

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5398682号
(P5398682)

(45) 発行日 平成26年1月29日 (2014. 1. 29)

(24) 登録日 平成25年11月1日 (2013. 11. 1)

(51) Int. Cl.

F I

G 0 6 T 7/00 (2006. 01)

G 0 6 T 7/00 3 5 0 B

G 0 6 N 3/00 (2006. 01)

G 0 6 N 3/00 5 6 0 A

請求項の数 5 外国語出願 (全 9 頁)

(21) 出願番号 特願2010-242078 (P2010-242078)
 (22) 出願日 平成22年10月28日 (2010. 10. 28)
 (65) 公開番号 特開2011-118883 (P2011-118883A)
 (43) 公開日 平成23年6月16日 (2011. 6. 16)
 審査請求日 平成25年8月27日 (2013. 8. 27)
 (31) 優先権主張番号 12/631, 590
 (32) 優先日 平成21年12月4日 (2009. 12. 4)
 (33) 優先権主張国 米国 (US)

早期審査対象出願

(73) 特許権者 597067574
 ミツビシ・エレクトリック・リサーチ・ラ
 ボラトリーズ・インコーポレイテッド
 アメリカ合衆国、マサチューセッツ州、ケ
 ンブリッジ、ブロードウェイ 201
 201 BROADWAY, CAMBR
 IDGE, MASSACHUSETTS
 02139, U. S. A.

(74) 代理人 100110423

弁理士 曾我 道治

(74) 代理人 100094695

弁理士 鈴木 憲七

(74) 代理人 100111648

弁理士 梶並 順

最終頁に続く

(54) 【発明の名称】 局所的学習のためのトレーニング点の近傍を選択するための方法

(57) 【特許請求の範囲】

【請求項 1】

局所的学習のためのトレーニング点のセットから、或るクエリポイントの近くのトレーニング点のサブセットを選択するための方法であって、

トレーニング点のセット、クエリポイント x_q を与えるステップと、

以下の式に従って、トレーニング点のセットから、前記クエリポイント x_q の近くのトレーニング点のサブセット X_N を求めるステップとを含み、

【数 1】

$$\operatorname{argmax}_{X \subset \mathcal{X}} G(X) = D_T(X) + \lambda e^{-H(X)}$$

10

ただし、関数 argmax は関数 G を最大にする X の値を返し、 $p = 1, 2$ の場合の

【数 2】

$$D_T(X) = \sum_{x \in X} \exp(-\|x - x_q\|_p)$$

は前記クエリポイントから前記トレーニング点のサブセット X_N への累積類似度であり、 $e^{-H(X)}$ は X によって誘導される分布の逆範囲を求め、 λ は 0 より大きな制御パラメータであり、 H はシャノンエントロピーであり、

前記与えるステップ及び前記求めるステップはプロセッサにおいて実行される、選択するための方法。

【請求項 2】

20

以下の式

【数 3】

$$H(\mathcal{N}(\mu, \Sigma)) = \frac{1}{2} \log |\Sigma| + \frac{C}{2} (\log 2\pi e)$$

に従って、前記シャノンエントロピーを推定するステップをさらに含み、ただし、 μ は前記サブセット X_N 内の前記トレーニング点の平均であり、 Σ は前記サブセット X_N 内の前記トレーニング点の共分散であり、 C は前記トレーニング点の前記セットの次元数である、請求項 1 に記載の方法。

【請求項 3】

10

前記トレーニング点の前記サブセット X_N を用いて分類法をトレーニングするステップをさらに含む、請求項 1 に記載の方法。

【請求項 4】

前記トレーニング点の前記サブセット X_N を用いて回帰法をトレーニングするステップをさらに含む、請求項 1 に記載の方法。

【請求項 5】

前記回帰法として異分散サポートベクトル回帰を用いるステップをさらに含む、請求項 4 に記載の方法。

【発明の詳細な説明】

【技術分野】

20

【0001】

この発明は包括的には局所的な教師あり学習に関し、より詳細には、単一のクエリポイントに基づいて、トレーニングデータセットからトレーニング点のサブセットを選択することに関する。

