

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7231829号
(P7231829)

(45)発行日 令和5年3月2日(2023.3.2)

(24)登録日 令和5年2月21日(2023.2.21)

(51)国際特許分類 F I
G 0 6 N 20/00 (2019.01) G 0 6 N 20/00

請求項の数 8 (全32頁)

(21)出願番号	特願2019-137027(P2019-137027)	(73)特許権者	000005223 富士通株式会社
(22)出願日	令和1年7月25日(2019.7.25)		神奈川県川崎市中原区上小田中4丁目1 番1号
(65)公開番号	特開2021-22051(P2021-22051A)	(74)代理人	110002918 弁理士法人扶桑国際特許事務所
(43)公開日	令和3年2月18日(2021.2.18)	(72)発明者	小林 健一 神奈川県川崎市中原区上小田中4丁目1 番1号 富士通株式会社内
審査請求日	令和4年4月7日(2022.4.7)	審査官	松平 英

最終頁に続く

(54)【発明の名称】 機械学習プログラム、機械学習方法および機械学習装置

(57)【特許請求の範囲】

【請求項1】

コンピュータに、

データ集合から抽出された複数の第1の訓練データを用いて、機械学習により前記複数の第1の訓練データに対応する複数の第1のモデルを学習し、

前記データ集合から抽出された第1のテストデータに含まれる2以上のレコードそれぞれを前記複数の第1のモデルに入力することで、前記複数の第1のモデルと前記2以上のレコードとの組み合わせ毎に算出された予測誤差を示す誤差情報を生成し、

前記誤差情報に基づいて、テストデータのサイズとテストデータを用いて算出されるモデルの精度の測定値が有する分散との間の対応関係を判定し、

前記データ集合から抽出された第2の訓練データを用いて学習された第2のモデルの精度を、前記データ集合から抽出される第2のテストデータを用いて測定する場合に、前記対応関係に基づいて、前記第2のモデルに対して算出される精度の測定値の分散が所定条件を満たすように前記第2のテストデータのサイズを決定する、

処理を実行させる機械学習プログラム。

【請求項2】

前記対応関係は、テストデータのサイズの増加に応じて分散が下限に漸近するように減少する非線形関係であり、前記第2のテストデータのサイズは、サイズの所定増加量に対する分散の減少度を示す効率性指標に基づいて決定される、

請求項1記載の機械学習プログラム。

10

20

【請求項 3】

前記所定条件は、前記効率性指標の値が閾値以上であることであり、前記第 2 のテストデータのサイズは、前記所定条件を満たす範囲で最大のサイズに決定される、

請求項 2 記載の機械学習プログラム。

【請求項 4】

前記対応関係の判定では、前記 2 以上のレコードそれぞれについて前記複数の第 1 のモデルに対して算出された予測誤差を平均化した予測バイアスを算出し、前記 2 以上のレコードの前記予測バイアスを合成して、前記対応関係を表すパラメータの値を決定する、

請求項 1 記載の機械学習プログラム。

【請求項 5】

前記対応関係の判定では、訓練データのサイズに依存しない第 1 のパラメータと訓練データのサイズに依存する第 2 のパラメータとテストデータのサイズを示す第 3 のパラメータとを用いて分散を算出する分散関数に対して、前記第 1 のパラメータの値を推定し、

前記第 2 のテストデータのサイズの決定では、前記第 2 のモデルの学習結果に基づいて前記第 2 のパラメータの値を推定し、前記第 3 のパラメータの値を変動させることで、分散が前記所定条件を満たすテストデータのサイズを探索する、

請求項 1 記載の機械学習プログラム。

【請求項 6】

前記第 2 のテストデータのサイズの決定では、前記第 1 のテストデータを前記第 2 のモデルに入力して算出される予測誤差に基づいて前記第 2 のパラメータの値を仮選択し、前記仮選択した第 2 のパラメータの値を用いてテストデータのサイズを仮選択し、前記データ集合から抽出された前記仮選択したサイズのテストデータを前記第 2 のモデルに入力して算出される予測誤差に基づいて前記第 2 のパラメータの値を決定する、

請求項 5 記載の機械学習プログラム。

【請求項 7】

コンピュータが、

データ集合から抽出された複数の第 1 の訓練データを用いて、機械学習により前記複数の第 1 の訓練データに対応する複数の第 1 のモデルを学習し、

前記データ集合から抽出された第 1 のテストデータに含まれる 2 以上のレコードそれぞれを前記複数の第 1 のモデルに入力することで、前記複数の第 1 のモデルと前記 2 以上のレコードとの組み合わせ毎に算出された予測誤差を示す誤差情報を生成し、

前記誤差情報に基づいて、テストデータのサイズとテストデータを用いて算出されるモデルの精度の測定値が有する分散との間の対応関係を判定し、

前記データ集合から抽出された第 2 の訓練データを用いて学習された第 2 のモデルの精度を、前記データ集合から抽出される第 2 のテストデータを用いて測定する場合に、前記対応関係に基づいて、前記第 2 のモデルに対して算出される精度の測定値の分散が所定条件を満たすように前記第 2 のテストデータのサイズを決定する、

機械学習方法。

【請求項 8】

データ集合を記憶する記憶部と、

前記データ集合から抽出された複数の第 1 の訓練データを用いて、機械学習により前記複数の第 1 の訓練データに対応する複数の第 1 のモデルを学習し、前記データ集合から抽出された第 1 のテストデータに含まれる 2 以上のレコードそれぞれを前記複数の第 1 のモデルに入力することで、前記複数の第 1 のモデルと前記 2 以上のレコードとの組み合わせ毎に算出された予測誤差を示す誤差情報を生成し、前記誤差情報に基づいて、テストデータのサイズとテストデータを用いて算出されるモデルの精度の測定値が有する分散との間の対応関係を判定し、前記データ集合から抽出された第 2 の訓練データを用いて学習された第 2 のモデルの精度を、前記データ集合から抽出される第 2 のテストデータを用いて測定する場合に、前記対応関係に基づいて、前記第 2 のモデルに対して算出される精度の測定値の分散が所定条件を満たすように前記第 2 のテストデータのサイズを決定する処理部

10

20

30

40

50

と、

を有する機械学習装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は機械学習プログラム、機械学習方法および機械学習装置に関する。

【背景技術】

【0002】

コンピュータを利用したデータ分析の1つとして、機械学習が行われることがある。機械学習では、幾つかの既知の事例を示す訓練データをコンピュータに入力する。コンピュータは、訓練データを分析して、要因（説明変数や独立変数と言うことがある）と結果（目的変数や従属変数と言うことがある）との間の関係を一般化したモデルを学習する。学習されたモデルを用いることで、未知の事例についての結果を予測することができる。

10

【0003】

機械学習では、学習されるモデルの正確さ、すなわち、未知の事例の結果を正確に予測する能力（モデルの精度、予測性能、性能などと言うことがある）が問題となる。モデルの精度は、分析対象とする事象の性質、モデルの学習に使用する訓練データのサイズ、機械学習アルゴリズムなど様々な側面に依存する。精度が不十分なモデルが生成された場合、訓練データのサイズを増加させる、機械学習アルゴリズムを変更するなどの調整を行って、機械学習を再実行することもある。そこで、データ母集合から訓練データとは異なる既知の事例を示すテストデータを抽出し、訓練データを用いて学習されたモデルに対してテストデータを入力することで、モデルの精度を評価することが多い。

20

【0004】

なお、同一のデータ母集合から、異なる分割方法によって訓練データとテストデータのペアを複数通り生成し、ペア毎に訓練データを用いたモデルの学習とテストデータを用いた精度の評価を行い、精度の平均を算出する汎化能力評価方法が提案されている。また、データベースから訓練データを抽出して回帰分析を行い、回帰モデルの精度を評価し、精度が不十分である場合には訓練データを追加して回帰分析を再度行う結果予測装置が提案されている。また、教師ラベルが付されたサンプルのデータ母集合から、訓練データと類似するテストデータを抽出し、訓練データを用いて学習された分類モデルの精度を、訓練データと類似するテストデータを用いて評価する情報処理システムが提案されている。

30

【先行技術文献】

【特許文献】

【0005】

【文献】特開平9 - 54764号公報

特開2014 - 13560号公報

国際公開第2017 / 183548号

【発明の概要】

【発明が解決しようとする課題】

【0006】

しかし、モデルの精度を評価するにあたり、テストデータのサイズをどの様に決定すればよいか問題となる。テストデータが少な過ぎると、テストデータとして選択されるサンプルの偶然性の影響を強く受けて、算出される精度が不正確になり信頼性が低下する。一方、テストデータが多過ぎると、精度の評価に長時間かかることになり非効率である。この点、従来の機械学習では、訓練データのサイズの2分の1から4分の1程度をテストデータのサイズとするなど、経験則に基づいてサイズを決定していた。そのため、テストデータを用いたモデルの精度の評価について改善の余地があった。

40

【0007】

1つの側面では、本発明は、機械学習のテストデータのサイズを適切に決定できる機械学習プログラム、機械学習方法および機械学習装置を提供することを目的とする。

50

【課題を解決するための手段】

【0008】

1つの態様では、コンピュータに以下の処理を実行させる機械学習プログラムが提供される。データ集合から抽出された複数の第1の訓練データを用いて、機械学習により複数の第1の訓練データに対応する複数の第1のモデルを学習する。データ集合から抽出された第1のテストデータに含まれる2以上のレコードそれぞれを複数の第1のモデルに入力することで、複数の第1のモデルと2以上のレコードとの組み合わせ毎に算出された予測誤差を示す誤差情報を生成する。誤差情報に基づいて、テストデータのサイズとテストデータを用いて算出されるモデルの精度の測定値が有する分散との間の対応関係を判定する。データ集合から抽出された第2の訓練データを用いて学習された第2のモデルの精度を、データ集合から抽出される第2のテストデータを用いて測定する場合に、対応関係に基づいて、第2のモデルに対して算出される精度の測定値の分散が所定条件を満たすように第2のテストデータのサイズを決定する。

10

【0009】

また、1つの態様では、コンピュータが実行する機械学習方法が提供される。また、1つの態様では、記憶部と処理部とを有する機械学習装置が提供される。

【発明の効果】

【0010】

1つの側面では、機械学習のテストデータのサイズが適切に決定される。

【図面の簡単な説明】

20

【0011】

【図1】第1の実施の形態の機械学習装置の例を説明する図である。

【図2】第2の実施の形態の機械学習装置のハードウェア例を示す図である。

【図3】訓練データサイズと予測性能の関係例を示すグラフである。

【図4】予測性能の測定値の分散例を示すグラフである。

【図5】予測性能の期待ロスおよび期待バイアスの例を示すグラフである。

【図6】機械学習装置の機能例を示すブロック図である。

【図7】誤差プロファイルテーブルの例を示す図である。

【図8】分散関数テーブルの例を示す図である。

【図9】機械学習の手順例を示すフローチャートである。

30

【図10】機械学習の手順例を示すフローチャート(続き)である。

【発明を実施するための形態】

【0012】

以下、本実施の形態を図面を参照して説明する。

[第1の実施の形態]

第1の実施の形態を説明する。

【0013】

図1は、第1の実施の形態の機械学習装置の例を説明する図である。

第1の実施の形態の機械学習装置10は、訓練データを用いて機械学習によりモデルを生成し、テストデータを用いてモデルの精度を測定する。機械学習装置10を、情報処理装置やコンピュータとすることもできる。機械学習装置10は、ユーザが操作するクライアント装置でもよいし、他の装置からアクセスされるサーバ装置でもよい。

40

【0014】

機械学習装置10は、記憶部11および処理部12を有する。記憶部11は、RAM(Random Access Memory)などの揮発性半導体メモリでもよいし、HDD(Hard Disk Drive)やフラッシュメモリなどの不揮発性ストレージでもよい。処理部12は、例えば、CPU(Central Processing Unit)、GPU(Graphics Processing Unit)、DSP(Digital Signal Processor)などのプロセッサである。ただし、処理部12は、ASIC(Application Specific Integrated Circuit)やFPGA(Field Programmable Gate Array)などの特定用途の電子回路を含んでもよい。プロセッサは、RAMなどのメ

50

メモリ（記憶部 11 でもよい）に記憶されたプログラムを実行する。複数のプロセッサの集合を「マルチプロセッサ」または単に「プロセッサ」と言うこともある。

【0015】

記憶部 11 は、データ集合 13 を記憶する。データ集合 13 は、既知の事例を示す複数のレコードを含む。レコードを、サンプルや行やデータブロックと言うこともできる。各レコードは、1 以上の説明変数の値と 1 つの目的変数の値とを含む。説明変数を、カラムと言うこともできる。目的変数の値は、ユーザから与えられる正解であり、教師ラベルと言うこともできる。データ集合 13 は、100 万レコード以上の多数のレコードを含んでもよく、ビッグデータと言われる大規模データであってもよい。

【0016】

処理部 12 は、データ集合 13 から、訓練データ 14a, 14b, 14c を含む複数セットの訓練データ（第 1 の訓練データ）を抽出する。ここで抽出する訓練データのセット数は少数でよく、各訓練データのサイズは小さくてよい。例えば、訓練データのセット数を 10 セット程度とし、各訓練データのサイズを 1 万レコード程度とする。各訓練データのサイズは、後述する訓練データ 18 のサイズの 100 分の 1 程度でよい。処理部 12 は、同一のレコードが異なる訓練データに含まれないようにデータ集合 13 からレコードを抽出してもよいし、同一のレコードが異なる訓練データに含まれることを許容してもよい。処理部 12 は、データ集合 13 からランダムにレコードを抽出してもよい。

