

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6253644号
(P6253644)

(45) 発行日 平成29年12月27日 (2017.12.27)

(24) 登録日 平成29年12月8日 (2017.12.8)

(51) Int. Cl.		F I	
G06F 17/30	(2006.01)	G06F 17/30	210D
G06F 19/24	(2011.01)	G06F 17/30	170F
G06N 99/00	(2010.01)	G06F 19/24	
		G06N 99/00	153

請求項の数 15 (全 22 頁)

(21) 出願番号	特願2015-517783 (P2015-517783)	(73) 特許権者	500586875
(86) (22) 出願日	平成25年6月21日 (2013.6.21)		フィリップ モリス プロダクツ エス
(65) 公表番号	特表2015-525413 (P2015-525413A)		アー
(43) 公表日	平成27年9月3日 (2015.9.3)		スイス国 2000 ヌーシャテル ケ
(86) 国際出願番号	PCT/EP2013/062980		ジャンルノー 3
(87) 国際公開番号	W02013/190084	(74) 代理人	100078282
(87) 国際公開日	平成25年12月27日 (2013.12.27)		弁理士 山本 秀策
審査請求日	平成28年6月20日 (2016.6.20)	(74) 代理人	100113413
(31) 優先権主張番号	61/662,792		弁理士 森下 夏樹
(32) 優先日	平成24年6月21日 (2012.6.21)	(74) 代理人	100181674
(33) 優先権主張国	米国 (US)		弁理士 飯田 貴敏
		(74) 代理人	100181641
			弁理士 石川 大輔
		(74) 代理人	230113332
			弁護士 山本 健策

最終頁に続く

(54) 【発明の名称】 統合バイアス補正およびクラス予測を用いてバイオマーキングネチャを生成するためのシステムおよび方法

(57) 【特許請求の範囲】

【請求項 1】

プロセッサによって実行される、2つ以上のクラスにデータセットを分類するコンピュータ実装方法であって、前記方法は、

(a) トレーニングデータセットおよびトレーニングクラスセットを受信することであって、前記トレーニングクラスセットは、既知のラベルのセットを含み、各既知のラベルは、前記トレーニングデータセット中の各要素と関連付けられるクラスを識別する、ことと、

(b) テストデータセットを受信することと、

(c) 前記トレーニングデータセットおよび前記トレーニングクラスセットに第1の機械学習技法を適用することによって、前記トレーニングデータセットについての第1の分類器を生成することと、

(d) 前記第1の分類器に従って、前記テストデータセット中の要素を分類することによって、第1のテストクラスセットを生成することと、

(e) トレーニングクラス重心のセットの中心に対応する量だけ前記トレーニングデータセット中の前記要素を偏移させることによって、前記トレーニングデータセットを変換することであって、各トレーニングクラス重心は、前記トレーニングデータセット中の要素のサブセットの中心を表す、ことと、

(f) 複数の反復の各々について、

(i) テストクラス重心のセットの中心に対応する量だけ前記テストデータセット中

10

20

の前記要素を偏移させることによって、前記テストデータセットを変換することによって、各テストクラス重心は、前記テストデータセット中の要素のサブセットの中心を表す、こと、

(i i) 第2の分類器に従って前記変換されたテストデータセット中の前記要素を分類することによって、第2のテストクラスセットを生成することによって、前記第2の分類器は、前記変換されたトレーニングデータセットおよび前記トレーニングクラスセットに第2の機械学習技法を適用することによって生成される、こと、および、

(i i i) 前記第1のテストクラスセットと前記第2のテストクラスセットとが異なる場合、前記第2のテストクラスセットを前記第1のテストクラスセットとして記憶し、前記変換されたテストデータセットを前記テストデータセットとして記憶し、ステップ(i)に戻ること
を行うこと

10

を含む、方法。

【請求項2】

前記第1のテストクラスセットと前記第2のテストクラスセットとが異なる場合に、前記第2のテストクラスセットを出力することをさらに含む、請求項1に記載の方法。

【請求項3】

前記トレーニングデータセットの前記要素は、疾患を有する患者についての、前記疾患に耐性がある患者についての、または、前記疾患がない患者についての遺伝子発現データを表す、請求項1～2のいずれかに記載の方法。

20

【請求項4】

前記トレーニングデータセットは、集約データセット中のサンプルのランダムなサブセットから形成され、前記テストデータセットは、前記集約データセット中のサンプルの残っているサブセットから形成される、請求項1～3のいずれかに記載の方法。

【請求項5】

ステップ(e)における前記偏移させることは、前記変換されたトレーニングデータセットを取得するように、前記トレーニングデータセットに回転、シアー、線形変換、または、非線形変換を適用することを含む、請求項1～4のいずれかに記載の方法。

