

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2009-282686
(P2009-282686A)

(43) 公開日 平成21年12月3日(2009.12.3)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 17/30 (2006.01)	G06F 17/30 210D	5B075
G06N 5/04 (2006.01)	G06F 17/30 180B	5L096
G06T 7/00 (2006.01)	G06N 5/04 580A	
	G06N 5/04 580P	
	G06T 7/00 350B	
審査請求 未請求 請求項の数 10 OL (全 16 頁) 最終頁に続く		

(21) 出願番号 特願2008-133224 (P2008-133224)
(22) 出願日 平成20年5月21日 (2008.5.21)

(71) 出願人 000003078
株式会社東芝
東京都港区芝浦一丁目1番1号
(74) 代理人 110000235
特許業務法人 天城国際特許事務所
(72) 発明者 中田 康太
東京都港区芝浦一丁目1番1号 株式会社東芝内
(72) 発明者 櫻井 茂明
東京都港区芝浦一丁目1番1号 株式会社東芝内
(72) 発明者 折原 良平
東京都港区芝浦一丁目1番1号 株式会社東芝内
Fターム(参考) 5B075 NR12 QT10
5L096 BA06 FA66 GA30 KA04 MA07

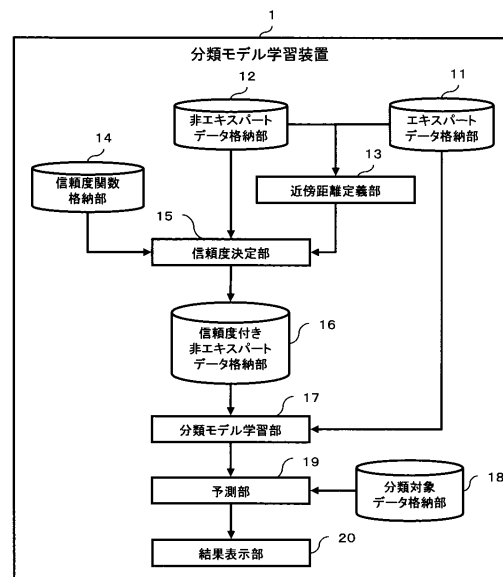
(54) 【発明の名称】 分類モデル学習装置および分類モデル学習方法

(57) 【要約】

【課題】 質の悪い教師データが含まれていても精度の良い分類モデルを構築する。

【解決手段】 ラベル付けの信頼度が所定の基準を満たすエキスパートデータおよびラベル付けの信頼度が不明な非エキスパートデータの各々が対応する座標を取得して非エキスパートデータからエキスパートデータまでの距離を各々算出し、所定の規則に当てはめて近傍距離を定義する。次に、選択した非エキスパートデータから近傍距離の範囲内にあるエキスパートデータを探索して同ラベル確率を算出し、付されたラベルが近傍距離の範囲内にあるエキスパートデータのラベルに一致する確率に基づく信頼度関数に当てはめて非エキスパートデータの信頼度を決定し、付加する。そして、エキスパートデータおよび信頼度が付加された非エキスパートデータに基づいて所望のデータにラベル付けを行う分類モデルを学習する。

【選択図】 図1



【特許請求の範囲】

【請求項 1】

機械学習におけるラベル付けの信頼度が所定の基準を満たす教師データをエキスパートデータとして格納するエキスパートデータ格納部と、

前記ラベル付けの信頼度が不明な教師データを非エキスパートデータとして格納する非エキスパートデータ格納部と、

前記エキスパートデータ格納部および前記非エキスパートデータ格納部に接続され、前記エキスパートデータおよび前記非エキスパートデータの各々が対応する座標を取得して前記非エキスパートデータから前記エキスパートデータまでの距離を各々算出すると共に、この算出された距離を所定の規則に当てはめて近傍距離を定義する近傍距離定義部と、

前記非エキスパートデータに付された前記ラベルが前記近傍距離の範囲内にある前記エキスパートデータに付された前記ラベルに一致する確率に基づく信頼度関数を格納する信頼度関数格納部と、

前記近傍距離定義部、前記信頼度関数格納部、前記エキスパートデータ格納部、および前記非エキスパートデータ格納部に接続され、選択した前記非エキスパートデータから前記近傍距離の範囲内にある前記エキスパートデータを探索して前記確率を算出すると共に、この算出された確率を前記信頼度関数に当てはめて前記非エキスパートデータにおける前記ラベル付けの信頼度を決定する信頼度決定部と、

前記エキスパートデータ格納部および前記信頼度決定部に接続され、前記エキスパートデータおよび前記信頼度が付加された非エキスパートデータに基づいて所望の分類対象データに前記ラベル付けを行う分類モデルを学習する分類モデル学習部と、
を有することを特徴とする分類モデル学習装置。

【請求項 2】

前記信頼度関数は、前記非エキスパートデータに付された前記ラベルが前記近傍距離の範囲内にある前記エキスパートデータに付された前記ラベルに一致する確率と前記ラベル付けの信頼度との関係を定義することを特徴とする請求項 1 記載の分類モデル学習装置。

【請求項 3】

前記近傍距離定義部は、前記非エキスパートデータから前記エキスパートデータまでの前記座標間の距離を各々算出して前記非エキスパートデータ毎に順位付けを行い、所望の順位についての距離を前記非エキスパートデータの各々から集計してその平均値を算出し、この平均値を前記近傍距離として定義することを特徴とする請求項 1 または請求項 2 記載の分類モデル学習装置。

【請求項 4】

前記信頼度関数は、分類問題における所望の評価基準に応じて予め作成されていることを特徴とする請求項 1 乃至請求項 3 のいずれか一項記載の分類モデル学習装置。

【請求項 5】

前記分類モデル学習部が、アンサンブル学習におけるデータ重みに対して前記信頼度決定部で付加された信頼度を反映させることにより前記分類モデルを学習することを特徴とする請求項 1 乃至請求項 4 のいずれか一項記載の分類モデル学習装置。

