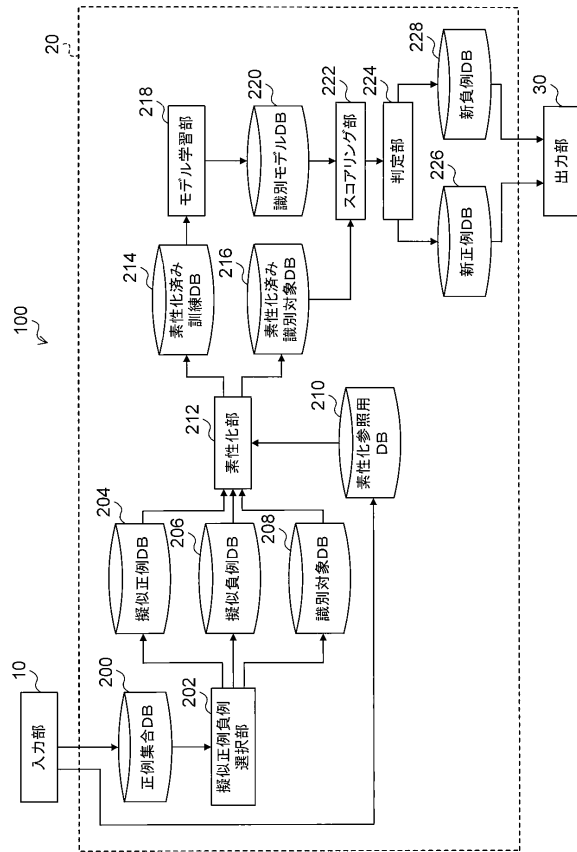
【0082】

10, 12 入力部  
 20 演算部  
 30 出力部  
 100 データ解析装置  
 200 正例集合データベース  
 202 擬似正例負例選択部  
 204 擬似正例データベース  
 206 擬似負例データベース  
 208 識別対象データベース  
 210 素性化参照用データベース  
 212 素性化部  
 214 素性化済み訓練データベース  
 216 素性化済み識別対象データベース  
 218 モデル学習部  
 220 識別モデルデータベース  
 222 スコアリング部  
 224, 2224 判定部  
 226 新正例データベース  
 228 新負例データベース

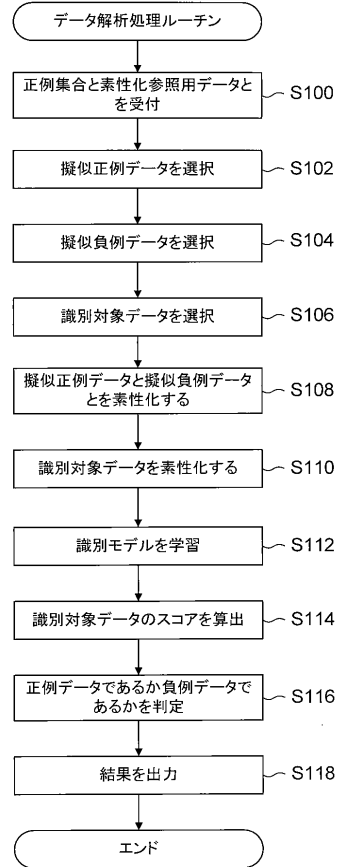
20

30

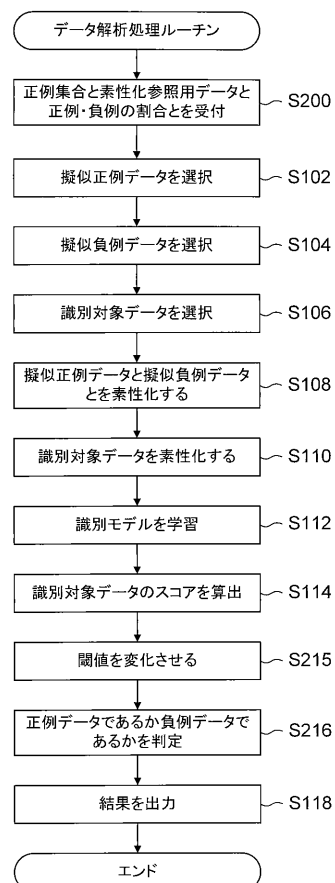
【 図 1 】



【 図 2 】



【圖 3】



---

フロントページの続き

審査官 多胡 滋

(56)参考文献 貞光九月，外4名，トピック情報を用いたブートストラップ法に基づく語彙獲得，自然言語処理  
，日本，言語処理学会，2012年 7月 6日，第19巻，第2号，pp.89-106

(58)調査した分野(Int.Cl.，DB名)  
G06N 99/00