【0017】

また、処理部 12 は、データ集合 13 からテストデータ 15（第 1 のテストデータ）を抽出する。テストデータ 15 のサイズは、訓練データ 14a, 14b, 14c のサイズより小さくてよい。テストデータ 15 のサイズは、訓練データ 14a, 14b, 14c のサイズの 2 分の 1 から 4 分の 1 程度でもよく、例えば、5000 レコード程度とする。処理部 12 は、テストデータ 15 に属するレコードを、訓練データ 14a, 14b, 14c と重複しないようにデータ集合 13 から抽出することが好ましい。

【0018】

処理部 12 は、訓練データ 14a, 14b, 14c を含む複数セットの訓練データを用いて、機械学習によりそれら複数セットの訓練データに対応する複数のモデルを学習する。訓練データ 14a から 1 つのモデルが学習され、それと独立に訓練データ 14b から 1 つのモデルが学習され、それと独立に訓練データ 14c から 1 つのモデルが学習される。

【0019】

複数のモデルの学習には、同一の機械学習アルゴリズムが使用される。使用する機械学習アルゴリズムは、ユーザにより指定されてもよい。機械学習アルゴリズムとして、回帰分析、サポートベクタマシン、ランダムフォレストなどが挙げられる。モデルは、説明変数と目的変数との間の関係を示し、通常、1 以上の説明変数と 1 以上の係数と 1 つの目的変数とを含む。係数は、機械学習を通じて訓練データに基づいて決定される。

【0020】

次に、処理部 12 は、テストデータ 15 および学習した複数のモデルを用いて、誤差情報 16 を生成する。誤差情報 16 を、誤差プロファイルと言うこともできる。誤差情報 16 は、複数セットの訓練データに対応する複数のモデルとテストデータ 15 に含まれる 2 以上のレコードとの組み合わせ毎に算出された予測誤差を示す。

【0021】

このとき、処理部 12 は、ある訓練データから学習された 1 つのモデルに、テストデータ 15 に含まれる 1 つのレコードを入力することで、当該 1 つのモデルと当該 1 つのレコードの組に対応する 1 つの予測誤差を算出する。例えば、処理部 12 は、テストデータ 15 のレコードに含まれる説明変数の値をモデルの説明変数に代入する。処理部 12 は、モデルによって算出される目的変数の値である予測値と、テストデータ 15 のレコードに含まれる目的変数の値である正解値とを比較し、両者の差を予測誤差として算出する。

【0022】

次に、処理部 12 は、誤差情報 16 に基づいて対応関係 17 を判定する。対応関係 17

10

20

30

40

50

は、テストデータのサイズと、テストデータを用いて算出されるモデルの精度の測定値が有するばらつきの程度である分散との間の対応関係を示す。モデルの精度は、未知の事例の結果を正確に予測する能力であり、予測性能や性能と言うこともできる。モデルの精度の指標として、正答率 (Accuracy)、適合率 (Precision)、平均二乗誤差 (MSE)、二乗平均平方根誤差 (RMSE) などが挙げられる。

【0023】

対応関係 17 は、例えば、テストデータのサイズの増加に応じて分散が下限に漸近するように減少する非線形関係である。一般に、データ集合 13 からのテストデータの抽出には、レコードの選択の偶然性がある。このため、テストデータのサイズが小さいと、レコードの選択の偶然性の影響を強く受けて、精度の測定値が真の値からずれるリスクが高くなる。テストデータのサイズを大きくすることで、分散を小さくすることができる。ただし、データ集合 13 からの訓練データの抽出にも、レコードの選択の偶然性がある。テストデータのサイズの増加だけでは、精度の測定値の分散は 0 にならない。

10

【0024】

対応関係 17 は、機械学習に使用するデータ集合 13 や機械学習アルゴリズムに依存し得る。そこで、処理部 12 は、誤差情報 16 に基づいて対応関係 17 を判定する。例えば、処理部 12 は、誤差情報 16 が示す予測誤差のうち、テストデータ 15 のレコードが同一でモデルが異なる予測誤差を平均化することで、テストデータ 15 のレコード毎に予測バイアスを算出する。処理部 12 は、テストデータ 15 の 2 以上のレコードの予測バイアスを合成して、対応関係 17 を規定するパラメータの値を決定する。

20

【0025】

対応関係 17 は、訓練データのサイズに依存しない第 1 のパラメータと、訓練データのサイズに依存する第 2 のパラメータと、テストデータのサイズを示す第 3 のパラメータとから分散を算出する分散関数であってもよい。この場合、処理部 12 は、誤差情報 16 を用いて第 1 のパラメータの値を推定してもよい。これにより、分散関数は、変数として第 2 のパラメータと第 3 のパラメータをもつ関数になる。

【0026】

次に、処理部 12 は、データ集合 13 から訓練データ 18 (第 2 の訓練データ) を抽出する。訓練データ 18 のサイズは、訓練データ 14a, 14b, 14c より十分に大きくてもよく、ユーザから指定されてもよい。例えば、訓練データ 18 のサイズを 100 万レコード程度とする。処理部 12 は、訓練データ 18 を用いてモデルを学習する。

30

【0027】

モデルが学習されると、処理部 12 は、データ集合 13 からテストデータ 19 (第 2 のテストデータ) を抽出する。処理部 12 は、テストデータ 19 に属するレコードを、訓練データ 18 と重複しないようにデータ集合 13 から抽出することが好ましい。処理部 12 は、訓練データ 18 から学習されたモデルの精度を、テストデータ 19 を用いて測定する。例えば、処理部 12 は、テストデータ 19 のレコードに含まれる説明変数の値をモデルの説明変数に代入し、モデルによって算出される目的変数の予測値とテストデータ 19 のレコードに含まれる目的変数の正解値とを比較して、精度を測定する。

【0028】

このとき、処理部 12 は、対応関係 17 に基づいて、モデルの精度の測定値の分散が所定条件を満たすように、テストデータ 19 のサイズを決定する。例えば、処理部 12 は、対応関係 17 において、サイズの所定増加量に対する分散の減少度を示す効率性指標を算出し、効率性指標に基づいてテストデータ 19 のサイズを決定する。対応関係 17 が、テストデータのサイズの増加に応じて分散が下限に漸近する非線形関係である場合、効率性指標の値は、テストデータのサイズの増加に応じて減少する。テストデータ 19 のサイズは、効率性指標の値が閾値以上である範囲で最大のサイズとしてもよい。

40

【0029】

また、例えば、処理部 12 は、訓練データ 18 を用いたモデルの学習結果に基づいて、分散関数に含まれる訓練データのサイズに依存する第 2 のパラメータの値を決定する。そ

50

して、処理部 12 は、決定された上記の第 1 のパラメータの値および第 2 のパラメータの値のもとで、テストデータのサイズを示す第 3 のパラメータの値を変動させることで、分散が所定条件を満たすテストデータのサイズを探索する。

【0030】

なお、誤差情報 16 の生成および対応関係 17 の判定は、訓練データ 18 を用いた機械学習の前に行ってもよいし後に行ってもよい。処理部 12 は、訓練データ 18 を用いて学習されたモデルと、テストデータ 19 を用いて測定された精度を出力する。処理部 12 は、学習されたモデルと測定された精度を、記憶装置に保存してもよいし、表示装置に表示してもよいし、他の情報処理装置に送信してもよい。

【0031】

第 1 の実施の形態の機械学習装置 10 によれば、小さいサイズの訓練データ 14 a , 14 b , 14 c を用いて複数のモデルが学習される。小さいサイズのテストデータ 15 を用いて、それら複数のモデルとテストデータ 15 の 2 以上のレコードとの組み合わせ毎に算出された予測誤差を示す誤差情報 16 が生成される。誤差情報 16 に基づいて、テストデータのサイズとモデルの精度の測定値が有する分散との間の対応関係 17 が判定される。そして、訓練データ 18 を用いて学習されたモデルの精度を、テストデータ 19 を用いて測定するにあたり、対応関係 17 に基づいて、精度の測定値の分散が所定条件を満たすようにテストデータ 19 のサイズが決定される。

【0032】

これにより、テストデータ 19 のサイズがモデル精度の測定値の分散に与える影響を考慮して、テストデータ 19 のサイズを適切に決定することができる。よって、テストデータ 19 のサイズが小さ過ぎることによる測定値の信頼性の低下を抑制できる。また、テストデータ 19 のサイズが大き過ぎることによる処理時間の増大を抑制できる。このため、機械学習により学習されたモデルの精度を、高信頼かつ短時間で測定することができ、モデルの精度の測定を効率化できる。特に、テストデータのサイズを訓練データのサイズの 2 分の 1 から 4 分の 1 程度とする経験則と比べて、テストデータのサイズを削減できる。

【0033】

[第 2 の実施の形態]

次に、第 2 の実施の形態を説明する。

図 2 は、第 2 の実施の形態の機械学習装置のハードウェア例を示す図である。

【0034】

機械学習装置 100 は、CPU 101、RAM 102、HDD 103、画像インタフェース 104、入力インタフェース 105、媒体リーダ 106 および通信インタフェース 107 を有する。機械学習装置 100 が有するこれらのユニットは、機械学習装置 100 の内部でバスに接続されている。機械学習装置 100 は、第 1 の実施の形態の機械学習装置 10 に対応する。CPU 101 は、第 1 の実施の形態の処理部 12 に対応する。RAM 102 または HDD 103 は、第 1 の実施の形態の記憶部 11 に対応する。

【0035】

CPU 101 は、プログラムの命令を実行するプロセッサである。CPU 101 は、HDD 103 に記憶されたプログラムやデータの少なくとも一部を RAM 102 にロードし、プログラムを実行する。CPU 101 は複数のプロセッサコアを備えてもよく、機械学習装置 100 は複数のプロセッサを備えてもよい。複数のプロセッサの集合を「マルチプロセッサ」または単に「プロセッサ」と言うことがある。

【0036】

RAM 102 は、CPU 101 が実行するプログラムや CPU 101 が演算に使用するデータを一時的に記憶する揮発性半導体メモリである。機械学習装置 100 は、RAM 以外の種類のメモリを備えてもよく、複数のメモリを備えてもよい。

【0037】

HDD 103 は、OS (Operating System) やミドルウェアやアプリケーションソフトウェアなどのソフトウェアのプログラム、および、データを記憶する不揮発性ストレージ

10

20

30

40

50

ジである。機械学習装置 100 は、フラッシュメモリや SSD (Solid State Drive) など他の種類のストレージを備えてもよく、複数のストレージを備えてもよい。

【0038】

画像インタフェース 104 は、CPU 101 からの命令に従って、機械学習装置 100 に接続された表示装置 111 に画像を出力する。表示装置 111 として、CRT (Cathode Ray Tube) ディスプレイ、液晶ディスプレイ (LCD: Liquid Crystal Display)、有機 EL (OEL: Organic Electro-Luminescence) ディスプレイ、プロジェクタなど、任意の種類の表示装置を使用することができる。機械学習装置 100 に、プリンタなど表示装置 111 以外の出力デバイスが接続されてもよい。

【0039】

入力インタフェース 105 は、機械学習装置 100 に接続された入力デバイス 112 から入力信号を受け付ける。入力デバイス 112 として、マウス、タッチパネル、タッチパッド、キーボードなど、任意の種類の入力デバイスを使用することができる。機械学習装置 100 に複数種類の入力デバイスが接続されてもよい。

【0040】

媒体リーダー 106 は、記録媒体 113 に記録されたプログラムやデータを読み取る読み取り装置である。記録媒体 113 として、フレキシブルディスク (FD: Flexible Disk) や HDD などの磁気ディスク、CD (Compact Disc) や DVD (Digital Versatile Disc) などの光ディスク、半導体メモリなど、任意の種類の記録媒体を使用することができる。媒体リーダー 106 は、例えば、記録媒体 113 から読み取ったプログラムやデータを、RAM 102 や HDD 103 などの他の記録媒体にコピーする。読み取られたプログラムは、例えば、CPU 101 によって実行される。なお、記録媒体 113 は可搬型記録媒体であってもよく、プログラムやデータの配布に用いられることがある。また、記録媒体 113 や HDD 103 を、コンピュータ読み取り可能な記録媒体とすることがある。

【0041】

通信インタフェース 107 は、ネットワーク 114 に接続され、ネットワーク 114 を介して他の情報処理装置と通信する。通信インタフェース 107 は、スイッチやルータなどの有線通信装置に接続される有線通信インタフェースでもよいし、基地局やアクセスポイントなどの無線通信装置に接続される無線通信インタフェースでもよい。

【0042】

次に、機械学習における訓練データサイズと予測性能について説明する。

第 2 の実施の形態の機械学習では、既知の事例を示す複数のレコードを含むデータ集合を予め収集しておく。レコードを、サンプル、インスタンス、行、データブロック、単位データなどと言うこともできる。機械学習装置 100 または他の情報処理装置が、センサデバイスなどの各種デバイスからネットワーク 114 経由でデータ集合を収集してもよい。収集されるデータ集合は、「ビッグデータ」と言われるサイズの大きなものであってもよい。各レコードは、通常、1 以上の説明変数の値と 1 つの目的変数の値とを含む。例えば、商品の需要予測を行う機械学習では、気温や湿度など商品需要に影響を与える要因を説明変数とし、商品需要量を目的変数とした実績データを収集する。