【請求項 6】

機械学習におけるラベル付けの信頼度が所定の基準を満たしている教師データをエキスパートデータ、前記ラベル付けの信頼度が不明な教師データを非エキスパートデータとして格納するコンピュータが行う分類モデル学習方法であって、

前記エキスパートデータおよび前記非エキスパートデータの各々が対応する座標を取得して前記非エキスパートデータから前記エキスパートデータまでの距離を各々算出すると共に、この算出された距離を所定の規則に当てはめて近傍距離を定義する近傍距離定義ステップと、

前記格納された非エキスパートデータから前記信頼度の付加対象となる非エキスパートデータを選択する選択ステップと、

前記選択された非エキスパートデータから前記近傍距離の範囲内にある前記エキスパー

10

20

30

40

50

トデータを探索して前記非エキスパートデータに付された前記ラベルが前記エキスパートデータに付された前記ラベルに一致する確率を算出すると共に、この算出された確率を予め定義された信頼度関数に当てはめて前記非エキスパートデータの前記ラベル付けの信頼度を決定する信頼度決定ステップと、

前記決定された信頼度が付加された非エキスパートデータおよび前記エキスパートデータに基づいて所望のデータに前記ラベル付けを行う分類モデルを学習する分類モデル学習ステップと、

を有することを特徴とする分類モデル学習方法。

【請求項 7】

前記信頼度関数は、前記非エキスパートデータに付された前記ラベルが前記近傍距離の範囲内にある前記エキスパートデータに付された前記ラベルに一致する確率と前記ラベル付けの信頼度との関係を定義することを特徴とする請求項 6 記載の分類モデル学習方法。

【請求項 8】

前記近傍距離定義ステップにおいて、前記非エキスパートデータから前記エキスパートデータまでの前記座標間の距離を各々算出して前記非エキスパートデータ毎に順位付けを行い、所望の順位についての距離を前記非エキスパートデータの各々から集計してその平均値を算出し、この平均値を前記近傍距離として定義することを特徴とする請求項 6 または請求項 7 記載の分類モデル学習方法。

【請求項 9】

前記信頼度関数は、分類問題における所望の評価基準に応じて予め作成されていることを特徴とする請求項 6 乃至請求項 8 のいずれか一項記載の分類モデル学習方法。

【請求項 10】

前記分類モデル学習ステップにおいて、アンサンブル学習におけるデータ重みに対して前記信頼度決定ステップにおいて付加された信頼度を反映させることにより前記分類モデルを学習することを特徴とする請求項 6 乃至請求項 9 のいずれか一項記載の分類モデル学習方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、機械学習において分類対象データにラベル付けを行う分類モデル学習装置および分類モデル学習方法に関する。

【背景技術】

【0002】

データマイニングにおいて重要な分野の一つに機械学習が挙げられる。機械学習は分類問題に使われることが多く、分類問題において計算機は人間のつけた評価を学習することで分類モデルを構築する。この機械学習の応用は画像認識、文字認識、テキスト分類など広い分野で多くの成果を上げている。このような学習は一般的に教師あり学習と呼ばれる。

【0003】

教師あり学習には計算機に正しい判断を教える「教師データ」、つまり人間の手によって「ラベル」が付けられたデータが必要である。計算機は教師データをもとにどのような分類をすれば良いかを学習し、新しいデータに対して自動的に判断を下せるようになる。現代ではIT環境の発展により、大量で詳細なデータが機械学習に利用可能となっており、これらを教師データとして用いればより正確な分類モデルの構築に繋がると期待されている。

【0004】

しかし、ここで「大量のデータ」から「大量の教師データ」を得る際のラベル付けが問題になっている。すなわち、得られたデータを教師データとして利用するためには、データに対して人間が判断したラベルを付与することが必要であるが、正確なラベル付けには、データが取られたドメインに対しての知識や経験などに基づく正確な判断が不可欠であ

10

20

30

40

50

る。

【0005】

理想的にはこれらの条件を満たす対象分野のエキスパートがラベル付けを行うことが望ましいが、全てのデータのラベル付けをエキスパートに依頼することは非常にコストが高くなってしまふ。しかし、現実的にはコストに制限があるため、大量の教師データが必要である場合には、非エキスパートがラベル付けを行った教師データを用いることになる。エキスパートによる高コストの教師データは少量になりがちなのに対し、非エキスパートによる低コストの教師データは比較的大量に獲得できるためである。その一方、非エキスパートによる教師データには判断のミスや知識の不正確さから、比較的多くの誤ったラベルが含まれてしまうことが考えられる。

10

【0006】

また、一般の機械学習においては、教師データの取得に関する情報は用いられず、エキスパートによるラベル付けのような「良質の教師データ」と非エキスパートによるラベル付けのような「ノイズを含む教師データ」が混在する状況においても、全てのデータを同列に扱い、等しく学習に使用する。

【0007】

したがって、エキスパートによる少量の教師データと非エキスパートによる大量の教師データを従来どおり同列とみなして学習に使用した場合、非エキスパートデータに含まれるノイズが学習に大きく影響し、精度の良い分類モデルが構築できないケースが考えられる。

20

【0008】

一方で、分類モデルを学習する際に、一部の教師データを選択的に使用して学習を行うことや、一部の教師データに重みを置いて学習を行うことが一般的に行われている。アンサンブル学習の代表的手法の1つである Ada Boost もその一つである。Ada Boost は、学習データに対して重みを与えて学習器を生成し、その際に誤った分類をしたデータに対して重みを増して再度学習器を生成することを繰り返して複数の弱学習器を得て、それらの弱学習器の重みつき投票により分類を行う手法である（例えば特許文献1、非特許文献1参照）。

【特許文献1】特開2002-133389号公報

【非特許文献1】Y. Freund and R.E. Shapire, "Experiments with a new boosting algorithm", Proc. of the 13th. Int. Conf. on Machine Learning, 1996, 148-156