【請求項6】

ステップ(i)における前記偏移させることは、前記変換されたテストデータセットを取得するように、前記テストデータセットに回転、シアー、線形変換、または、非線形変換を適用することを含む、請求項1～5のいずれかに記載の方法。

30

【請求項7】

前記テストデータセットは、既知のラベルのテストセットを含み、各既知のラベルは、前記テストデータセット中の各要素と関連付けられるクラスを識別し、
前記第1のテストクラスセットは、前記テストデータセットについての予測されるラベルのセットを含み、
前記第2のテストクラスセットは、前記変換されたテストデータセットについての予測されるラベルのセットを含む、
請求項1～6のいずれかに記載の方法。

40

【請求項8】

前記複数の反復の各々について、前記第1のテストクラスセットを前記第2のテストクラスセットと比較することをさらに含む、請求項1～7のいずれかに記載の方法。

【請求項9】

前記第1の機械学習技法と前記第2の機械学習技法は同一である、請求項1～8のいずれかに記載の方法。

【請求項10】

ステップ(e)での前記変換は、ステップ(i)の同一の変換を適用することによって行われる、請求項1～9のいずれかに記載の方法。

【請求項11】

50

前記第 2 のテストクラスセットを表示デバイス、印刷デバイス、または、記憶デバイスに提供することをさらに含む、請求項 1 ~ 10 のいずれかに記載の方法。

【請求項 12】

前記第 1 のテストクラスセットおよび前記第 2 のテストクラスセットは、前記第 1 のテストクラスセットの任意の要素が前記第 2 のテストクラスセットの対応する要素と異なる場合に、異なる、請求項 1 ~ 11 のいずれかに記載の方法。

【請求項 13】

前記第 2 のテストクラスセットは、前記変換されたテストデータセットについての予測されるラベルのセットを含み、前記方法は、予測されるラベルの総数によって除算された前記第 2 のテストクラスセット中の正確な予測されるラベルの数を表す性能測定基準を計算することによって、前記第 2 の分類器を評価することをさらに含む、請求項 1 ~ 12 のいずれかに記載の方法。

10

【請求項 14】

コンピュータ可読命令を記憶したコンピュータ可読記憶媒体であって、前記コンピュータ可読命令は、少なくとも 1 つのプロセッサを備えるコンピュータ化システムにおいて実行される場合、前記少なくとも 1 つのプロセッサに請求項 1 ~ 13 のいずれかに記載の方法を実行させる、コンピュータ可読記憶媒体。

【請求項 15】

非一時的なコンピュータ可読命令を伴って構成される少なくとも 1 つのプロセッサを備えるコンピュータ化システムであって、前記非一時的なコンピュータ可読命令は、実行される場合、前記プロセッサに請求項 1 ~ 13 のいずれかに記載の方法を実行させる、コンピュータ化システム。

20

【発明の詳細な説明】

【背景技術】

【0001】

関連出願への参照

本願は、米国仮特許出願第 61 / 662 , 792 号 (発明の名称「Systems and Methods for Generating Biomarker Signatures with Integrated Bias Correction and Class Prediction」、2012年6月21日出願) に対する 35 U.S.C. § 119 の下での優先権を主張し、それは、本明細書にその全体が援用される。

30

【0002】

背景

生物医学分野において、特定の生物学的状態を示す物質、すなわち、バイオマーカを識別することが重要である。ゲノミクスおよびプロテオミクスの新しい技術が出現するにつれて、バイオマーカは、生物学的発見、薬剤開発、および、ヘルスケアにおいてますます重要になりつつある。バイオマーカは、多くの疾患の診断および予後のためだけでなく、治療法の開発のための基礎を理解するためにも有用である。バイオマーカの成功した効果的な識別は、新薬開発プロセスを加速させることができる。診断および予後と治療法との組み合わせによって、バイオマーカ識別はまた、現在の薬物治療の品質を向上し、したがって、薬理遺伝学、薬理ゲノム学、および、薬理プロテオミクスの使用において重要な役割を果たす。

40

【0003】

高スループットスクリーニングを含むゲノムおよびプロテオームの分析は、細胞において発現させられるタンパク質の数および形態に関する豊富な情報を供給し、各細胞について、特定の細胞状態の特性を示す発現させられたタンパク質のプロファイルを識別する潜在的な可能性を提供する。特定の場において、この細胞状態は、疾患と関連付けられる異常生理学的反応の特性を示し得る。結果として、疾患を有する患者からの細胞状態を識別し、それを正常な患者からの対応する細胞の細胞状態と比較することによって、疾患を診