【0043】

機械学習装置 100 は、収集されたデータ集合の中から一部のレコードを訓練データとしてサンプリングし、訓練データを用いてモデルを学習する。モデルは、説明変数と目的変数との間の関係を示し、通常、1 以上の説明変数と 1 以上の係数と 1 つの目的変数とを含む。モデルは、例えば、線形式、二次以上の多項式、指数関数、対数関数などの各種数式によって表されてもよい。数式の形は、機械学習の前にユーザによって指定されてもよい。係数は、機械学習によって訓練データに基づいて決定される。

【0044】

学習されたモデルを用いることで、未知の事例の説明変数の値 (要因) から、未知の事例の目的変数の値 (結果) を予測することができる。例えば、来期の気象予報から来期の商品需要量を予測できる。モデルによって予測される結果は、0 以上 1 以下の確率などの

10

20

30

40

50

連続量であってもよいし、YES / NOの2値などの離散値であってもよい。

【0045】

学習されたモデルに対しては「予測性能」を算出することができる。予測性能は、未知の事例の結果を正確に予測する能力であり、「精度」と言うこともできる。機械学習装置100は、収集されたデータ集合の中から訓練データ以外のレコードをテストデータとしてサンプリングし、テストデータを用いて予測性能を算出する。機械学習装置100は、テストデータに含まれる説明変数の値をモデルに入力し、モデルが出力する目的変数の値（予測値）とテストデータに含まれる目的変数の値（実績値）とを比較する。なお、学習したモデルの予測性能を検証することを「バリデーション」と言うことがある。

【0046】

予測性能の指標としては、正答率（Accuracy）、適合率（Precision）、平均二乗誤差（MSE）、二乗平均平方根誤差（RMSE）などが挙げられる。例えば、結果がYES / NOの2値で表されるとする。また、n件のテストデータのレコードのうち、予測値 = YESかつ実績値 = YESの件数をTp、予測値 = YESかつ実績値 = NOの件数をFp、予測値 = NOかつ実績値 = YESの件数をFn、予測値 = NOかつ実績値 = NOの件数をTnとする。正答率は予測が当たった割合であり、 $(T_p + T_n) / n$ と算出される。適合率は「YES」の予測を間違えない確率であり、 $T_p / (T_p + F_p)$ と算出される。平均二乗誤差MSEは、各事例の実績値をYと表し予測値をyと表すと、 $\sum (Y - y)^2 / n$ と算出される。二乗平均平方根誤差RMSEは、 $(\sum (Y - y)^2 / n)^{1/2}$ と算出される。MSE = RMSE²である。

【0047】

ここで、訓練データからモデルを学習する手順（機械学習アルゴリズム）には様々なものが存在する。機械学習装置100が使用する機械学習アルゴリズムは、ユーザから指定されてもよいし、機械学習装置100が所定の評価方法に従って選択するようにしてもよい。機械学習装置100が使用できる機械学習アルゴリズムの数は、数十～数百程度であってもよい。機械学習アルゴリズムの一例として、ロジスティック回帰分析、サポートベクタマシン、ランダムフォレストなどを挙げることができる。

【0048】

ロジスティック回帰分析は、目的変数yの値と説明変数 x_1, x_2, \dots, x_d の値をS字曲線にフィッティングする回帰分析である。目的変数yおよび説明変数 x_1, x_2, \dots, x_d は、 $\log(y / (1 - y)) = a_1 x_1 + a_2 x_2 + \dots + a_d x_d + b$ の関係を満たすと仮定される。 a_1, a_2, \dots, a_d, b は係数であり、回帰分析によって決定される。

【0049】

サポートベクタマシンは、空間に配置されたレコードの集合を、2つのクラスに最も明確に分割するような境界面を算出する機械学習アルゴリズムである。境界面は、各クラスとの距離（マージン）が最大になるように算出される。

【0050】

ランダムフォレストは、複数の単位データを適切に分類するためのモデルを生成する機械学習アルゴリズムである。ランダムフォレストでは、データ集合からレコードをランダムにサンプリングする。説明変数の一部をランダムに選択し、選択した説明変数の値に応じてサンプリングしたレコードを分類する。説明変数の選択とレコードの分類を繰り返すことで、複数の説明変数の値に基づく階層的な決定木を生成する。レコードのサンプリングと決定木の生成を繰り返すことで複数の決定木を取得し、それら複数の決定木を合成することで、レコードを分類するための最終的なモデルを生成する。

【0051】

あるデータ集合に1つの機械学習アルゴリズムを適用する場合、訓練データとしてサンプリングするレコードの数（訓練データサイズ）が大きいほど予測性能は高くなる。

図3は、訓練データサイズと予測性能の関係例を示すグラフである。

【0052】

曲線31は、モデルの予測性能と訓練データサイズとの間の関係を示す。訓練データサ

10

20

30

40

50

イズ s_1, s_2, s_3, s_4, s_5 の間の大小関係は、 $s_1 < s_2 < s_3 < s_4 < s_5$ である。例えば、 s_2 は s_1 の 2 倍または 4 倍であり、 s_3 は s_2 の 2 倍または 4 倍であり、 s_4 は s_3 の 2 倍または 4 倍であり、 s_5 は s_4 の 2 倍または 4 倍である。

【0053】

曲線 31 が示すように、訓練データサイズが s_2 の場合の予測性能は s_1 の場合よりも高い傾向にある。同様に、訓練データサイズが s_3 の場合の予測性能は s_2 の場合よりも高い傾向にある。訓練データサイズが s_4 の場合の予測性能は s_3 の場合よりも高い傾向にある。訓練データサイズが s_5 の場合の予測性能は s_4 の場合よりも高い傾向にある。このように、訓練データサイズが大きくなるほど予測性能も高くなる傾向にある。ただし、予測性能が低い場合は、訓練データサイズの増加に応じて予測性能が大きく上昇する。一方で、予測性能には上限があり、予測性能が上限に近づくと、訓練データサイズの増加量に対する予測性能の上昇量の比は逓減する。すなわち、曲線 31 は、訓練データサイズの増加に応じて、ある上限に漸近するように予測性能が増加することを示している。

10

【0054】

このような訓練データサイズと予測性能との間の関係は、使用する機械学習アルゴリズムによって異なり、収集したデータ集合の性質（データ集合の種類）によっても異なる。このため、曲線 31 に示すような予測性能の上限や各訓練データサイズにおける予測性能を、機械学習を開始する前に推定することは容易でない。

【0055】

次に、予測性能の測定値の信頼性について説明する。

20

図 4 は、予測性能の測定値の分散例を示すグラフである。

ある訓練データサイズのもとで学習されたモデルの予測性能の測定値は、機械学習アルゴリズムとデータ集合の性質とから決まる期待値から乖離するリスクがある。すなわち、同じデータ集合を使用しても、訓練データおよびテストデータとして選択するレコードの偶然性によって、予測性能の測定値にばらつきが生じる。測定値の「ばらつき」は、分散や標準偏差などと解釈することもできる。分散は、訓練データサイズが小さいほど大きく、訓練データサイズが大きいほど小さくなる傾向にある。また、分散は、テストデータサイズが小さいほど大きく、テストデータサイズが大きいほど小さくなる傾向にある。

【0056】

グラフ 32 は、訓練データサイズと予測性能との間の関係を示す。ここでは、同じ機械学習アルゴリズムおよび同じデータ集合を用いて、訓練データサイズ 1 つ当たり 50 回ずつモデルの生成および予測性能の測定を行っている。また、テストデータサイズは、訓練データサイズの 2 分の 1 または 4 分の 1 とするなど、訓練データサイズに比例するようにして訓練データサイズに従属させている。グラフ 32 は、1 つの訓練データサイズにつき 50 個の予測性能の測定値をプロットしたものである。なお、グラフ 32 では、予測性能の指標として、値が大きいほど予測性能が高いことを示す正答率を用いている。

30

【0057】

グラフ 32 では、訓練データサイズ = 100 の場合の予測性能の測定値は、約 0.58 ~ 0.68 であり広範囲に広がっている。訓練データサイズ = 500 の場合の予測性能の測定値は、約 0.69 ~ 0.75 であり、訓練データサイズ = 100 の場合よりも範囲が狭くなっている。以降、訓練データサイズが大きくなるに従って測定値の範囲は狭くなる。訓練データサイズが十分に大きくなると、測定値は約 0.76 に収束している。

40

【0058】

以下では、予測性能の測定値の分散について更に検討する。

まず、バイアス・バリエーション分解について説明する。バイアス・バリエーション分解は、ある機械学習アルゴリズムの良否を評価するために用いられることがある。バイアス・バリエーション分解では、ロス（損失）とバイアスとバリエーションという 3 つの指標が用いられる。ロス = バイアスの二乗 + バリエーションという関係が成立する。

【0059】

ロスは、機械学習によって生成されるモデルが予測を外す度合いを示す指標である。ロ

50

スの種類には 0 - 1 ロスや二乗ロスなどがある。0 - 1 ロスは、予測に成功すれば 0 を付与し予測に失敗すれば 1 を付与することで算出されるロスであり、その期待値は予測が失敗する確率を示す。予測が外れることが少ないほど 0 - 1 ロスの期待値は小さく、予測が外れることが多いほど 0 - 1 ロスの期待値は大きい。二乗ロスは、予測値と真の値との差（予測誤差）の二乗である。予測誤差が小さいほど二乗ロスは小さく、予測誤差が大きいほど二乗ロスは大きい。期待ロス（ロスの期待値）と予測性能とは相互に変換できる。

【0060】

予測性能が正答率（Accuracy）でありロスが 0 - 1 ロスである場合、期待ロス = 1 - 予測性能である。予測性能が平均二乗誤差（MSE）でありロスが二乗ロスである場合、期待ロス = MSE である。予測性能が二乗平均平方根誤差（RMSE）でありロスが二乗ロスである場合、期待ロス = RMSE の二乗である。

10

【0061】

バイアスは、機械学習によって生成されるモデルの予測値が真の値に対して偏る程度を示す指標である。バイアスが小さいほど精度の高いモデルであると言える。バリエーションは、機械学習によって生成されるモデルの予測値がばらつく程度を示す指標である。バリエーションが小さいほど精度の高いモデルであると言える。ただし、バイアスとバリエーションの間にはトレードオフの関係があることが多い。

【0062】

次数の小さい多項式など複雑性の低いモデル（表現力の低いモデルと言うこともできる）では、モデルの係数をどの様に調整しても、複数のレコードの全てについて真の値に近い予測値を出力するようにすることは難しい。すなわち、複雑性の低いモデルを用いると複雑な事象を表現できない。よって、複雑性の低いモデルのバイアスは大きくなる傾向にある。この点、次数の大きい多項式など複雑性の高いモデル（表現力の高いモデルと言うこともできる）では、モデルの係数を適切に調整することで、複数のレコードの全てについて真の値に近い予測値を出力することができる余地がある。よって、複雑性の高いモデルのバイアスは小さくなる傾向にある。

20

【0063】

一方で、複雑性の高いモデルでは、訓練データとして使用するレコードの特徴に過度に依存したモデルが生成されるという過学習が生じるリスクがある。過学習によって生成されたモデルは、他のレコードについて適切な予測値を出力できないことが多い。例えば、 d 次の多項式を用いると、 $d + 1$ 個のレコードについて真の値と完全に一致する予測値を出力するモデル（残差が 0 のモデル）を生成できる。

30

【0064】

しかし、あるレコードについて残差が 0 になるモデルは、通常は過度に複雑なモデルであり、他のレコードについて予測誤差が著しく大きい予測値を出力してしまうリスクが高くなる。よって、複雑性の高いモデルのバリエーションは大きくなる傾向にある。この点、複雑性の低いモデルでは、予測誤差が著しく大きい予測値を出力してしまうリスクは低く、バリエーションは小さくなる傾向にある。このように、ロスの成分としてのバイアスとバリエーションは、モデルを生成する機械学習アルゴリズムの特性に依存している。

【0065】

次に、ロスとバイアスとバリエーションの形式的定義を説明する。ここでは、二乗ロスをバイアスとバリエーションに分解する例について説明する。

40

同一のデータ集合から m 個の訓練データ D_k ($k = 1, 2, \dots, m$) が抽出され、 m 個のモデルが生成されたとする。また、上記のデータ集合から i 個のレコードを含むテストデータ T が抽出されたとする。 i 番目のレコード（テストケースと言うこともできる）は、説明変数の値 X_i と目的変数の真の値 Y_i とを含む ($i = 1, 2, \dots, n$)。 k 番目のモデルからは説明変数の値 X_i に対して目的変数の予測値 y_{ik} が算出される。

【0066】

すると、 k 番目のモデルと i 番目のレコードとの間で算出される予測誤差 e_{ik} は $e_{ik} = Y_i - y_{ik}$ と定義され、そのロス（ここでは二乗ロス）は e_{ik}^2 と定義される。 i 番