30

【発明の開示】

【発明が解決しようとする課題】

【0009】

しかしながら、従来技術は、あくまで所定のアルゴリズムに即した形で教師データに対してデータ重みをつけるものであり、教師データの精度の差異という学習過程を開始する前の知識・情報を含んだものではなく、例えばエキスパートによる少量の教師データ（以下、「エキスパートデータ」という。）と非エキスパートによる大量の教師データ（以下、「非エキスパートデータ」という。）のような、質の異なる教師データを従来どおり同列として学習に使用した場合、質の劣る教師データに含まれるノイズが学習に大きく影響し、精度の良い分類モデルが構築できないという問題があった。

40

【0010】

このような問題に対して、本出願人は特許出願2007-278893においてエキスパートによる少量の教師データを利用することで精度の良い分類モデルの学習を行う手法を提案している。この手法は、エキスパートによる教師データを基にして非エキスパートによる教師データのラベルに信頼度を付加し、分類モデルの学習にその信頼度を反映することで分類モデルを学習するものである。この信頼度は、エキスパートデータおよび非エキスパートデータの各々を所定の規則に基づいて対応付けた座標の間の距離（例えばユークリッド距離やコサイン距離）に応じて求められている。

【0011】

50

対象の非エキスパートデータから距離の近いN個のエキスパートデータを探索し、もしラベルが同じであればそのエキスパートデータから信頼度を得る。この信頼度は例えば距離に反比例する形で与えられ、非エキスパートデータの近くのエキスパートデータが同じラベルであれば、その非エキスパートデータは高い信頼度を得られるようになっている。これは、信頼できるデータが近くにあるほど信頼度は高いという直感的な信頼度付けを表していると言える。

【0012】

しかしながら、上記の信頼度付け方法は、エキスパートデータには誤ラベルが含まれていないことを暗に仮定している。エキスパートデータに全て適切なラベルが与えられているならば、それらを参照して与えられた非エキスパートデータの信頼度も適切な値になることが期待できる。その反面、エキスパートデータに誤ラベルが含まれている場合には、このような信頼度の付加は必ずしも適切とは言えない。図10および図11は、エキスパートデータのラベル付けと非エキスパートデータのラベル付けに対する信頼度の関係を説明する図である。

10

【0013】

図10では、ある非エキスパートデータ x_1 の非常に近傍にエキスパートデータ X_1 が存在している。この X_1 は非常に近傍にあるため、 X_1 と x_1 のラベルが同じであれば x_1 の信頼度は高く、異なれば低くなる。ここで X_1 、 x_1 に本来付与されるべきラベルは L_1 であるとする。エキスパートデータ X_1 に、正確なラベル L_1 が付与されているとすると、非エキスパートデータ x_1 に L_1 が付与されている場合には信頼度は高く、異なったラベル L_2 が付与されている場合には信頼度は低くなる。これは、適切な信頼度であるといえる。

20

【0014】

図11では、エキスパートデータ X_1 に誤ラベル L_2 が付与されている場合を考える。このとき、非エキスパートデータ x_1 に本来付与されるべきラベル L_1 が付与されていたときは信頼度が低く、反対に付与されるべきでないラベル L_2 が付与されていたときに信頼度が高くなってしまふ。これは明らかに適切な信頼度とは反対の傾向である。

【0015】

すなわち、エキスパートデータ中に誤ラベルが含まれている場合、非エキスパートデータに適切でない信頼度が付加され、その信頼度を反映して生成される分類モデルの性能が劣化してしまう。現実にはエキスパートデータ中にも少量の誤ラベルが含まれると考えられるため、エキスパートデータ中の誤ラベルに頑健な信頼度付加が必要である。

30

【0016】

そこで、本発明は、従来技術の問題に鑑み、質の悪い教師データが含まれている状況であっても精度の良い分類モデルの構築が可能な分類モデル学習装置および分類モデル学習方法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】

【0017】

本発明に係る分類モデル学習装置は、機械学習におけるラベル付けの信頼度が所定の基準を満たす教師データをエキスパートデータとして格納するエキスパートデータ格納部と、前記ラベル付けの信頼度が不明な教師データを非エキスパートデータとして格納する非エキスパートデータ格納部と、前記エキスパートデータ格納部および前記非エキスパートデータ格納部に接続され、前記エキスパートデータおよび前記非エキスパートデータの各々が対応する座標を取得して前記非エキスパートデータから前記エキスパートデータまでの距離を各々算出すると共に、この算出された距離を所定の規則に当てはめて近傍距離を定義する近傍距離定義部と、前記非エキスパートデータに付された前記ラベルが前記近傍距離の範囲内にある前記エキスパートデータに付された前記ラベルに一致する確率に基づく信頼度関数を格納する信頼度関数格納部と、前記近傍距離定義部、前記信頼度関数格納部、前記エキスパートデータ格納部、および前記非エキスパートデータ格納部に接続され、選択した前記非エキスパートデータから前記近傍距離の範囲内にある前記エキスパートデータを探索して前記確率を算出すると共に、この算出された確率を前記信頼度関数に当

40

50

てはめて前記非エキスパートデータにおける前記ラベル付けの信頼度を決定する信頼度決定部と、前記エキスパートデータ格納部および前記信頼度決定部に接続され、前記エキスパートデータおよび前記信頼度が付加された非エキスパートデータに基づいて所望の分類対象データに前記ラベル付けを行う分類モデルを学習する分類モデル学習部と、を有することを特徴とする。