【背景技術】

【0002】

教師あり学習では、クエリポイントから成る入力ベクトルに基づいて未知の出力を予測する関数を推定するために、トレーニングデータが用いられる。局所的学習法では、所与のクエリポイントに対して、そのクエリポイントの「近くにある」トレーニング点によって関数が求められる。近さは、或る距離メトリックによって判断することができる。

30

【0003】

局所的学習法の例は、最近傍回帰及び分類、並びに局所重み付け回帰を含む。2つの応用例として、過去の値に基づいて、或る時系列の未来の値を予測すること、及びピクセル値に基づいて、或る画像内に或る特定の物体が存在するか否かを検出することが挙げられる。

【0004】

そのような問題において、トレーニングデータセット D は複数のペアから成るセットであり、

【0005】

【数 1】

40

$$D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_M, y_M)\} \subset X \times \mathcal{R}$$

【0006】

である。ただし、 X は入力パターンを表し、たとえば、 $X = \mathbb{R}^d$ である。各ペアは入力ベクトル \mathbf{x}_i 及び出力 y_i を含む。関数

【0007】

【数 2】

$$\hat{y} = F(\mathbf{x})$$

【0008】

50

は対応する入力ベクトルから出力を推定し、その関数はトレーニングデータセットから学習される。

【0009】

局所的学習法において、クエリポイント x_q 毎に、入力クエリポイント x_q の近くにある、トレーニングセット内のトレーニングデータ点だけに基いて、局所関数 $F(x)$ が学習される。クエリポイントの近くにあるトレーニング点は通常、距離メトリックに従って、トレーニングデータセット内の k 個の最も近い点から選択される。代替的には、選択されるトレーニング点は、クエリポイントからの或る距離しきい値 d 未満である。

【0010】

局所的学習法の背景にある概念は、データは入力空間の異なる部分において異なる特性を有することができること、及びクエリポイントの近くにあるデータは、所与の入力から所望の出力を予測する関数を学習するのに最も有用なはずであるということである。

10

【0011】

1つの応用例において、毎日の電力需要を予測することが望まれる。1年のうちの異なる時点では、異なる要因が需要負荷に影響を及ぼす可能性がある。クエリポイントが夏日に対応する場合には、トレーニングデータセット内の夏日だけに基いて、関数 $F()$ を学習することが好都合である可能性がある。

【0012】

しかしながら、 k 個の最近傍点、又は或る距離 d 内にある全ての近傍点を用いることが、必ずしも最良の性能を与えるとは限らない。

20

【発明の概要】

【発明が解決しようとする課題】

【0013】

いずれのトレーニング点が近傍に属するかを求めるための方法と共に、局所近傍の新たな概念を提供することが望まれる。

【課題を解決するための手段】

【0014】

方法は、累積類似度を最大にすることによって、トレーニング点のセットから、1つのクエリポイントに近いトレーニング点のサブセットを選択し、その累積類似度は、クエリポイントとサブセット内の各点との類似度、及びサブセット内の点の互いの類似度を評価する。

30

【図面の簡単な説明】

【0015】

【図1A】トレーニングデータ、及びこの発明の実施の形態による方法によって選択されるトレーニングデータのサブセットのプロット図である。

【図1B】トレーニングデータ、及び従来技術の方法によって選択されるトレーニングデータのサブセットのプロット図である。

【図2A】トレーニングデータ、及びこの発明の実施の形態による方法によって選択されるトレーニングデータのサブセットのプロット図である。

【図2B】トレーニングデータ、及び従来技術の方法によって選択されるトレーニングデータのサブセットのプロット図である。

40

【図3】この発明の実施の形態による、トレーニングデータを選択するための方法のフローチャートである。

【発明を実施するための形態】

【0016】

図3に示されるように、この発明の実施の形態は、単一の入力クエリポイント x_q 305の近くにある、複数の点から成る局所近傍 X_N 302を選択するための方法を提供する。近傍点の全てがコンパクトであることが望ましい。この明細書において定義されるように、複数の点が或る距離メトリックに従って互いに相対的に近い場合には、そのセットはコンパクトである。この方法のステップは、当該技術分野において既知のように、メモリ