50

目のレコードに対しては、バイアス B_i とバリエーション V_i とロス L_i が定義される。バイアス B_i は $B_i = E_D [e_{ik}]$ と定義される。 $E_D []$ は m 個の訓練データの間の平均値（期待値）を表す。バリエーション V_i は $V_i = V_D [e_{ik}]$ と定義される。 $V_D []$ は m 個の訓練データの間の分散を表す。ロス L_i は $L_i = E_D [e_{ik}^2]$ と定義される。前述のロスとバイアスとバリエーションの関係から $L_i = B_i^2 + V_i$ が成立する。

【0067】

テストデータ T 全体に対しては、期待バイアス EB_2 と期待バリエーション EV と期待ロス EL が定義される。期待バイアス EB_2 は $EB_2 = E_x [B_i^2]$ と定義される。 $E_x []$ は n 個のレコードの間の平均値（期待値）を表す。期待バリエーション EV は $EV = E_x [V_i]$ と定義される。期待ロス EL は $EL = E_x [L_i]$ と定義される。前述のロスとバイアスとバリエーションの関係から $EL = EB_2 + EV$ が成立する。

10

【0068】

バイアス・バリエーション分解の考え方を応用して、予測性能の測定値に生じる分散を推定することができる。測定値の分散は、次の数式によって近似される。 $VL = C \times (EL + EB_2) \times (EL - EB_2)$ 。 VL は訓練データサイズ s における予測性能の測定値の分散を表す。 C は定数である。 EL は訓練データサイズ s における期待ロスを表す。 EB_2 は期待バイアスを表す。以下、この数式の意味について説明を加える。

【0069】

図5は、予測性能の期待ロスおよび期待バイアスの例を示すグラフである。

曲線33は、訓練データサイズとロスの推定値との間の関係を示すロス曲線である。図3では縦軸が予測性能であるのに対し、図5では縦軸がロスに変換されている。前述のように予測性能とロスは、予測性能の指標とロスの指標に応じて相互に変換可能である。曲線33は、訓練データサイズの増加に応じてロスが単調に減少し一定の下限ロ스에漸近する非線形曲線である。訓練データサイズが小さいうちはロスの減少量が大きく、訓練データサイズが大きくなるとロスの減少量が小さくなる。

20

【0070】

訓練データサイズ s_p における曲線33上の点のロス（ロス = 0 から曲線33上の点までの距離）は、訓練データサイズ s_p の期待ロス EL_p に相当する。曲線33によって特定される下限ロスは、図3の曲線31によって特定される予測性能の上限に対応しており、0より大きい値である。例えば、予測性能の上限を c とおくと、予測性能が正答率である場合、下限ロスは $1 - c$ となる。予測性能が平均二乗誤差（ MSE ）である場合、下限ロスは c となる。予測性能が二乗平均平方根誤差（ $RMS E$ ）である場合、下限ロスは c^2 となる。下限ロスは、この機械学習アルゴリズムにとっての期待バイアス EB_2 に相当する。訓練データサイズが十分大きくなると、機械学習に用いる訓練データの特徴がデータ集合の特徴に一致し、期待バリエーションが0に近づくためである。

30

【0071】

期待ロス EL_p と期待バイアス EB_2 の差は、訓練データサイズ s_p におけるギャップとすることができる。ギャップは、訓練データサイズを大きくすることでその機械学習アルゴリズムがロスを低減できる余地を表している。ギャップは、図3の曲線31上の点と予測性能の上限との間の距離に対応し、訓練データサイズを大きくすることでその機械学習アルゴリズムが予測性能を改善できる余地を表しているとも言える。ギャップは、訓練データサイズ s_p における期待バリエーションの影響を受ける。

40

【0072】

次に、予測性能の測定値の分散を示す数式の数学的根拠について説明する。

(a) 問題の形式的な記述

同一のデータ集合から m セットの訓練データ D_1, D_2, \dots, D_m とテストデータ T が抽出されたとする。ある機械学習アルゴリズムに訓練データ D_k を与えて学習されたモデルを f_k とする（ $k = 1, 2, \dots, m$ ）。テストデータ T をレコード $\langle Y_i, X_i \rangle$ の集合とする（ $i = 1, 2, \dots, n$ ）。 X_i は説明変数の値（入力値）であり、 Y_i は入力値 X_i に対応する目的変数の既知の値（真値）である。入力値 X_i に対してモデル f_k が予測した

50

値（予測値）を $y_{ik} = f_k(X_i)$ とする。入力値 X_i に対するモデル f_k による予測の誤差は $e_{ik} = Y_i - y_{ik}$ と定義される。テストデータ T に含まれるレコードの数、すなわち、テストデータ T のサイズは n である。以下では主に、 i, j はテストデータ T のレコードを識別する添え字、 k はモデルを識別する添え字として使用する。

【0073】

機械学習アルゴリズムが回帰を目的とする場合、予測値は連続量であり、ロスの指標として数式（1）の二乗ロスが用いられることが多い。この二乗ロスをテストデータ T の全てのレコードについて平均したものが数式（2）の MSE （平均二乗誤差）である。

【0074】

【数1】

$$\text{loss}_{\text{sq}}(e) = e^2 \quad (1)$$

10

【0075】

【数2】

$$MSE = E_X[\text{loss}_{\text{sq}}(e_{ik})] = E_X[e_{ik}^2] \quad (2)$$

【0076】

ここで、 $E[\cdot]$ は期待値を求める演算子であり、 $V[\cdot]$ は分散を求める演算子である。ここで、 $E[\cdot]$ 、 $V[\cdot]$ に付加する添え字 X は、この演算子がテストデータ T の複数のレコードの間の演算であることを示す。また、 $E[\cdot]$ 、 $V[\cdot]$ に付加する添え字 M は、この演算子が複数のモデルの間の演算であることを示す。すなわち、 $E_X[\cdot]$ はテストデータ T の複数のレコードの間で平均化した期待値を示し、 $E_M[\cdot]$ は複数のモデルの間で平均化した期待値を示す。また、 $V_X[\cdot]$ はテストデータ T の複数のレコードの間の分散を示し、 $V_M[\cdot]$ は複数のモデルの間の分散を示す。また、 $\text{cov}(\cdot, \cdot)$ は共分散を求める共分散関数であり、 $\text{cor}(\cdot, \cdot)$ は相関係数を求める相関係数関数である。また、 $\text{cov}(\cdot, \cdot)$ 、 $\text{cor}(\cdot, \cdot)$ にも添え字 X, M が付加される。

20

【0077】

機械学習アルゴリズムが二値分類を目的とする場合、予測値は $\{-1, 1\}$ のような二値の離散値であり、ロスの指標として数式（3）の 0-1 ロスが用いられることが多い。この 0-1 ロスをテストデータ T の全てのレコードについて平均して 1 から引いたものが、数式（4）の正答率（Accuracy）である。

30

【0078】

【数3】

$$\text{loss}_{01}(e) = \begin{cases} 0 & \text{if } e = 0 \\ 1 & \text{otherwise} \end{cases} \quad (3)$$

40

【0079】

【数4】

$$\text{Accuracy} = 1 - E_X[\text{loss}_{01}(e_{ik})] = 1 - E_X[e_{ik}^2] \quad (4)$$

【0080】

MSE は値が小さいほど予測性能が高いことを示し、正答率は値が大きいほど予測性能が高いことを示す。ただし、両者ともテストデータ T 全体について平均化したロスがモデルの予測性能の良否を表している点で共通しており、モデルロスと行うことができる。モデル f_k に対するモデルロス ML_k は数式（5）の通りである。予測性能の分散は、数式

50

(6) のように複数のモデルの間のモデルロスの分散として表される。

【0081】

【数5】

$$ML_k = E_X[e_{ik}^2] \quad (5)$$

【0082】

【数6】

$$V_M[E_X[e_{ik}^2]] = V_M[ML_k] \quad (6)$$

10

【0083】

(b) バイアス・バリエンス分解

モデルによる予測で生じるロスはバイアスとバリエンスに分解できる。バイアスはモデルの予測値の偏りを示す量である。バイアスが低いモデルほど正確なモデルであると言える。表現力の低いモデル（調整可能な係数が少ないような複雑性の低いモデルなど）はバイアスが高くなる傾向にある。バリエンスはモデルの予測値のばらつきを示す量である。バリエンスが低いほど正確なモデルであると言える。表現力の高いモデル（調整可能な係数が多いような複雑性の高いモデルなど）はバリエンスが高くなる傾向にある。表現力の高いモデルには、訓練データに過剰適合するという過学習のリスクがある。

20

【0084】

テストデータTの入力値 X_i に対するロス L_i 、バイアス B_i およびバリエンス V_i は、数式(7)～(9)のように定義される。ロス L_i は複数のモデルの間の二乗誤差の期待値であり、バイアス B_i は複数のモデルの間の誤差の期待値であり、バリエンス V_i は複数のモデルの間の誤差の分散である。ロス L_i とバイアス B_i とバリエンス V_i との間には、数式(10)の関係（バイアス・バリエンス分解）が成立する。

【0085】

【数7】

$$L_i = E_M[e_{ik}^2] \quad (7)$$

30

【0086】

【数8】

$$B_i = E_M[e_{ik}] \quad (8)$$

【0087】

【数9】

$$V_i = V_M[e_{ik}] \quad (9)$$

40

【0088】

【数10】

$$L_i = V_i + B_i^2 \quad (10)$$

【0089】

様々な入力値 X_i に対するロス L_i の期待値を期待ロス E_L 、バイアス B_i の二乗の期待値を期待バイアス E_{B^2} 、バリエンス V_i の期待値を期待バリエンス E_V とする。期待ロス E_L 、期待バイアス E_{B^2} 、期待バリエンス E_V は、数式(11)～(13)のように

50

定義される。期待ロス EL と期待バイアス $EB2$ と期待バリエーション EV との間には、数式 (14) の関係 (バイアス・バリエーション分解) が成立する。

【0090】

【数11】

$$EL = E_X[L_i] \quad (11)$$

【0091】

【数12】

$$EB2 = E_X[B_i^2] \quad (12)$$

10

【0092】

【数13】

$$EV = E_X[V_i] \quad (13)$$

【0093】

【数14】

$$EL = EV + EB2 \quad (14)$$

20

【0094】

ここでの目的は、 EL 、 $EB2$ 、 EV とモデルロスの分散との間の関係を導出することである。なお、期待ロス EL とモデルロス ML_k の期待値とは、数式 (15) に示すように等価である。一方、ロス L_i の分散とモデルロス ML_k の分散とは等価でない。以下では、予測性能の分散を推定する数式を次の流れで導出する。第1に、ロスの分散をバイアスとバリエーションで記述する。第2に、モデルロスの分散をインスタンス成分と相互作用成分に分解する。第3に、インスタンス成分を算出する。第4に、相互作用成分を算出する。第5に、モデルロスの分散をバイアスとバリエーションで記述する。

30

【0095】

【数15】

$$EL = E_X[L_i] = E_X[E_M[e_{ik}^2]] = E_M[E_X[e_{ik}^2]] = E_M[ML_k] \quad (15)$$

【0096】

(c) ロスの分散をバイアスとバリエーションで記述

テストデータ T の入力値 X_i を固定して複数のモデルの誤差を並べた誤差ベクトルを考える。誤差 e を確率変数とみなしてその分布が正規分布に従うと仮定すると、複数のモデルの間のロスの分散は数式 (16) のように定義され、バイアス B_i とバリエーション V_i の組またはロス L_i とバイアス B_i の組によって記述することができる。数式 (16) の1行目から2行目への変形では、数式 (17) に示す統計学上の性質 (確率変数の4乗の期待値) が利用されている。数式 (17) において X は確率変数であり、 S は歪度であり、 K は尖度である。正規分布の場合は $S = 0$ かつ $K = 3$ である。

40

【0097】

【数16】

$$\begin{aligned} V_M[e_{ik}^2] &= E_M[e_{ik}^4] - (E_M[e_{ik}^2])^2 \\ &= 2V_i^2 + 4V_i B_i^2 = 2L_i^2 - 2B_i^4 \end{aligned} \quad (16)$$

50

【 0 0 9 8 】

【数 1 7】

$$E[X^4] = K(V[X])^2 + 4S(V[X])^{1.5}E[X] + 6V[X](E[X])^2 + (E[X])^4 \quad (17)$$

【 0 0 9 9 】

(d) モデルロスの分散をインスタンス成分と相互作用成分に分解

分散の基本的性質から、予測性能の分散（複数のモデルの間のモデルロスの分散）について数式（18）が成立する。これを $n \times n$ 行列の成分の平均と考えると、 $i = j$ である対角成分は入力値 X_i に対するロスの分散を表しており、その相関係数は 1 になる。一方、 $i \neq j$ である非対角成分の相関係数は異なる入力値の間のロスの相関を表している。異なる入力値に対する誤差の発生状況は共通点が少ないため、その相関係数の絶対値は十分に小さくなることが多く、予測性能の高いモデルほどその相関係数は 0 に近づく。対角成分と非対角成分とは性質が異なるため、数式（19）のように両者を分離して考える。

10

【 0 1 0 0 】

【数 1 8】

$$\begin{aligned} V_M[E_X[e_{ik}^2]] &= E_{X_i}[E_{X_j}[\text{cov}_M(e_{ik}^2, e_{jk}^2)]] \\ &= E_{X_i}[E_{X_j}[\text{cor}_M(e_{ik}^2, e_{jk}^2) \cdot (V_M[e_{ik}^2])^{0.5} (V_M[e_{jk}^2])^{0.5}]]] \end{aligned} \quad (18)$$