【0018】

本発明に係る分類モデル学習方法は、機械学習におけるラベル付けの信頼度が所定の基準を満たしている教師データをエキスパートデータ、前記ラベル付けの信頼度が不明の教師データを非エキスパートデータとして格納するコンピュータが行う分類モデル学習方法であって、前記エキスパートデータおよび前記非エキスパートデータの各々が対応する座標を取得して前記非エキスパートデータから前記エキスパートデータまでの距離を各々算出すると共に、この算出された距離を所定の規則に当てはめて近傍距離を定義する近傍距離定義ステップと、前記格納された非エキスパートデータから前記信頼度の付加対象となる非エキスパートデータを選択する選択ステップと、前記選択された非エキスパートデータから前記近傍距離の範囲内にある前記エキスパートデータを探索して前記非エキスパートデータに付された前記ラベルが前記エキスパートデータに付された前記ラベルに一致する確率を算出すると共に、この算出された確率を予め定義された信頼度関数に当てはめて前記非エキスパートデータの前記ラベル付けの信頼度を決定する信頼度決定ステップと、前記決定された信頼度が付加された非エキスパートデータおよび前記エキスパートデータに基づいて所望のデータに前記ラベル付けを行う分類モデルを学習する分類モデル学習ステップと、を有することを特徴とする。

【発明の効果】

【0019】

本発明によれば、質の悪い教師データが含まれている状況であっても精度の良い分類モデルの構築が可能な分類モデル学習装置および分類モデル学習方法が提供される。

【発明を実施するための最良の形態】

【0020】

以下、本発明の実施形態について図面を用いて説明する。図1は、本発明の一実施形態に係る分類モデル学習装置1の全体構成例を示すブロック図である。同図に示されるように、本実施形態に係る分類モデル学習装置1は、エキスパートデータ格納部11、非エキスパートデータ格納部12、近傍距離定義部13、信頼度関数格納部14、信頼度決定部15、信頼度付き非エキスパートデータ格納部16、分類モデル学習部17、分類対象データ格納部18、予測部19、および結果表示部20から構成されている。

【0021】

エキスパートデータ格納部11は、エキスパートデータを格納する記憶装置である。「エキスパートデータ」とは、知識、経験の豊富な専門家が機械学習においてデータを分類するためのラベル付けを行われており、ラベル付けの精度（信頼性）が高い教師データを示すものとする。

【0022】

非エキスパートデータ格納部12は、非エキスパートデータを格納する記憶装置である。「非エキスパートデータ」とは、ラベル付けは行われているが、その精度（信頼性）が不明確な教師データを示すものとする。

【0023】

近傍距離定義部13は、非エキスパートデータからエキスパートデータまでの座標間距離を各々算出し、この座標間距離に基づいてデータ間の類似度が基準値以上の範囲を表す近傍距離を定義するプログラムである。ここでは、算出された座標間距離の中から所定の規則に基づいて複数の距離を選択し、これらの距離の平均値から近傍距離を算出するが、算出方法はこれに限られない。

【0024】

図2は、エキスパートデータおよび非エキスパートデータを2次元で具体的に説明する

図である。同図において、丸印はエキスパートデータ、四角印は非エキスパートデータを表し、各印の色はラベルを表している。これらの座標は各データを所定の規則に基づいて変換することで得られる。例えば、電子メールの分類においては、多数の迷惑メールを解析することによって特徴語リストを予め作成しておき、この特徴語リストと受信メール本文内の単語を比較することで座標化を行う。具体的には、特徴語リストに含まれるN個の単語との比較結果を受信メール内に含まれる場合を1、含まれない場合を0として表すことにより、受信メールのデータをN次元の座標（例えば(1, 0, 1, ..., 1)）に変換できる。ここでは、説明のためにメールデータを座標化したN次元のデータを擬似的に2次元で表しているものとする。すなわち、受信メール本文の内容が近似する場合には、座標が近似するので迷惑メールか否かのラベル付け等に用いることができる。

10

【0025】

また、図2においては、近傍距離定義部13が非エキスパートデータを選択し、この選択された非エキスパートデータから各エキスパートデータまでの距離を順次求めることが示されている。例えば、近傍距離を“非エキスパートデータから4番目に近いエキスパートデータまでの距離の平均値”とする規則が予め定められている場合には、距離 r_4 を非エキスパートデータ毎に求め、その平均値を算出する。

【0026】

信頼度関数格納部14は、分類問題に適した信頼度関数を格納する記憶装置である。この信頼度関数は、非エキスパートデータから近傍距離内にあるエキスパートデータの同ラベル確率に基づいて信頼度を定義する関数であり、この関数は種々の分類問題に対応させて予め複数作成しておくことが好適である。具体的な定義方法は後述する。

20

【0027】

信頼度決定部15は、近傍距離定義部13により定義された近傍距離に基づいて非エキスパートデータの近傍にあるエキスパートデータを探索すると共に非エキスパートデータとの同ラベル確率を算出し、この同ラベル確率を信頼度関数格納部14から取得される信頼度関数に当てはめて非エキスパートデータの信頼度を決定するプログラムである。尚、複数の信頼度関数の内、どの関数を用いるか選択する方法としては、モデル作成時にユーザが入力装置（図示省略する）から入力した情報に基づいて選択する方法や使用する関数を予め設定しておく方法などが挙げられる。

【0028】

30

信頼度付き非エキスパートデータ格納部16は、信頼度決定部15における処理によって信頼度が付与された非エキスパートデータ（以下、「信頼度付き非エキスパートデータ」という。）を格納する記憶装置である。

【0029】

分類モデル学習部17は、エキスパートデータと信頼度付き非エキスパートデータを用いて分類モデルを学習するプログラムである。

【0030】

分類対象データ格納部18は、新たに分類の対象となるデータ、すなわち、ラベルが付与されていないデータ（以下、「分類対象データ」という。）を格納する記憶装置である。

40

【0031】

予測部19は、分類モデル学習部17で得られた分類モデルを用いて分類対象データ格納部18に格納されている分類対象データにラベル付けを行うプログラムである。尚、AdaBoostを用いた場合、予測部19での分類手法は、一般的なAdaBoostにおける手法と同様であるので説明は省略する。