50

及び入力／出力インターフェースを備えるプロセッサにおいて実行することができる。

【 0 0 1 7 】

従来の最近傍法は、結果として得られる近傍がコンパクトであるか否かを考慮することなく、クエリポイントの近くにある近傍点を選択する。

【 0 0 1 8 】

この発明の実施の形態による方法は、入力トレーニングデータ(3 0 1 参照)が入力空間内で不均一に分布するときに性能を改善するために、コンパクトであることの判定基準を含む。

【 0 0 1 9 】

この発明による局所近傍点のサブセット X_N 3 0 2 は(3 1 0 参照)以下の通りである。

【 0 0 2 0 】

【数 3】

$$X_N = \operatorname{argmax}_{X \subset \mathcal{X}} G(X) = D_T(X) + \lambda e^{-H(X)}, \lambda > 0 \quad (1)$$

【 0 0 2 1 】

ただし、関数 argmax 3 1 0 は関数 G を最大にするパラメータ X の値を返し、 $p = 1, 2$ の場合に、

【 0 0 2 2 】

【数 4】

$$D_T(X) = \sum_{x \in X} \exp(-\|x - x_q\|_p) \quad 20$$

【 0 0 2 3 】

はクエリポイント 3 0 5 からトレーニングサブセット

【 0 0 2 4 】

【数 5】

$$X \subset \mathcal{X}$$

【 0 0 2 5 】

への累積類似度であり、 $e^{-H(X)}$ は X によって誘導される分布の逆範囲(inverse range)を求め、 λ は制御パラメータである。分布の「範囲」によって、この発明では、サンプル空間が離散的である場合には、サンプル空間内の点の数を意味し、サンプル空間が指数関数的なエントロピーの特性に従う実線である場合には、確率密度関数が 0 でない区間の長さを意味する。 H はシャノンエントロピーである。シャノンエントロピーは、確率変数の値が未知であるときに失う平均情報量の指標である。

【 0 0 2 6 】

この発明では、以下のようにガウス分布を仮定してシャノンエントロピーを推定する。

【 0 0 2 7 】

【数 6】

$$H(\mathcal{N}(\mu, \Sigma)) = \frac{1}{2} \log |\Sigma| + \frac{C}{2} (\log 2\pi e) \quad (2) \quad 40$$

【 0 0 2 8 】

ただし、 μ はサブセット内の点の平均であり、 Σ はその共分散である。 C は入力トレーニング点 3 0 1 の次元数である。

【 0 0 2 9 】

式(1)の目的は、或るクエリポイントの近くの最も大きなクラスタへの累積類似度を最大にすることである。この発明の目的は、トレーニングデータパターンの基本的な構造に適応できるような方法で、トレーニングデータのサブセット 3 0 2 を見つけることである。

【 0 0 3 0 】

10

20

30

40

50

この問題の組み合わせ最適化特性は、従来の最近傍法において用いられる貪欲法との重要な相違点である。この発明において最大にする、式(1)において定義される目的関数は、優モジュラ性として知られている数学的特性を有する。全ての $x, y \in \mathbb{R}^k$ に対して、

【0031】

【数7】

$$f(x \vee y) + f(x \wedge y) \geq f(x) + f(y)$$

【0032】

である場合には、関数

【0033】

$$f: \mathbb{R}^k \rightarrow \mathbb{R}$$

【0034】

は優モジュラ性である。ただし、

【0035】

【数8】

$$x \vee y$$

【0036】

は x 及び y の成分毎の最大値を表し、

【0037】

【数9】

$$x \wedge y$$

【0038】

は成分毎の最小値を表す。

【0039】

優モジュラ関数を最大にすることは、劣モジュラ関数を最小にすることと同じである。それゆえ、この発明では、この関数を最適化するために、劣モジュラ性最適化の従来の手段を適用することができる。

【0040】

上記の手順を用いて、複数の点から成る最適なサブセットが求められた後に、任意の分類又は回帰法をトレーニングするために(320)、その点のサブセットを用いることができる。

【0041】

図1A、1B、2A、2Bを用いて、2つの単純な例を説明し、 $k=15$ の場合の、この発明の劣モジュラモデル選択法と、従来の k 最近傍点選択法とを比較する。それらの図は、単一のクエリポイント101と、サブセット102とを示す。これらの図において、時間間隔 $[-8, 8]$ 内に100個の合成入力データ点 x_i が均一な間隔で配置されており、対象となる出力データ点 $y_i = 2 \sin(2x_i)$ が、 $x=0$ における0.5から $x=1$ における1.5まで線形に増加する標準偏差を有するガウス雑音で汚染される。