20

【 0 1 0 1 】

【数 1 9】

$$\begin{aligned} E_{X_i}[E_{X_j}[\text{cov}_M(e_{ik}^2, e_{jk}^2)]] &= \frac{n}{n^2} E_{X_i}[E_{X_j, j=i}[\text{cov}_M(e_{ik}^2, e_{jk}^2)]] + \frac{n(n-1)}{n^2} E_{X_i}[E_{X_j, j \neq i}[\text{cov}_M(e_{ik}^2, e_{jk}^2)]] \\ &= \frac{1}{n} E_X[V_M[e_{ik}^2]] + \frac{n-1}{n} E_{X_i}[E_{X_j, j \neq i}[\text{cov}_M(e_{ik}^2, e_{jk}^2)]] \end{aligned} \quad (19)$$

30

【 0 1 0 2 】

数式（19）では、モデルロスの分散を第 1 項のインスタンス成分と第 2 項の相互作用成分とに分解している。第 1 項はロスの分散の期待値を表しており、モデルロスの分散の大部分を占めることが多い。第 2 項は異なる入力値の間の共分散の期待値を表しており、モデルロスの分散に対する寄与は小さいことが多い。第 1 項はテストデータ T のサイズ n に反比例するため、テストデータ T のレコードを増やすことでモデルロスの分散を低減できる。ただし、第 2 項が存在することから低減効果には限界がある。

40

(e) インスタンス成分を算出

数式（19）の第 1 項について検討する。上記の数式（16）より数式（20）が成立する。ここで、数式（20）の第 1 項と第 2 項を算出するために幾つかの仮定をおく。多くの機械学習アルゴリズムは不偏推定量を出力するようにモデルを学習することから、数式（21）のように誤差の期待値が 0 になるという仮定をおく。数式（21）からバイアス B_i について数式（22）の性質が導出される。

【 0 1 0 3 】

【数 2 0】

$$E_X[V_M[e_{ik}^2]] = 2E_X[L_i^2] - 2E_X[B_i^4] \quad (20)$$

50

【 0 1 0 4 】

【数 2 1】

$$E_X[e_{ik}] = 0 \quad (21)$$

【 0 1 0 5 】

【数 2 2】

$$E_X[B_i] = 0 \quad (22)$$

10

【 0 1 0 6 】

確率分布の中には、訓練データサイズや訓練データのサンプリング方法に依存して期待値や分散が変化することはあっても、確率分布の形状を示す歪度や尖度は変化しない（または、変化が非常に緩やかである）ものがあると仮定する。具体的には、入力値 X_i に対する複数のモデルの間の誤差の分布は正規分布を形成し、尖度 = 3 かつ歪度 = 0 になることを仮定する。また、バイアス B_i の分布の尖度 K_1 は変化しないことを仮定する。バイアス B_i の分布の尖度 K_1 は、数式 (23) のように定義される。数式 (23) と上記の数式 (12) から数式 (24) が算出される。

【 0 1 0 7 】

【数 2 3】

$$K_1 = \frac{E_X[B_i^4]}{(E_X[B_i^2])^2} \quad (23)$$

20

【 0 1 0 8 】

【数 2 4】

$$E_X[B_i^4] = K_1 \cdot EB^2^2 \quad (24)$$

【 0 1 0 9 】

また、モデル f_k に対する複数の入力値の間の誤差の分布の尖度 K_2 は、モデル間で共通でありかつ変化しないことを仮定する。尖度 K_2 は数式 (25) のように定義される。 K_1 , K_2 の値はそれぞれ 3 ~ 10 の範囲内であることが多く、両者は近いことが多い。

30

【 0 1 1 0 】

【数 2 5】

$$K_2 = \frac{E_M[E_X[e_{ik}^4]]}{E_M[(E_X[e_{ik}^2])^2]} \quad (25)$$

【 0 1 1 1 】

数式 (25) から数式 (26) が導出される。数式 (26) を数式 (18) , (19) に代入することで数式 (27) が算出される。ここで、尖度 K_2 はサイズ n より十分に小さいため、 $1 - K_2 / n$ は 1 に近似される。数式 (20) , (23) を数式 (18) , (19) に代入することで数式 (28) が算出される。数式 (28) から数式 (27) を減算して数式 (29) が算出される。そして、数式 (20) , (24) , (29) から数式 (30) が算出される。これが、数式 (19) の第 1 項の主要成分である。

40

【 0 1 1 2 】

【数 2 6】

50

$$\begin{aligned}
E_X[V_M[e_{ik}^2]] &= E_X[E_M[e_{ik}^4]] - E_X[(E_M[e_{ik}^2])^2] \\
&= K2 \cdot E_M[(E_X[e_{ik}^2])^2] - K2(E_M[E_X[e_{ik}^2]])^2 \\
&\quad + K2(E_M[E_X[e_{ik}^2]])^2 - E_X[(E_M[e_{ik}^2])^2] \quad (26) \\
&= K2 \cdot V_M[E_X[e_{ik}^2]] + K2 \cdot EL^2 - E_X[L_i^2]
\end{aligned}$$

【 0 1 1 3 】

【 数 2 7 】

10

$$\begin{aligned}
\left(1 - \frac{K2}{n}\right) V_M[E_X[e_{ik}^2]] &\approx V_M[E_X[e_{ik}^2]] \\
= \frac{1}{n}(K2 \cdot EL^2 - E_X[L_i^2]) &+ \frac{n-1}{n} E_{X_i}[E_{X_{j,j \neq i}}[\text{cov}_M(e_{ik}^2, e_{jk}^2)]] \quad (27)
\end{aligned}$$

【 0 1 1 4 】

【 数 2 8 】

20

$$\begin{aligned}
V_M[E_X[e_{ik}^2]] &\quad (28) \\
= \frac{1}{n}(2E_X[L_i^2] - 2K1 \cdot EB2^2) &+ \frac{n-1}{n} E_{X_i}[E_{X_{j,j \neq i}}[\text{cov}_M(e_{ik}^2, e_{jk}^2)]]
\end{aligned}$$

【 0 1 1 5 】

【 数 2 9 】

$$E_X[L_i^2] = \frac{1}{3}K2 \cdot EL^2 + \frac{2}{3}K1 \cdot EB2^2 \quad (29)$$

【 0 1 1 6 】

【 数 3 0 】

30

$$E_X[V_M[e_{ik}^2]] = \frac{2}{3}K2 \cdot EL^2 - \frac{2}{3}K1 \cdot EB2^2 \quad (30)$$

【 0 1 1 7 】

(f) 相互作用成分を算出

不動点 $Cor1v$ を数式 (3 1) のように定義する。不動点 $Cor1v$ は、訓練データサイズを変化させても値が変化しないかまたは非常に緩やかに変化することが多いため、ここでは訓練データサイズに依存しないと仮定する。不動点 $Cor1v$ の値は 0 . 0 0 1 ~ 0 . 1 程度であることが多い。

40

【 0 1 1 8 】

【 数 3 1 】

$$\begin{aligned}
Cor1v &= V_{i,j,j \neq i}[\text{cor}_M(e_{ik}, e_{jk})] \quad (31) \\
&= E_{X_i}[E_{X_{j,j \neq i}}[\text{cor}_M(e_{ik}, e_{jk})^2]] - (E_{X_i}[E_{X_{j,j \neq i}}[\text{cor}_M(e_{ik}, e_{jk})]])^2
\end{aligned}$$

【 0 1 1 9 】

ここで、数式 (3 2) に示す統計学上の性質 (誤差の相関係数の期待値) を利用する。

50

誤差の期待値が 0 であるとき、2 つの誤差の相関係数の期待値は 0 に近似する。この性質から数式 (3 3) が成立し、上記の数式 (3 1) から数式 (3 4) が算出される。

【 0 1 2 0 】

【数 3 2】

$$E_X[e_{ik}] = 0 \Rightarrow E_{X_i}[E_{X_j}[\text{cor}_M(e_{ik}, e_{jk})]] \approx 0 \quad (32)$$

【 0 1 2 1 】

【数 3 3】

$$E_{X_i}[E_{X_j, j \neq i}[\text{cor}_M(e_{ik}, e_{jk})]] = \frac{n^2 E_{X_i}[E_{X_j}[\text{cor}_M(e_{ik}, e_{jk})]] - n}{n(n-1)} \quad (33)$$

$$\approx \frac{-1}{n-1}$$

【 0 1 2 2 】

【数 3 4】

$$E_{X_i}[E_{X_j, j \neq i}[\text{cor}_M(e_{ik}, e_{jk})^2]] \approx \text{Cor1}v + \frac{1}{(n-1)^2} \quad (34)$$

【 0 1 2 3 】

また、数式 (3 5) が成立する。数式 (3 5) の 2 行目から 3 行目への変形では、相関係数 cor_M とバリエーション V_i, V_j とは互いに独立であることを仮定している。数式 (3 5) の 3 行目から 4 行目への変形では、上記の数式 (3 4) を利用しており、 $V_i \cdot V_j$ の期待値が EV^2 に近似することを利用している。数式 (3 5) の 4 行目の近似では、テストデータサイズ n が 1 より十分に大きいため $1 / (n - 1)^2$ を無視している。

【 0 1 2 4 】

【数 3 5】

$$E_{X_i}[E_{X_j, j \neq i}[\text{cov}_M(e_{ik}, e_{jk})^2]]$$

$$= E_{X_i}[E_{X_j, j \neq i}[\text{cor}_M(e_{ik}, e_{jk})^2 V_i \cdot V_j]]$$

$$= E_{X_i}[E_{X_j, j \neq i}[\text{cor}_M(e_{ik}, e_{jk})^2]] \cdot E_{X_i}[E_{X_j, j \neq i}[V_i \cdot V_j]] \quad (35)$$

$$\approx \left(\text{Cor1}v + \frac{1}{(n-1)^2} \right) EV^2 \approx \text{Cor1}v \cdot EV^2 \approx \text{Cor1}v(EL - EB2)^2$$

【 0 1 2 5 】

ここで、数式 (3 6) に示す統計学上の性質 (共分散の二乗と二乗の共分散の関係) を利用する。確率変数 X, Y の結合確率が二次元正規分布に従うならば数式 (3 6) が成立する。誤差の分散が正規分布に従うため、数式 (3 6) を利用して数式 (3 7) が算出される。また、数式 (3 8) が成立する。数式 (3 8) の 1 行目から 2 行目への変形では、共分散 cov_M とバイアス B_i, B_j は概ね独立であることを仮定している。数式 (3 8) の 2 行目の近似では、 $B_i B_j$ の期待値はバイアス B_i の期待値の二乗に近似しその結果 0 に近似するという性質を利用している。数式 (3 5) , (3 8) を数式 (3 7) に代入することで数式 (3 9) が算出される。これが、数式 (1 9) の第 2 項の主要成分である。

【 0 1 2 6 】

【数 3 6】

10

20

30

40

50

$$\text{cov}(X^2, Y^2) = 2\text{cov}(X, Y)^2 + 4E[X]E[Y]\text{cov}(X, Y) \quad (36)$$

【 0 1 2 7 】

【 数 3 7 】

$$\begin{aligned} E_{X_i}[E_{X_{j,j \neq i}}[\text{cov}_M(e_{ik}^2, e_{jk}^2)]] & \quad (37) \\ = E_{X_i}[E_{X_{j,j \neq i}}[2\text{cov}_M(e_{ik}, e_{jk})^2]] + E_{X_i}[E_{X_{j,j \neq i}}[4B_i B_j \text{cov}_M(e_{ik}, e_{jk})]] \end{aligned}$$

10

【 0 1 2 8 】

【 数 3 8 】

$$\begin{aligned} E_{X_i}[E_{X_{j,j \neq i}}[B_i B_j \text{cov}_M(e_{ik}, e_{jk})]] & \\ = E_{X_i}[E_{X_{j,j \neq i}}[\text{cov}_M(e_{ik}, e_{jk})]] \cdot E_{X_i}[E_{X_{j,j \neq i}}[B_i B_j]] \approx 0 & \quad (38) \end{aligned}$$

【 0 1 2 9 】

【 数 3 9 】

$$E_{X_i}[E_{X_{j,j \neq i}}[\text{cov}_M(e_{ik}^2, e_{jk}^2)]] \approx 2 \text{Cor1v}(EL - EB2)^2 \quad (39)$$

20

【 0 1 3 0 】

(g) モデルロスの分散をバイアスとバリエーションで記述

上記の数式(18), (19), (30), (39)より数式(40)の近似式が成立する。尖度K2は尖度K1に近似するため、数式(40)は数式(41)のように近似される。典型的にはK1(EL + EB2)はcor1v(EL - EB2)より十分に大きい。ため、数式(41)は更に数式(42)のように近似される。尖度K1は事前には不明であるが、分散の比が判明すれば実用上十分であることも多い。そこで、数式(42)は比例定数Cを用いて数式(43)のように単純化できる。これにより、予測性能の測定値の分散が、期待ロスELと期待バイアスEB2の差に比例し、かつ、期待ロスELと期待バイアスEB2の和に比例するという数式が導出される。