【0032】

結果表示部20は、予測部19における予測結果を表示するディスプレイなどの表示装置である。

【0033】

以下、分類モデル学習装置1における動作を図面に基づいて説明する。尚、本実施形態

50

においては、エキスパートデータおよび非エキスパートデータを2次元のデータとして具体的に説明する。図3は、近傍距離定義部13における処理の具体例を示すフローチャートである。

【0034】

S301においては、未だ選択されていない非エキスパートデータが存在するか否かを判断する。ここで、全ての非エキスパートデータが選択済みであればS305へ進む。これに対し、選択されていない非エキスパートデータが存在する場合にはS302へ進む。

【0035】

S302においては、非エキスパートデータ格納部12から未だ選択されていない非エキスパートデータを一つ選択する。

S303においては、選択された非エキスパートデータから全てのエキスパートデータへの距離を各々算出する。

【0036】

S304においては、選択された非エキスパートデータからk番目に近いエキスパートデータまでの距離をバッファ領域(図示省略する)に保持する。尚、最適な整数kは問題によって異なるが、ここでは整数kをユーザが予め設定した値とする。例えば、S303で算出された距離の分布を解析し、各非エキスパートデータからの距離が所定の範囲内にあるように整数kを設定することができる。また、信頼度の付加にあたって複数の近傍エキスパートデータを考慮したい場合などには整数kを大きくすれば良い。

S305においては、保持していた全ての距離の平均をとり、その値を近傍距離として信頼度決定部15へ出力し、処理を終了する。

【0037】

以上の処理により、k番目に近いエキスパートデータまでの平均距離が求められる。問題に適した整数kを設定すれば、この距離は近傍を定義する典型的な値をとることができる。

【0038】

図4は、信頼度決定部15における処理の具体例を示すフローチャートである。S401においては、選択する非エキスパートデータが存在するか否かを判断する。ここで、全ての非エキスパートデータに信頼度が付与されており選択する非エキスパートデータがなければ処理を終了する。これに対し、信頼度が付与されていない非エキスパートデータが存在する場合にはS402へ進む。

【0039】

S402においては、非エキスパートデータ格納部12から未だ信頼度が付与されていない非エキスパートデータを1つ選択する。ここでは、下記の式(1)で表されるj番目の非エキスパートデータが選択されているとする。尚、xは座標、yはラベルを表すものとする。

【数1】

$$[\bar{x}_j, y_j] = [(1.0, 2.0), 1] \quad \dots (1)$$

【0040】

S403においては、選択された非エキスパートデータの近傍に含まれるエキスパートデータをエキスパートデータ格納部11から探索して保持する。この例では、「近傍」とは近傍距離定義部13において定義された近傍距離rを用いて、上記式(1)で表される非エキスパートデータを中心とした半径rの円の中の領域を指すものとする。したがって、近傍距離rが0.5ときは、下記の式(2)のエキスパートデータ X_{j_1} は近傍に含まれるが、式(3)のエキスパートデータ X_{j_2} は近傍には含まれない。

10

20

30

40

【数 2】

$$[\vec{X}_{j_1}, Y_{j_1}] = [(1.1, 2.3), 1] \quad \dots (2)$$

$$[\vec{X}_{j_2}, Y_{j_2}] = [(1.5, 2.3), 1] \quad \dots (3)$$

【0041】

S404においては、探索されたN個のエキスパートデータから同ラベル確率を算出する。この例では、対象の非エキスパートデータ x_j と同ラベルの近傍エキスパートデータの数をK個とし、同ラベル確率 P_j を下記の式(4)で定義する。

$$P_j = K / N \quad \dots (4)$$

【0042】

S405においては、式(4)を入力とする信頼度関数を用いて非エキスパートデータのラベルの信頼度を算出する。信頼度関数は分類問題によって適した形が考えられる。図5乃至図7は、分類問題の評価基準に応じた信頼度関数の具体例を説明する図である。この信頼度関数の性質の直感的な理解のために、対象となっている式(1)が表す非エキスパートデータの近傍にエキスパートデータが10例含まれており、さらにノイズのため本来は9例が同ラベルであるところ8例が同ラベルとなっている状況を考える。

【0043】

この状況下で、例えば、非エキスパートデータのラベル付けが近傍エキスパートデータの8割以上と一致するならば、そのラベル付けに高信頼度を与えたい場合には、下記の式(5)のような信頼度関数を用いると好適である。尚、aは関数の形を決定するパラメータである。

$$c_j = C(P_j) = \frac{\exp(a(P_j - 1/2))}{\exp(a(P_j - 1/2)) + \exp(-a(P_j - 1/2))} / \frac{\exp(1/2a)}{\exp(1/2a) + \exp(-1/2a)} \quad \dots (5)$$

【0044】

図5は、式(5)の信頼度関数を説明する図である。ここでは、横軸を同ラベル確率(P_j)、縦軸を信頼度(c_j)とし、 $a = 2.0$ の場合に式(5)によって求められる点を結んだ曲線で示されている。同ラベル数が9例から8例に変化するときノイズによる信頼度 c_j の変化は $C(9/10) = 0.98$ から $C(8/10) = 0.96$ となり、信頼度 c_j への影響は小さい。すなわち、近傍の10例中の同ラベルが9例、8例のいずれの場合であっても、その非エキスパートデータのラベルの信頼度は高く維持されるという結果が得られる設定になっており、直感的にも妥当な信頼度関数であると言える。

【0045】

また、誤ラベルの混入に対して厳しい設定としたい場合には、下記の式(6)のような信頼度関数を用いると好適である。

【数 4】

$$c_j = C(P_j) = \frac{\exp(aP_j) - 1.0}{\exp(a) - 1.0} \quad \dots (6)$$

【0046】

図6は、式(6)の信頼度関数を示す図である。ここでは、横軸を同ラベル確率(P_j)、縦軸を信頼度(c_j)とし、 $a = 5.0$ の場合に式(6)によって求められる点を結んだ曲線で示されている。この関数を用いる場合には、一つでも誤ラベルがあると信頼度が大