【0042】

これらの図は、この発明の方法が従来の方法よりも優れていることを明示する。劣モジュラ法では、近傍点が適応的に選択され、ここで図1A及び図1Bにおいて示されるように分布の先端では点の数が少なく、図2A及び図2Bにおいて示されるように、末端に近いほど数が増え、それは、この発明の例においてガウス雑音の汚染が増加することと一致する。

【0043】

異分散サポートベクトル回帰

上記のように選択される近傍トレーニングデータを用いて、任意の回帰又は分類法をトレーニングすることができる。ここで、異分散サポートベクトル回帰のための1つのその

10

20

30

40

50

ような技法を説明する。異分散サポートベクトル回帰は、サポートベクトル回帰を拡張したものであり、局所近傍を用いて局所回帰関数を見つける。統計学において、一連の確率変数が異なる分散を有するときに、それらの確率変数は異分散である。

【 0 0 4 4 】

異分散サポートベクトル回帰は、形式

【 0 0 4 5 】

$$F(x) = w^T x + b$$

【 0 0 4 6 】

の関数 $F(x)$ を推定する。ただし、 w^T は転置演算子 T を伴うベクトルであり、 b はスカラーである。

【 0 0 4 7 】

その関数は、以下の最適化問題を解くことによって求められる。

【 0 0 4 8 】

【数 1 0】

$$\min_{w, b, \epsilon, \xi, \xi^*} \frac{1}{2N} \sum_{i=1}^N w^T (2NI + \Sigma_i) w + C \sum_{i=1}^N (\xi_i + \xi_i^* + \epsilon)^p$$

$$\text{s.t. } y_i - (w^T x_i + b) \leq \epsilon + \xi_i^*$$

$$(w^T x_i + b) - y_i \leq \epsilon + \xi_i$$

$$\xi_i^*, \xi_i \geq 0, \forall i,$$

(3)

10

20

【 0 0 4 9 】

ただし、 I は $N \times N$ 恒等行列であり、 N は入力ベクトルの次元数であり、 ξ_i 及び ξ_i^* はスラック変数であり、 ϵ は誤差許容範囲であり、 $p \in \{1, 2\}$ はペナルティタイプを決定し、そして

【 0 0 5 0 】

【数 1 1】

$$\Sigma_i = \frac{1}{k_i + 1} \sum_{x \in X_i} (x - \bar{x}_i)(x - \bar{x}_i)^T$$

【 0 0 5 1 】

は x_i の近傍内のトレーニング点のための実験的な共分散である。ただし、 X_i は近傍点のサブセットであり、 k_i は X_i 内の点の数であり、

【 0 0 5 2 】

【数 1 2】

$$\bar{x}_i$$

【 0 0 5 3 】

はこれらの近傍点の平均である。複数の点から成るこの近傍は、上記の劣モジュラ性最適化技法を用いて選択することができる。

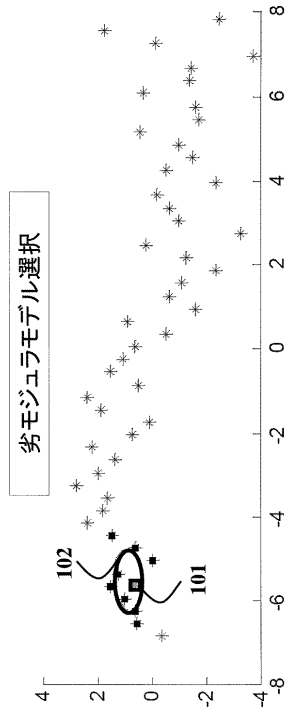
【 0 0 5 4 】

この発明を、好ましい実施の形態の例として説明してきたが、この発明の精神及び範囲内で他のさまざまな適合及び変更を行えることが理解されるべきである。したがって、この発明の真の精神及び範囲内に入るすべての変形及び変更を包含することが、添付の特許請求の範囲の目的である。