30

【 0 1 3 1 】

【 数 4 0 】

$$\begin{aligned} V_M[E_X[e_{ik}^2]] \approx \frac{1}{n} \left(\frac{2}{3} K2 \cdot EL^2 - \frac{2}{3} K1 \cdot EB2^2 \right) & \quad (40) \\ + \frac{n-1}{n} \cdot 2 \text{Cor1v}(EL - EB2)^2 & \end{aligned}$$

【 0 1 3 2 】

【 数 4 1 】

$$\begin{aligned} V_M[E_X[e_{ik}^2]] & \quad (41) \\ \approx \left(\frac{1}{n} \frac{2}{3} K1(EL + EB2) + \frac{n-1}{n} \cdot 2 \text{Cor1v}(EL - EB2) \right) (EL - EB2) \end{aligned}$$

40

【 0 1 3 3 】

【 数 4 2 】

50

$$V_M[E_X[e_{ik}^2]] \approx \frac{1}{n} \frac{2}{3} K1(EL + EB2)(EL - EB2) \quad (42)$$

【 0 1 3 4 】

【 数 4 3 】

$$\begin{aligned} V_M[E_X[e_{ik}^2]] &\approx C(EL - EB2)(EL + EB2) \\ &\approx C \cdot EV(EV + 2 EB2) \end{aligned} \quad (43)$$

10

【 0 1 3 5 】

予測性能の測定値の分散を示す数式(41)は、数式(44)のように変形することができる。テストデータサイズnに着目すると、数式(44)の第1項は、予測性能の測定値の分散のうち、テストデータサイズnの増大に応じて減少するテストデータ依存成分に相当する。一方、数式(44)の第2項は、予測性能の測定値の分散のうち、テストデータサイズnの増大によつては減少しない訓練データ依存成分に相当する。このため、数式(44)は、予測性能の測定値の分散が、テストデータサイズnの増大によつて減少するものの、0より大きい下限が存在することを示している。

【 0 1 3 6 】

20

【 数 4 4 】

$$\begin{aligned} V_M[E_X[e_{ik}^2]] &\approx \frac{1}{n} \frac{2}{3} K1(EL^2 - EB2^2) \\ &\quad + \frac{n-1}{n} \cdot 2 \text{Cor1v}(EL - EB2)^2 \end{aligned} \quad (44)$$

【 0 1 3 7 】

前述の図4では、テストデータサイズを訓練データサイズの2分の1や4分の1とするなど、テストデータサイズを訓練データサイズに比例するように決定していた。しかし、このようなテストデータサイズの決定方法は、予測性能の測定値の信頼性と予測性能の測定の負荷とを両立させる観点から、改善の余地がある。テストデータサイズが小さ過ぎると、予測性能の測定値が有する潜在的な分散が大きくなり、算出される測定値の信頼性が低下する。一方、テストデータサイズが大き過ぎると、予測性能の測定値の分散があまり減少せず、測定値の信頼性の向上にあまり寄与しないにもかかわらず、テスト処理を無駄に繰り返すことになり、テスト処理の負荷が増大する。

30

【 0 1 3 8 】

この点、数式(44)が示す予測性能の測定値の分散とテストデータサイズnとの間の対応関係によれば、測定値の信頼性とテスト負荷とを両立させるような効率的なテストデータサイズnが存在することになる。そこで、第2の実施の形態の機械学習装置100は、数式(44)に基づいて、適切なテストデータサイズを決定する。

40

【 0 1 3 9 】

データ集合と機械学習アルゴリズムが特定されると、機械学習装置100は、数式(44)の尖度K1と不動点Cor1vと期待バイアスEB2を決定する。これにより、機械学習装置100は、テストデータサイズnと期待ロスELを引数として有し、予測性能の測定値の分散を推定する分散関数f(n, EL)を生成する。尖度K1と不動点Cor1vと期待バイアスEB2は、訓練データサイズに依存しないパラメータである。そのため、データ集合と機械学習アルゴリズムが同じであれば、訓練データサイズが異なっても、同じ分散関数を用いて予測性能の測定値の分散を推定することができる。

【 0 1 4 0 】

50

あるデータ集合と機械学習アルゴリズムの組に対する尖度 K_1 と不動点 $Cor_1 v$ と期待バイアス $E B_2$ を決定するには、前述のように、 m セットの訓練データと 1 セットのテストデータの間の網羅的な誤差を示す誤差プロファイルを用意することが好ましい。そこで、機械学習装置 100 は、予測性能を測定したいモデルの訓練データサイズよりも十分に小さい訓練データを、同一のデータ集合から m セット抽出し、 m セットの訓練データを用いて機械学習により m 個のモデルを生成する。また、機械学習装置 100 は、十分に小さいテストデータを当該データ集合から抽出し、テストデータに含まれる複数のレコードと m 個のモデルとの間で網羅的に誤差を算出する。

【0141】

例えば、予測性能を測定したいモデルの訓練データサイズが 100 万レコードであるとすると、この場合、誤差プロファイル生成のための訓練データの個数を 10 セットとし、各訓練データのサイズを 1 万レコードとする。また、テストデータサイズを、訓練データサイズの 2 分の 1 である 5000 レコードとする。これにより、10 個のモデルとテストデータの 5000 レコードとの間で、 10×5000 個の誤差が算出される。機械学習装置 100 は、この誤差プロファイルを用いて、数式 (23) の尖度 K_1 と、数式 (31) の不動点 $Cor_1 v$ と、数式 (12) の期待バイアス $E B_2$ を算出する。

10

【0142】

分散関数 $f(n, EL)$ が生成されると、機械学習装置 100 は、予測性能を測定したいモデルに対応する期待ロス EL を分散関数に代入する。期待ロス EL は、図 5 に示すように、データ集合と機械学習アルゴリズムが同じでも訓練データサイズに応じて変化する。そのため、対象のモデルに対応する期待ロス EL を使用することになる。

20

【0143】

ある訓練データサイズに対応する期待ロス EL は、測定せずに与えられることもあるし対象のモデルから測定して求めることもある。測定しない場合として、データ集合および機械学習アルゴリズムが同一であり訓練データサイズが異なる複数のモデルの予測性能が、既に測定済みである場合が考えられる。その場合、回帰分析などの統計的方法により、それら測定値から未知の期待ロス EL を推定することが考えられる。未知の期待ロス EL の推定には、図 3 や図 5 の非線形曲線を利用することができる。

【0144】

対象のモデルから測定する場合、例えば、機械学習装置 100 は、誤差プロファイルの生成に使用した小さなテストデータを対象のモデルに入力し、テストデータに含まれる複数のレコードに対応する誤差を算出する。そして、機械学習装置 100 は、それら誤差から数式 (11) の期待ロス EL を算出する。例えば、5000 レコードのテストデータから 5000 個の誤差が算出され、期待値としての期待ロス EL が算出される。

30

【0145】

上記の方法で測定される期待ロス EL は、対象のモデルが大きい訓練データサイズで学習されているため、当該大きい訓練データサイズに対応した測定値になる。ただし、小さいテストデータを使用するため、大きいテストデータを使用して測定される本来の期待ロス EL と比較すると、測定値の分散が大きくなる。その点で、小さなテストデータで測定される期待ロス EL は、近似値または推定値であると言える。

40

【0146】

分散関数 $f(n, EL)$ に入力する期待ロス EL の精度を上げるため、機械学習装置 100 は、期待ロス EL の推定とテストデータサイズ n の選択を 2 回繰り返してもよい。例えば、機械学習装置 100 は、小さなテストデータで測定した期待ロス EL を分散関数 $f(n, EL)$ に入力し、以下で説明する方法でテストデータサイズ n を仮選択する。機械学習装置 100 は、データ集合から当該仮選択したサイズのテストデータを抽出し、抽出したテストデータを用いて期待ロス EL を再測定する。機械学習装置 100 は、再測定した期待ロス EL を分散関数 $f(n, EL)$ に入力し、以下で説明する方法でテストデータサイズ n を再選択し、これを最終的なテストデータサイズと決定する。

【0147】

50

期待ロス $E L$ を分散関数 $f(n, E L)$ に入力して期待ロス $E L$ を固定すると、分散関数は、テストデータサイズ n と分散の推定値とを 1 対 1 に対応付ける対応関係を表す。機械学習装置 100 は、分散関数のテストデータサイズ n を変動させながら分散の推定値を評価することで、適切なテストデータサイズ n を決定する。

【0148】

テストデータサイズ n と分散の推定値との対応関係は、テストデータサイズ n の増加に応じて、分散の推定値が下限に漸近するように減少する非線形曲線に相当する。テストデータサイズ n が小さいうちは、テストデータサイズ n の単位増加量あたりの分散の推定値の減少量が大きい。テストデータサイズ n が大きいほど、テストデータサイズ n の単位増加量あたりの分散の推定値の減少量が小さくなる。予測性能の測定値の信頼性を維持しつつ

10

【0149】

例えば、機械学習装置 100 は、効果指標として $f(n, E L) / f(2 * n, E L)$ を算出する。この効果指標は、テストデータサイズ n を 2 倍にした場合の分散の減少率に相当し、分散減少効果の評価指標である。効果指標の値が大きいほど分散減少効果が大きいことを示し、効果指標の値が小さいほど分散減少効果が小さいことを示す。テストデータサイズ n と分散の推定値の関係から、 n が大きいほど効果指標の値は小さくなる。

【0150】

機械学習装置 100 は、小さいテストデータサイズ n で効果指標の値を算出し、閾値と比較する。閾値は、1.1 などと予め決めておく。効果指標の値が閾値以上である場合、機械学習装置 100 は、テストデータサイズ n を 2 倍にし、効果指標の値が閾値未満になるまで上記を繰り返す。効果指標の値が閾値未満になると、機械学習装置 100 は、その時点のテストデータサイズ n を適切なテストデータサイズとして決定する。

20

【0151】

なお、上記の方法におけるテストデータサイズ n の増加速度である「2倍」や閾値の「1.1」は調整可能パラメータであり、ユーザがこれらのパラメータを変更することも可能である。また、分散関数 $f(n, E L)$ から適切なテストデータサイズ n を探索する他の方法として、例えば、機械学習装置 100 は、テストデータサイズ n を無限大にした場合の分散の推定値の下限を算出する。そして、機械学習装置 100 は、分散の推定値が下限の所定倍（例えば、1.1倍）になるようなテストデータサイズ n を選択する。

30

【0152】

このようにして機械学習装置 100 によって決定されるテストデータサイズは、訓練データサイズの 2 分の 1 または 4 分の 1 をテストデータサイズとする慣習的方法と比べて、十分に小さいサイズとなる。例えば、訓練データサイズが 100 万レコードである場合、慣習的方法では、テストデータサイズが 50 万レコードまたは 25 万レコードとなる。これに対して、第 2 の実施の形態の方法によれば、予測性能の測定値の分散を慣習的方法と同程度に維持しつつ、テストデータサイズを数万レコード程度に削減できる。よって、予測性能の測定値の信頼性を維持しつつ、テスト処理を高速化できる。

【0153】

なお、第 2 の実施の形態で決定される最終的なテストデータサイズ n は、慣習的方法よりも十分に小さい。そのため、テストデータを用いて期待ロス $E L$ を算出することを 1 回または 2 回行って、全体のテスト処理の負荷は慣習的方法よりも十分に小さくなる。

40

【0154】

次に、機械学習装置 100 の機能および処理手順について説明する。

図 6 は、機械学習装置の機能例を示すブロック図である。

機械学習装置 100 は、データ記憶部 121、制御情報記憶部 122、学習結果記憶部 123、モデル生成部 124、テスト実行部 125、テストサイズ決定部 126 および機械学習制御部 127 を有する。データ記憶部 121、制御情報記憶部 122 および学習結果記憶部 123 は、例えば、RAM 102 または HDD 103 の記憶領域を用いて実現さ

50

れる。モデル生成部 1 2 4、テスト実行部 1 2 5、テストサイズ決定部 1 2 6 および機械学習制御部 1 2 7 は、例えば、CPU 1 0 1 が実行するプログラムを用いて実現される。

【 0 1 5 5 】

データ記憶部 1 2 1 は、訓練データまたはテストデータに使用可能な多数のレコードを含むデータ集合を記憶する。各レコードは、説明変数の値と教師ラベルである目的変数の値とを含む。データ集合は、数百万レコードなどサイズの大きなものであってもよい。機械学習装置 1 0 0 は、ユーザからデータ集合を受け付けてもよいし、他の情報処理装置からデータ集合を受信してもよいし、センサデバイスからデータ集合を収集してもよい。

【 0 1 5 6 】

制御情報記憶部 1 2 2 は、訓練データを用いたモデルの学習やテストデータを用いたモデルの予測性能の測定の過程で生成される各種の制御情報を記憶する。制御情報には、分散関数の生成に用いられる誤差プロファイルや分散関数のパラメータが含まれる。

10

【 0 1 5 7 】

学習結果記憶部 1 2 3 は、機械学習の結果を記憶する。機械学習の結果には、学習されたモデルおよび当該モデルの予測性能の測定値が含まれる。

モデル生成部 1 2 4 は、機械学習によりモデルを生成する。モデル生成部 1 2 4 は、機械学習制御部 1 2 7 から機械学習アルゴリズムの指定と訓練データを受け付ける。モデル生成部 1 2 4 は、指定された機械学習アルゴリズムに従って、訓練データに含まれるレコードを用いてモデルの係数を決定することでモデルを学習する。機械学習アルゴリズムには、回帰分析、サポートベクタマシン、ランダムフォレストなどが含まれる。モデル生成部 1 2 4 は、学習されたモデルを機械学習制御部 1 2 7 に提供する。