10

20

30

40

50

幅に下がる。例えば、医療などの高い信頼度が要求される分野において特に有用である。

【 0 0 4 7 】

更に、誤ラベルの混入に対して寛容な設定としたい場合には、下記の式 (7) のような信頼度関数を用いると好適である。

【 数 5 】

$$c_j = C(P_j) = \frac{\log(aP_j + 1.0)}{\log(a - 1.0)} \quad \dots (7)$$

【 0 0 4 8 】

図 7 は、式 (7) の信頼度関数を示す図である。ここでは、横軸を同ラベル確率 (P_j)、縦軸を信頼度 (c_j) とし、 $a = 10.0$ の場合に式 (7) によって求められる点を結んだ曲線で示されている。この関数を用いる場合には、誤ラベルが多く含まれていても信頼度が大幅に下がることはなく、誤ラベルの増加に応じて信頼度が緩やかに低下する。

【 0 0 4 9 】

S 4 0 6 において、S 4 0 5 で得られた信頼度 c_j を対象の非エキスパートデータに付加し、下記の式 (8) のような形で信頼度付き非エキスパートデータ格納部 1 6 に格納する。

【 数 6 】

$$[(\vec{x}_1, y_1, c_1), \dots, (\vec{x}_j, y_j, c_j), \dots, (\vec{x}_n, y_n, c_n)] \quad \dots (8)$$

【 0 0 5 0 】

前述の 2 次元データの例 (式 (1) の非エキスパートデータ) であれば、下記の式 (9) の形で信頼度付き非エキスパートデータ格納部 1 6 に格納される。

【 数 7 】

$$[\vec{x}_j, y_j, c_j] = [(1.0, 2.0), 1, 0.96] \quad \dots (9)$$

【 0 0 5 1 】

尚、エキスパートデータの信頼度は常に 1 としているので、エキスパートデータは擬似的に下記の式 (1 0) の形でエキスパートデータ格納部 1 1 に格納されているとみなすことができる。

【 数 8 】

$$[(\vec{X}_1, Y_1, 1), \dots, (\vec{X}_j, Y_j, 1), \dots, (\vec{X}_N, Y_N, 1)] \quad \dots (10)$$

【 0 0 5 2 】

このように、近傍距離内における同ラベル確率を考慮した信頼度関数を用いることで、最近傍にあるエキスパートデータに誤ラベルが与えられていたとしても、他のラベルが正確であれば非エキスパートデータに適切な信頼度を付加することが可能になる。このような信頼度付けはデータ間の距離の長短のみに基づく信頼度付けよりもエキスパートデータの誤ラベルに対して頑健であると言える。

【 0 0 5 3 】

図 8 は、分類モデル学習部 1 7 における処理の具体例を示すフローチャートである。学習器については信頼度を反映する形のものであれば、どのような学習器でも機能すると考えられるが、ここではデータ重みに対する信頼度の組み込み易さを考慮して Ada B o o s t の手法に即した形で処理を行うものとする。尚、B a g g i n g などの他の手法を用いても良い。

【 0 0 5 4 】

10

20

30

40

50

S 8 0 1 においては、読み込まれた信頼度付き非エキスパートデータとエキスパートデータに、AdaBoostの手法に即して均等のデータ重み w_j を付ける。本発明では、AdaBoostにおける従来のデータ重み w_j に加え、信頼度決定部15で得られた信頼度 c_j が教師データに付加されているため、ここでは読み込まれた n 個の非エキスパートデータは下記の式(11)、 N 個のエキスパートデータはそれぞれ下記の式(12)の形で処理されるものとする。

【数9】

$$[(\vec{x}_1, y_1, c_1, w_1), \dots, (\vec{x}_j, y_j, c_j, w_j), \dots, (\vec{x}_n, y_n, c_n, w_n)] \dots (11)$$

10

$$[(\vec{X}_1, Y_1, 1, w_{n+1}), \dots, (\vec{X}_j, Y_j, 1, w_{n+j}), \dots, (\vec{X}_N, Y_N, 1, w_{n+N})] \dots (12)$$

【0055】

S 8 0 2 においては、非エキスパートデータに付与された信頼度 c_j をデータ重みに反映させる。ここでは、AdaBoostにおけるデータ重み w_j に対して信頼度 c_j を反映させたデータ重み w'_j を下記の式(13)により設定する。

$$w'_j = c_j w_j \dots (13)$$

【0056】

このように設定することにより、データ重み w_j が大きく学習に大きな影響を及ぼすと考えられる非エキスパートデータに関しても、その非エキスパートデータの信頼度 c_j が低ければデータ重み w'_j の値は小さくなり、非エキスパートデータに含まれる信頼度 c_j の低い教師データの影響を自然な形で小さくすることができる。

20

【0057】

S 8 0 3 においては、S 8 0 2 で得られたデータ重み w'_j を用いて弱学習器を生成する。AdaBoostに用いられる弱学習器には決定木など様々なものが考えられる。

S 8 0 4 においては、AdaBoostのアルゴリズムに従い、データ重みと弱学習器の性能に依るコスト関数の更新を行う。

【0058】

S 8 0 5 においては、終了条件を満たしているか否かを判定する。ここで、終了条件を満たすと判定された場合にはS 8 0 6へ進む。これに対し、終了条件を満たさないと判定された場合はS 8 0 2に戻る。尚、一般的なAdaBoostの手法における終了条件は、弱学習器の数が所定数を満たすことである。例えばユーザが弱学習器を100個作るという設定にすれば、S 8 0 2 ~ S 8 0 5を100回繰り返すことが終了条件である。

30

S 8 0 6 においては、生成された弱学習器を組合せることにより精度の高い分類モデルである強学習器を生成し、処理を終了する。

【0059】

このように、教師データの精度の差異という学習過程を開始する前の知識を利用して非エキスパートデータに信頼度を付与し、分類モデルの学習に組み込むことで、エキスパートデータが少ない場合であっても精度の良い分類モデルを得ることができる。