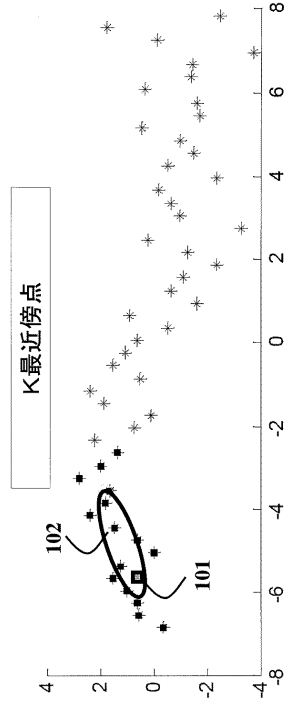
40

30

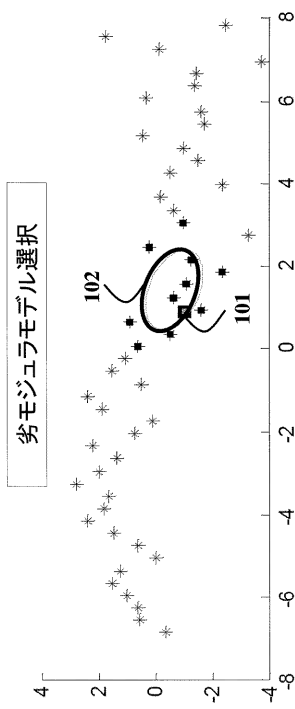
【図 1 A】



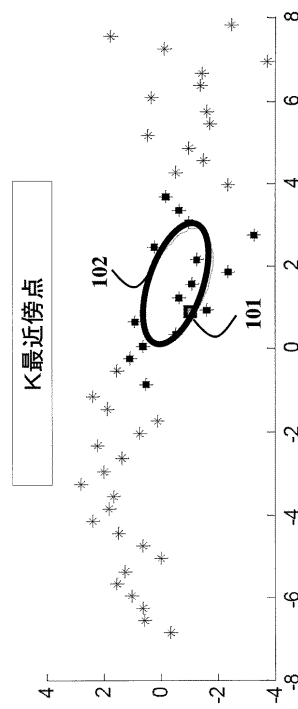
【図 1 B】



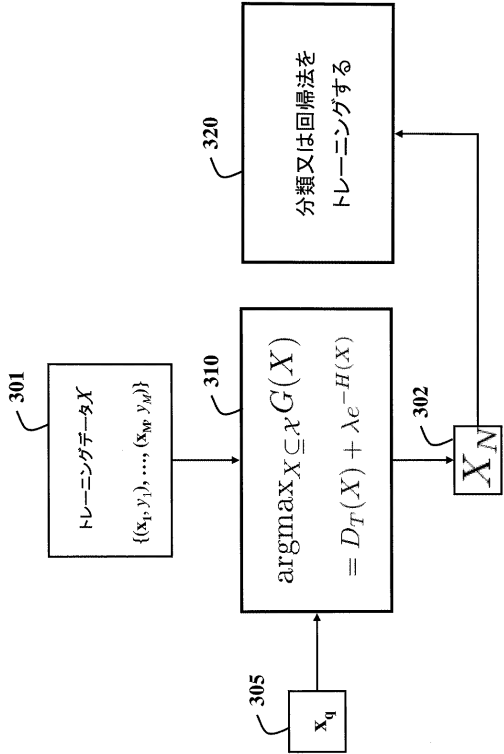
【図 2 A】



【図 2 B】



【図 3】



フロントページの続き

(74)代理人 100122437

弁理士 大宅 一宏

(74)代理人 100147566

弁理士 上田 俊一

(74)代理人 100161171

弁理士 吉田 潤一郎

(74)代理人 100161115

弁理士 飯野 智史

(72)発明者 ケヴィン・ダブリュ・ウィルソン

アメリカ合衆国、マサチューセッツ州、ケンブリッジ、ウィンター・ストリート 80、ユニット
1

(72)発明者 シャオハン・ジアン

アメリカ合衆国、マサチューセッツ州、ケンブリッジ、ヴァッサー・ストリート 32

審査官 佐藤 実

(56)参考文献 特開平8-106295(JP,A)

特開2009-259109(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06T 7/00

G06N 3/00