20

【 0 1 5 8 】

テスト実行部 1 2 5 は、モデルのテストを行う。テスト実行部 1 2 5 は、機械学習制御部 1 2 7 からモデルとテストデータを受け付ける。テスト実行部 1 2 5 は、テストデータのレコードに含まれる説明変数の値をモデルに入力し、モデルに従って目的変数の予測値を算出する。テスト実行部 1 2 5 は、テストデータのレコードに含まれる目的変数の真値とモデルから算出された予測値とを比較して、誤差を算出する。そして、テスト実行部 1 2 5 は、誤差を列挙した誤差プロファイルを生成する。

【 0 1 5 9 】

テスト実行部 1 2 5 は、誤差プロファイルを機械学習制御部 1 2 7 に提供する。または、テスト実行部 1 2 5 は、誤差プロファイルを予測性能または期待ロスに変換し、予測性能または期待ロスを機械学習制御部 1 2 7 に提供する。予測性能の指標には、正答率、適合率、平均二乗誤差、二乗平均平方根誤差などが含まれる。予測性能または期待ロスは、テストデータに含まれる複数のレコードに対応する誤差から算出することができる。機械学習制御部 1 2 7 に提供される情報は、機械学習制御部 1 2 7 の要求に応じて変わる。

30

【 0 1 6 0 】

テストサイズ決定部 1 2 6 は、テストデータサイズを決定する。まず、テストサイズ決定部 1 2 6 は、機械学習制御部 1 2 7 から誤差プロファイルを受け付ける。この誤差プロファイルは、 m セットの小さな訓練データを用いて学習された m 個のモデルに対して、小さなテストデータを用いて測定された誤差を列挙したものである。テストサイズ決定部 1 2 6 は、この誤差プロファイルを用いて、予測性能の測定値の分散を推定するための分散関数のパラメータを決定する。分散関数のパラメータには、尖度 K_1 と不動点 Cor_1v と期待バイアス EB_2 が含まれる。テストサイズ決定部 1 2 6 は、分散関数の式や分散関数のパラメータの決定方法を予め知っている。テストサイズ決定部 1 2 6 は、分散関数のパラメータを機械学習制御部 1 2 7 に提供する。

40

【 0 1 6 1 】

また、テストサイズ決定部 1 2 6 は、機械学習制御部 1 2 7 から、先に算出した分散関数のパラメータと、対象のモデルの訓練データサイズに対応する期待ロス EL を受け付ける。テストサイズ決定部 1 2 6 は、分散関数 $f(n, EL)$ に期待ロス EL を代入し、テストデータサイズ n を変えながら分散の推定値を算出する。そして、テストサイズ決定部

50

126は、適切なテストデータサイズ n を決定して機械学習制御部127に提供する。例えば、テストサイズ決定部126は、分散の推定値から算出される効果指標の値が閾値以上である範囲で、最大のテストデータサイズ n を検出する。

【0162】

機械学習制御部127は、機械学習を制御する。まず、機械学習制御部127は、モデルの学習および予測性能の測定の対象とする機械学習アルゴリズムおよび訓練データサイズを特定する。対象の機械学習アルゴリズムおよび訓練データサイズは、ユーザから指定されてもよいし、所定の規則に従って機械学習制御部127が選択してもよい。

【0163】

次に、機械学習制御部127は、テストサイズ決定部126に分散関数のパラメータを決定させる。ただし、分散関数のパラメータの決定は、予測性能を測定する対象のモデルが学習された後に行うようにすることも可能である。

10

【0164】

分散関数のパラメータの決定では、機械学習制御部127は、 m セットの小さな訓練データと1セットの小さなテストデータを、データ記憶部121に記憶されたデータ集合から抽出する。機械学習制御部127は、 m セットの訓練データをモデル生成部124に提供し、 m 個のモデルをモデル生成部124から取得する。機械学習制御部127は、 m 個のモデルと1セットのテストデータをテスト実行部125に提供し、誤差プロファイルをテスト実行部125から取得する。そして、機械学習制御部127は、誤差プロファイルをテストサイズ決定部126に提供し、分散関数のパラメータをテストサイズ決定部126から取得し、制御情報として制御情報記憶部122に格納する。

20

【0165】

次に、機械学習制御部127は、モデル生成部124に対象のモデルを学習させる。機械学習制御部127は、先に特定したサイズの訓練データを、データ記憶部121に記憶されたデータ集合から抽出する。機械学習制御部127は、抽出した訓練データをモデル生成部124に提供し、学習されたモデルをテスト実行部125から取得する。機械学習制御部127は、モデルを学習結果記憶部123に格納する。

【0166】

次に、機械学習制御部127は、対象のモデルの予測性能を測定するための適切なテストデータサイズをテストサイズ決定部126に決定させる。まず、機械学習制御部127は、学習結果記憶部123に記憶されたモデルと、分散関数のパラメータの決定の際に使用した小さなテストデータとを、テスト実行部125に提供する。機械学習制御部127は、このために小さなテストデータを保存しておいてもよい。また、機械学習制御部127は、分散関数のパラメータの決定の際に使用したテストデータに代えて、同等のサイズのテストデータを、データ記憶部121に記憶されたデータ集合から抽出してもよい。

30

【0167】

機械学習制御部127は、テスト実行部125から期待ロスを取得し、制御情報記憶部122に記憶された分散関数のパラメータと期待ロスをテストサイズ決定部126に提供する。ただし、機械学習制御部127は、対象のモデルを用いて期待ロスを測定する代わりに、回帰分析などの統計的方法によって期待ロスを推定してもよい。機械学習制御部127は、テストサイズ決定部126からテストデータサイズを取得する。

40

【0168】

すると、機械学習制御部127は、データ記憶部121に記憶されたデータ集合から、決定されたサイズのテストデータを抽出する。テストデータに含まれるレコードは訓練データと重複しないことが好ましい。機械学習制御部127は、抽出したテストデータと学習結果記憶部123に記憶されたモデルとをテスト実行部125に提供する。機械学習制御部127は、テスト実行部125から予測性能の測定値を取得し、学習結果記憶部123に格納する。ただし、機械学習制御部127は、上記のテストデータに対して、更新された期待ロスをテスト実行部125から取得し、更新された期待ロスに基づいて、更新されたテストデータサイズをテストサイズ決定部126から取得してもよい。

50

【 0 1 6 9 】

モデルの学習と予測性能の測定が完了すると、機械学習制御部 1 2 7 は、モデルおよび予測性能の測定値を出力する。例えば、機械学習制御部 1 2 7 は、表示装置 1 1 1 にモデルおよび予測性能の測定値を表示する。機械学習制御部 1 2 7 は、他の出力デバイスにモデルおよび予測性能の測定値を出力してもよい。また、例えば、機械学習制御部 1 2 7 は、他の情報処理装置にモデルおよび予測性能の測定値を送信する。

【 0 1 7 0 】

図 7 は、誤差プロファイルテーブルの例を示す図である。

誤差プロファイルテーブル 1 3 1 は、制御情報記憶部 1 2 2 に記憶される。誤差プロファイルテーブル 1 3 1 は、 m セットの訓練データと n レコードのテストデータとの間で網羅的に算出された $m \times n$ 個の誤差を記憶する。誤差プロファイルテーブル 1 3 1 の列は、訓練データ D_1, D_2, \dots, D_m に対応する。誤差プロファイルテーブル 1 3 1 の行は、テストデータの n 個のレコードに含まれる入力値 X_1, X_2, \dots, X_n に対応する。1 つの訓練データ D_k から学習された 1 つのモデルに、テストデータの 1 つのレコードに含まれる入力値 X_i を入力することで、予測値と真値との差である誤差 e_{ik} が算出される。

10

【 0 1 7 1 】

図 8 は、分散関数テーブルの例を示す図である。

分散関数テーブル 1 3 2 は、制御情報記憶部 1 2 2 に記憶される。分散関数テーブル 1 3 2 は、尖度 K_1 、不動点 $Cor 1 v$ および期待バイアス $E B 2$ の 3 つのパラメータに対応する値を記憶する。これら 3 つのパラメータは、数式 (4 4) に含まれるパラメータであって、訓練データサイズに依存しないパラメータである。分散関数テーブル 1 3 2 に記憶される値は、誤差プロファイルテーブル 1 3 1 から算出される。

20

【 0 1 7 2 】

図 9 は、機械学習の手順例を示すフローチャートである。

(S 1 0) 機械学習制御部 1 2 7 は、機械学習アルゴリズムと訓練データサイズを指定する。機械学習アルゴリズムと訓練データサイズの指定はユーザから受け付けてもよい。

【 0 1 7 3 】

(S 1 1) 機械学習制御部 1 2 7 は、データ記憶部 1 2 1 から m セットの小さいサイズの訓練データと 1 セットの小さいサイズのテストデータを抽出する。例えば、1 万レコードの訓練データが 1 0 セット抽出され、5 0 0 0 レコードのテストデータが 1 セット抽出される。

30

【 0 1 7 4 】

(S 1 2) モデル生成部 1 2 4 は、ステップ S 1 0 で指定された機械学習アルゴリズムに従って、 m セットの訓練データから m 個のモデルを学習する。

(S 1 3) テスト実行部 1 2 5 は、ステップ S 1 2 で学習された m 個のモデルに、ステップ S 1 1 のテストデータの各レコードを入力して誤差を算出し、算出した誤差を列挙した誤差プロファイルテーブル 1 3 1 を生成する。具体的には、テスト実行部 1 2 5 は、1 つのモデルとテストデータの 1 つのレコードの組毎に、レコードに含まれる説明変数の値をモデルに入力し、モデルから算出された目的変数の予測値とレコードに含まれる真値との差を誤差として算出する。例えば、1 0 個のモデルと 5 0 0 0 レコードのテストデータから、1 0 \times 5 0 0 0 個の誤差を含む誤差プロファイルテーブル 1 3 1 が生成される。

40

【 0 1 7 5 】

(S 1 4) テストサイズ決定部 1 2 6 は、誤差プロファイルテーブル 1 3 1 から、所定の数式に従って、分散関数 $f(n, EL)$ を規定するパラメータの値を決定する。パラメータには、尖度 K_1 と不動点 $Cor 1 v$ と期待バイアス $E B 2$ が含まれる。ここで決定されるパラメータの値は、使用するデータ集合と指定された機械学習アルゴリズムに依存するものである一方、訓練データサイズに依存しないものである。

【 0 1 7 6 】

(S 1 5) 機械学習制御部 1 2 7 は、データ記憶部 1 2 1 から、ステップ S 1 0 で指定されたサイズの訓練データを抽出する。

(S 1 6) モデル生成部 1 2 4 は、ステップ S 1 0 で指定された機械学習アルゴリズム

50

に従って、ステップ S 1 5 で抽出された訓練データからモデルを学習する。

【 0 1 7 7 】

図 1 0 は、機械学習の手順例を示すフローチャート（続き）である。

（ S 1 7 ）テスト実行部 1 2 5 は、ステップ S 1 6 で学習されたモデルに、ステップ S 1 1 で抽出された小サイズのテストデータの各レコードを入力して誤差を算出する。ただし、ステップ S 1 1 で抽出されたものとは異なるテストデータを使用してもよい。

【 0 1 7 8 】

（ S 1 8 ）テスト実行部 1 2 5 は、ステップ S 1 7 で算出された誤差から、所定の数式に従って、ステップ S 1 6 で学習されたモデルの期待ロス $E L$ を推定する。

（ S 1 9 ）テストサイズ決定部 1 2 6 は、ステップ S 1 4 で決定されたパラメータの値をもつ分散関数 $f(n, EL)$ に、ステップ S 1 8 で推定された期待ロス $E L$ を代入する。テストサイズ決定部 1 2 6 は、分散関数 $f(n, EL)$ により算出される分散が所定条件を満たす範囲で、最大のテストデータサイズ n_1 を判定する。例えば、テストサイズ決定部 1 2 6 は、テストデータサイズ n を 2 倍にした場合の分散の減少率を示す効果指標の値と所定の閾値とを比較し、効果指標の値が閾値未満になるまでテストデータサイズ n を 2 倍にすることを繰り返す。これにより、最大のテストデータサイズ n_1 が選択される。

【 0 1 7 9 】

（ S 2 0 ）機械学習制御部 1 2 7 は、データ記憶部 1 2 1 から、ステップ S 1 9 で判定されたサイズ n_1 のテストデータを抽出する。

（ S 2 1 ）テスト実行部 1 2 5 は、ステップ S 1 6 で学習されたモデルに、ステップ S 2 0 で抽出されたテストデータの各レコードを入力して誤差を算出する。

【 0 1 8 0 】

（ S 2 2 ）テスト実行部 1 2 5 は、ステップ S 2 1 で算出された誤差から、所定の数式に従って、ステップ S 1 6 で学習されたモデルの期待ロス $E L$ を再推定する。