40

【0060】

図9は、予測部19における処理の具体例を示すフローチャートである。S 9 0 1 においては、分類対象データ格納部18における分類対象データの有無を判定する。ここで、分類対象データが有ると判定された場合には、S 9 0 2へ進む。これに対し、分類対象データが無いと判定された場合には、処理を終了する。

【0061】

S 9 0 2 においては、分類対象データ格納部18から分類対象データを一つ選択する。

S 9 0 3 においては、選択した分類対象データを分類モデルに当てはめてラベル付けを行い、S 9 0 1へ戻る。S 9 0 1 ~ S 9 0 3までの処理は全ての分類対象データに対してラベル付けが完了するまで繰り返し行われる。

50

【0062】

上記のように構成することで、高信頼度とされる教師データ（エキスパートデータ）の中にノイズが含まれる場合においても、同ラベル確率を入力とする信頼度関数と、各教師データの精度という事前知識を利用して非エキスパートデータに信頼度を付与し、分類モデルの学習に組み込むことで、精度の良い分類モデルを得ることができる。

【0063】

尚、本発明は上記実施形態そのままに限定されるものではなく、実施段階ではその要旨を逸脱しない範囲で構成要素を変形して具体化できる。また、上記実施形態に開示されている複数の構成要素の適宜な組み合わせにより、種々の発明を形成できる。例えば、実施形態に示される全構成要素から幾つかの構成要素を削除してもよい。更に、異なる実施形態にわたる構成要素を適宜組み合わせてもよい。

10

【0064】

また、上記実施形態においてはメールデータのようなテキストデータを例として説明したが、対象データの種類はこれに限定されない。すなわち、画像データや音声データなどのデータにおいても所定の規則に基づいて座標化することで、分類モデルを作成可能である。例えば、2次元のレントゲン画像データはM行N列に分割し、これを1行M×N列のデータに変換すればM×N次元の座標が得られる。この場合、画像データにおける色彩区分（例えば16段階のグレースケールなど）を行列の成分とすると好適である。そして、経験豊富な医師によって病変の有無が判定（ラベル付け）されたレントゲン画像データをエキスパートデータ、経験の浅い医師によって判定されたデータを非エキスパートデータとし、上記実施形態と同様に信頼度付けを行うことで精度の高い分類モデルを作成できる。

20

【図面の簡単な説明】

【0065】

【図1】本発明の一実施形態に係る分類モデル学習装置1の全体構成例を示すブロック図。

【図2】エキスパートデータおよび非エキスパートデータを2次元で具体的に説明する図。

【図3】近傍距離定義部13における処理の具体例を示すフローチャート。

【図4】信頼度決定部15における処理の具体例を示すフローチャート。

30

【図5】分類問題の評価基準に応じた信頼度関数の具体例を説明する図。

【図6】分類問題の評価基準に応じた信頼度関数の具体例を説明する図。

【図7】分類問題の評価基準に応じた信頼度関数の具体例を説明する図。

【図8】分類モデル学習部17における処理の具体例を示すフローチャート。

【図9】予測部19における処理の具体例を示すフローチャート。

【図10】エキスパートデータのラベル付けと非エキスパートデータのラベル付けに対する信頼度の関係を説明する図。

【図11】エキスパートデータのラベル付けと非エキスパートデータのラベル付けに対する信頼度の関係を説明する図。

【符号の説明】

40

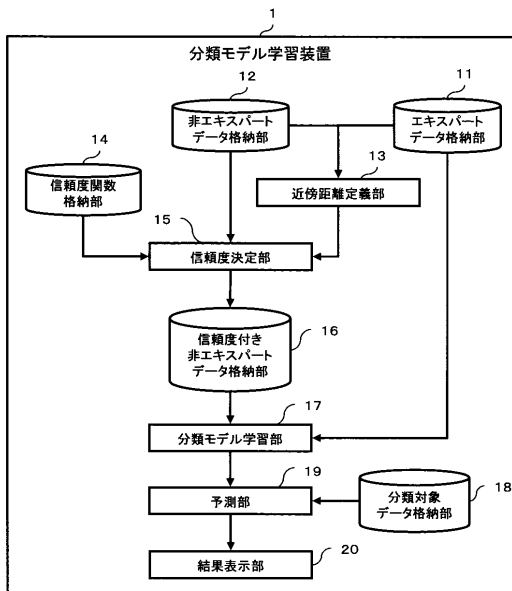
【0066】

- 1 ... 分類モデル学習装置、
- 11 ... エキスパートデータ格納部、
- 12 ... 非エキスパートデータ格納部、
- 13 ... 近傍距離定義部、
- 14 ... 信頼度関数格納部、
- 15 ... 信頼度決定部、
- 16 ... 信頼度付き非エキスパートデータ格納部、
- 17 ... 分類モデル学習部、
- 18 ... 分類対象データ格納部、

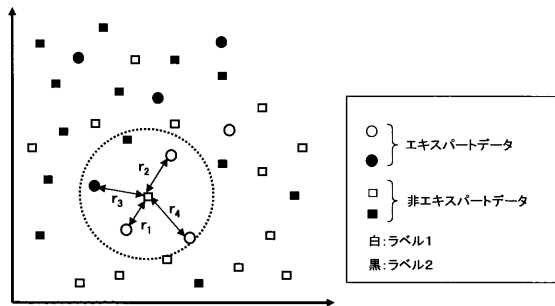
50

- 19 ... 予測部、
- 20 ... 結果表示部。

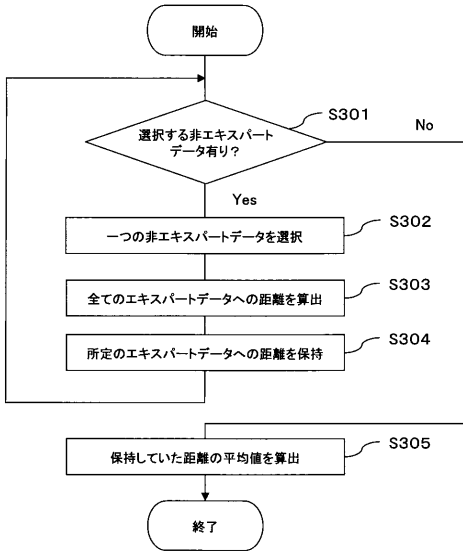
【 図 1 】



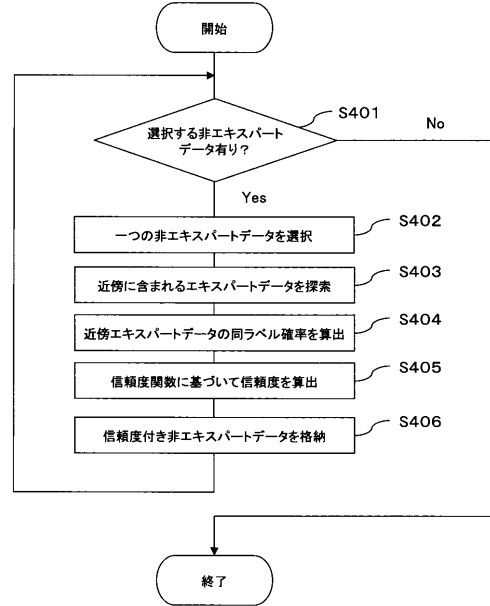
【 図 2 】



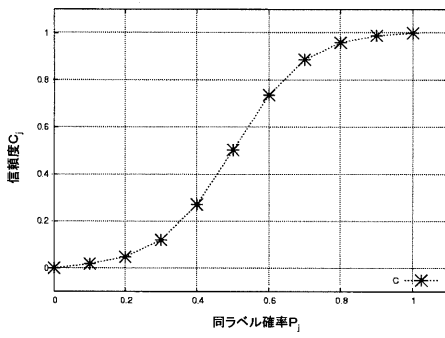
【 図 3 】



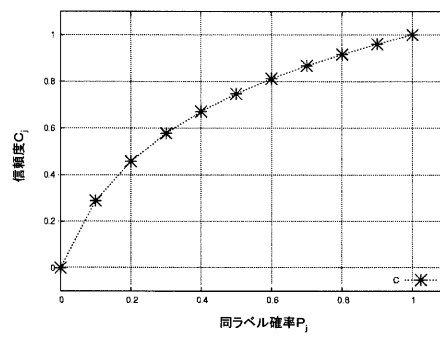
【 図 4 】



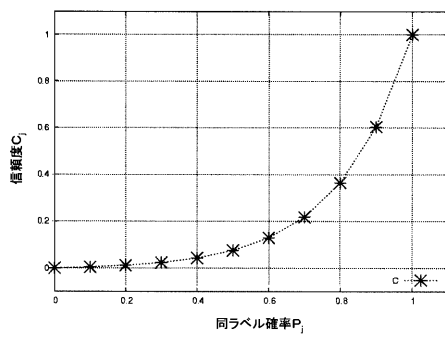
【 図 5 】



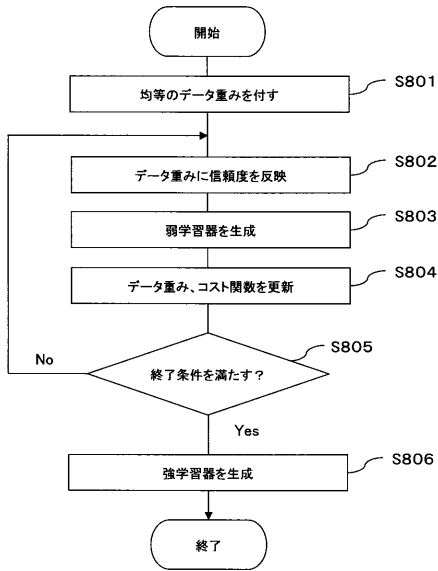
【 図 7 】



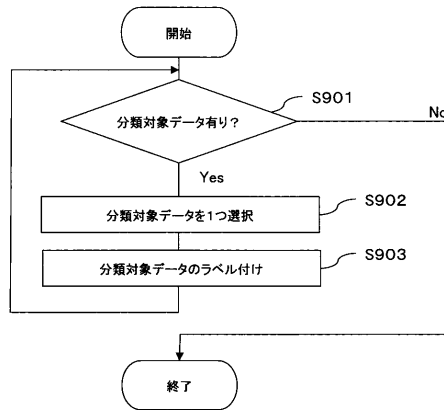
【 図 6 】



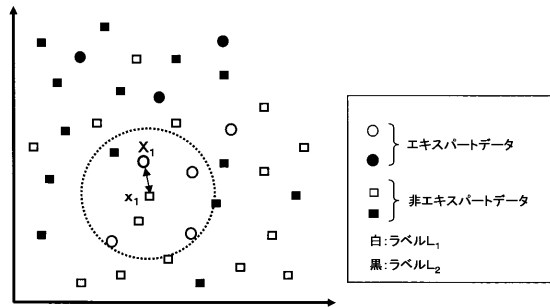
【 図 8 】



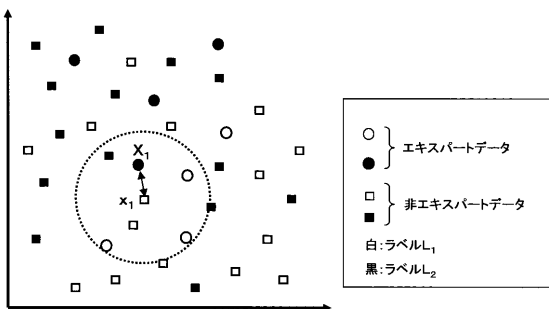
【 図 9 】



【 図 10 】



【 図 11 】



フロントページの続き

(51)Int.Cl.

F I

テーマコード(参考)

G 0 6 N 5/04 5 5 0 N