（ S 2 3 ）テストサイズ決定部 1 2 6 は、ステップ S 1 4 で決定されたパラメータの値をもつ分散関数 $f(n, EL)$ に、ステップ S 2 2 で再推定された期待ロス $E L$ を代入する。テストサイズ決定部 1 2 6 は、分散関数 $f(n, EL)$ により算出される分散が所定条件を満たす範囲で、最大のテストデータサイズ n_2 を判定する。テストデータサイズ n_2 の判定方法は、ステップ S 1 9 と同様の方法でよい。

【 0 1 8 1 】

（ S 2 4 ）機械学習制御部 1 2 7 は、データ記憶部 1 2 1 から、ステップ S 2 3 で判定されたサイズ n_2 のテストデータを抽出する。

（ S 2 5 ）テスト実行部 1 2 5 は、ステップ S 1 6 で学習されたモデルに、ステップ S 2 4 で抽出されたテストデータの各レコードを入力して誤差を算出する。テスト実行部 1 2 5 は、算出された誤差から、当該モデルの予測性能の測定値を算出する。

【 0 1 8 2 】

（ S 2 6 ）機械学習制御部 1 2 7 は、ステップ S 1 6 で学習されたモデルとステップ S 2 5 で算出された予測性能の測定値を、学習結果記憶部 1 2 3 に保存する。また、機械学習制御部 1 2 7 は、モデルおよび予測性能の測定値を表示装置 1 1 1 に表示する。

【 0 1 8 3 】

なお、上記のフローチャートでは、対象となるモデルの期待ロス $E L$ の推定を 2 回繰り返している。期待ロス $E L$ の推定を 1 回だけ行う場合、上記のステップ S 1 9 ~ S 2 2 を省略することができる。また、対象となるモデルを使用せずに統計的方法により期待ロス $E L$ を推定する場合、上記のステップ S 1 7 ~ S 2 2 を省略することができる。

【 0 1 8 4 】

第 2 の実施の形態の機械学習装置 1 0 0 によれば、複数セットの小さい訓練データと 1 セットの小さいテストデータを用いて、同一のデータ集合および機械学習アルゴリズムのもとで生じる誤差の分布を示す誤差プロファイルが生成される。誤差プロファイルに基づいて、期待ロスとテストデータサイズを引数としてもち、予測性能の測定値の分散を算出する分散関数が決定される。そして、大きい訓練データを用いて学習された対象モデルの

10

20

30

40

50

期待ロスが推定され、分散関数が示すテストデータサイズと分散の対応関係に基づいて、対象モデルの予測性能を測定するための適切なテストデータサイズが決定される。

【0185】

テストデータサイズは、予測性能の測定値が実用上十分な信頼性をもつ範囲、すなわち、その分散が許容できる範囲で、できる限り小さいサイズに決定される。これにより、テストデータサイズが小さ過ぎることにより予測性能の測定値の信頼性が低下することを抑制できる。また、テストデータサイズが大き過ぎることにより予測性能の測定値の信頼性向上に寄与しない無駄なテスト処理が発生することを抑制でき、テスト処理の負荷を軽減してテスト時間を短縮できる。よって、学習されたモデルの予測性能を高信頼かつ短時間で測定することができ、テスト処理を効率化することができる。例えば、テストデータサイズを訓練データサイズの2分の1から4分の1程度とする慣習的方法と比べて、測定値の分散を同程度に抑えつつ、テストデータサイズを削減することができる。

10

【符号の説明】

【0186】

- 10 機械学習装置
- 11 記憶部
- 12 処理部
- 13 データ集合
- 14 a, 14 b, 14 c, 18 訓練データ
- 15, 19 テストデータ
- 16 誤差情報
- 17 対応関係

20

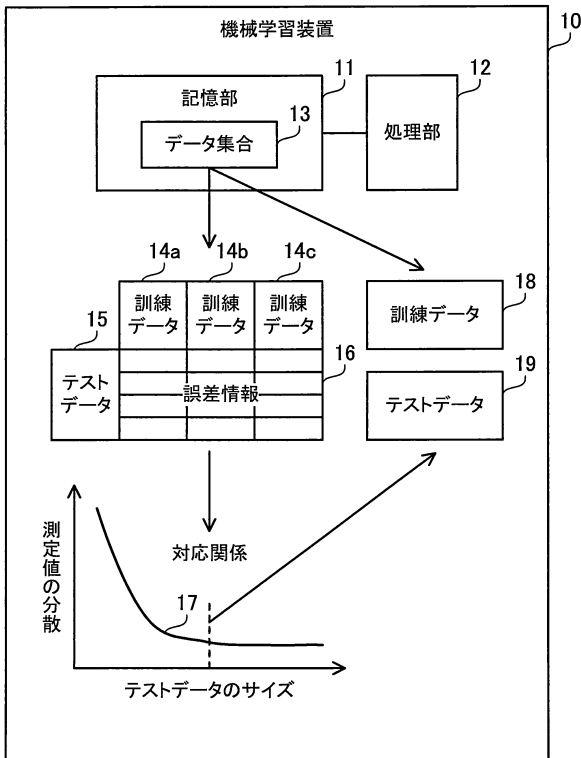
30

40

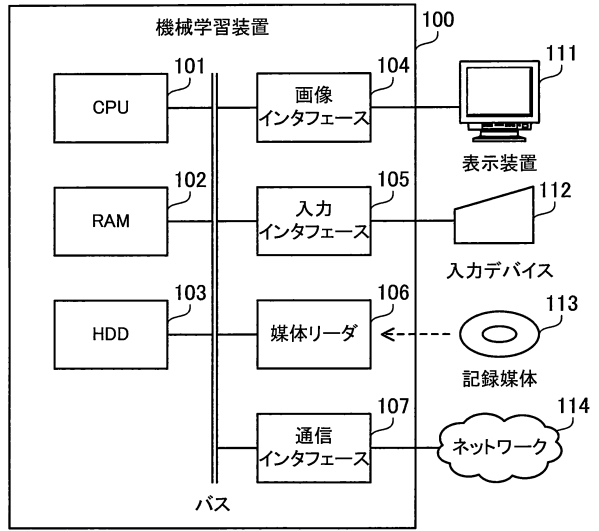
50

【図面】

【図 1】



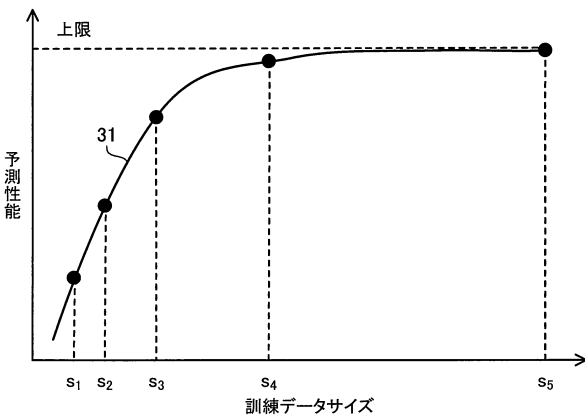
【図 2】



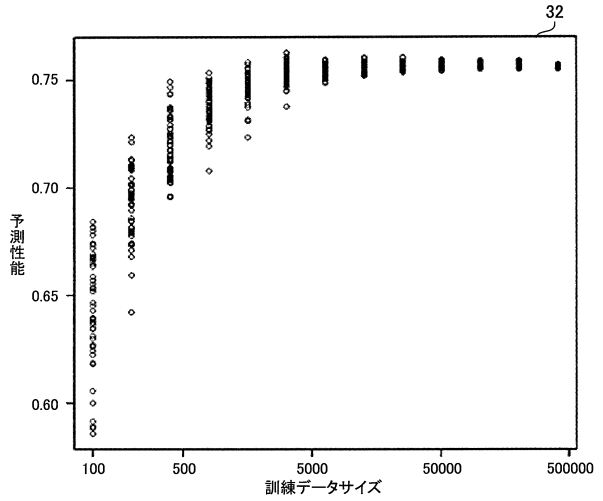
10

20

【図 3】



【図 4】

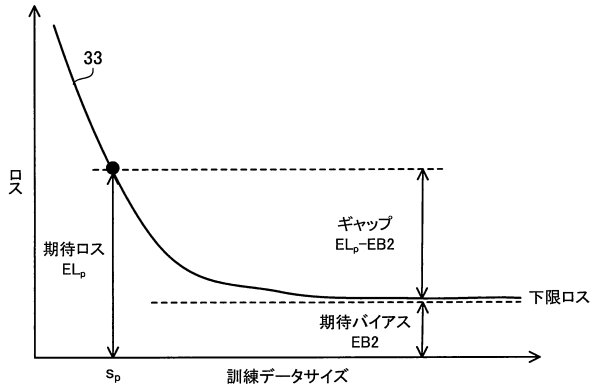


30

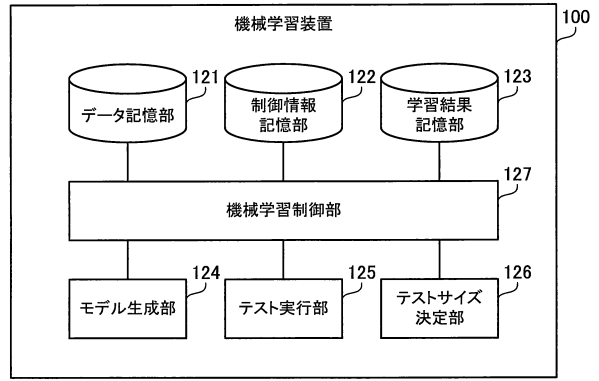
40

50

【図 5】



【図 6】



10

【図 7】

誤差プロファイルテーブル 131

		訓練データ					
		D_1	D_2	...	D_k	...	D_m
テストデータ	X_1	e_{11}	e_{12}	...	e_{1k}	...	e_{1m}
	X_2	e_{21}	e_{22}	...	e_{2k}	...	e_{2m}
	⋮	⋮	⋮		⋮		⋮
	X_i	e_{i1}	e_{i2}	...	e_{ik}	...	e_{im}
	⋮	⋮	⋮		⋮		⋮
	X_n	e_{n1}	e_{n2}	...	e_{nk}	...	e_{nm}

【図 8】

分散関数テーブル 132

パラメータ	値
尖度 $K1$...
不動点 $Cor1v$...
期待バイアス $EB2$...

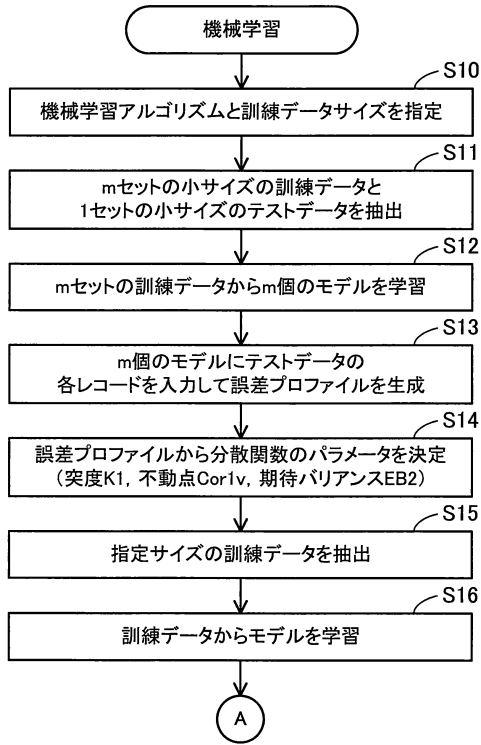
20

30

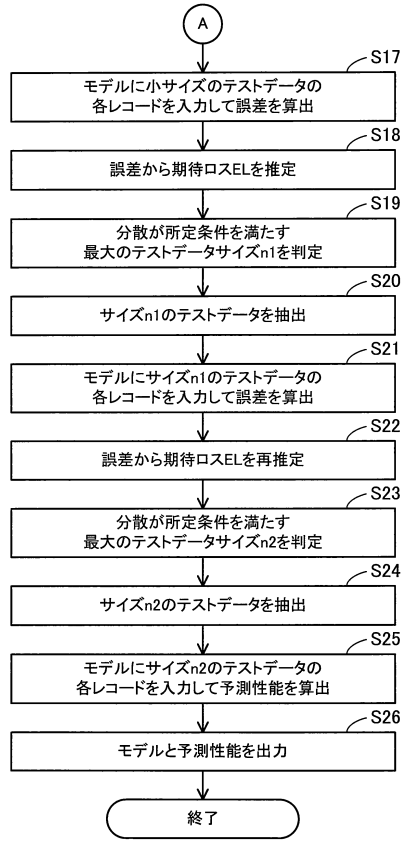
40

50

【 図 9 】



【 図 1 0 】



10

20

30

40

50

フロントページの続き

- (56)参考文献 特表 2 0 1 5 - 5 2 5 4 1 3 (J P , A)
特開平 9 - 5 4 7 6 4 (J P , A)
特開 2 0 1 9 - 1 1 3 9 1 5 (J P , A)
特開 2 0 1 7 - 4 9 6 7 4 (J P , A)

- (58)調査した分野 (Int.Cl. , D B 名)
G 0 6 F 8 / 0 0 - 8 / 3 8
8 / 6 0 - 8 / 7 7
9 / 4 4 - 9 / 4 4 5
9 / 4 5 1
G 0 6 N 3 / 0 0 - 3 / 1 2
7 / 0 8 - 9 9 / 0 0