50

断して治療する機会を提供することができる。

【0004】

これらの高スループットスクリーニング技法は、遺伝子発現情報の大量のデータセットを提供する。研究者らは、個人の多様な集団について再現可能に診断するパターンにこれらのデータセットを組織化するための方法を開発しようとしてきた。1つのアプローチは、複合データセットを形成するように複数のソースからのデータをプールし、次いで、データセットを発見/トレーニングセットおよびテスト/検証セットに分割することであった。しかしながら、転写プロファイリングデータおよびタンパク質発現プロファイリングデータは両方とも、しばしば、利用可能な数のサンプルに対する多数の変数によって特徴付けられる。

10

【0005】

患者または対照の群からの検体の発現プロファイルの間の観察された差異は、典型的に、疾患または対照の集団内の生物学的変動または未知のサブ表現型、研究プロトコルにおける差異による部位特異的なバイアス、検体の取り扱い、器具条件（例えば、チップバッチ等）における差異によるバイアス、および、測定誤差による変動を含むいくつかの要因によって、弱められる。いくつかの技法は、データサンプルにおけるバイアスを補正しようとする（例えば、別のクラスよりもむしろ、データセットにおいて表されるサンプルの1つのクラスを有することに起因し得る）。

【0006】

いくつかのコンピュータベースの方法が、疾患および対照のサンプルの間の差異を最も良く説明する一組の特徴（マーカ）を見出すために開発されてきた。いくつかの初期の方法は、LIMMA、乳癌に関するバイオマーカを識別するためのFDA承認マンマプリント技法、ロジスティック回帰技法、および、サポートベクトルマシン（SVM）等の機械学習方法のような統計的テストを含んでいた。概して、機械学習の視点から、バイオマーカの選択は、典型的に、分類タスクについての特徴選択問題である。しかしながら、これらの初期の解決策は、いくつかの不利点に直面した。これらの技法によって生成されるシグネチャは、しばしば、対象の包含および除外が異なるシグネチャにつながり得るので、再現可能ではなかった。これらの初期の解決策はまた、多くの偽陽性シグネチャを生成し、小サンプルサイズおよび高次元を有するデータセットに作用するので、ロバストではなかった。

20

30

【0007】

したがって、臨床的な診断および/または予後についてのバイオマーカを識別するため、より一般的には、データセットの中の要素を2つ以上のクラスに分類するために使用されることができるデータマーカを識別するための改良型技法の必要性がある。

【発明の概要】

【課題を解決するための手段】

【0008】

出願人らは、既存のコンピュータベースの方法が、クラス予測技法とは別にバイアス補正技法を不利に適用することを認識している。本明細書で説明されるコンピュータシステムおよびコンピュータプログラム製品は、バイオマーカおよび他のデータ分類適用において改善された分類性能を達成し得る、バイアス補正およびクラス予測への統合アプローチを適用する方法を実装する。特定すると、本明細書で開示されるコンピュータ実装方法は、バイアス補正およびクラス予測への反復アプローチを採用する。コンピュータ実装方法の種々の実施形態において、システム中の少なくとも1つのプロセッサが、トレーニングデータセットおよびトレーニングクラスセットを受信し、そのトレーニングクラスセットは、トレーニングデータセットの中の要素の各々と関連付けられるクラスを識別する。システム中のプロセッサはまた、テストデータセットを受信する。プロセッサは、機械学習技法をトレーニングデータセットおよびトレーニングクラスセットに適用することによって、トレーニングデータセットについての第1の分類器を生成し、第1の分類器に従ってテストデータセット中の要素を分類することによって、第1のテストクラスセットを生成

40

50

する。複数の反復の各々について、プロセッサは、トレーニングクラスセットとテストクラスセットとのうちの少なくとも1つに基づいて、トレーニングデータセットを変換し、以前のステップの変換を適用することによって、テストデータセットを変換し、変換されたトレーニングデータセットおよびトレーニングクラスセットに機械学習技法を適用することによって、変換されたトレーニングデータセットについての第2の分類器を生成し、第2の分類器に従って、変換されたテストデータセット中の要素を分類することによって、第2のテストクラスセットを生成する。プロセッサはまた、第1のテストクラスセットと第2のテストクラスセットとを比較し、第1のテストクラスセットと第2のテストクラスセットとが異なる場合、プロセッサは、第2のクラスセットを第1のクラスセットとして記憶し、変換されたテストデータセットをテストデータセットとして記憶し、反復の開始に戻る。本発明のコンピュータシステムは、上記で説明されるような方法およびその種々の実施形態を実装するための手段を備える。

例えば、本発明は、下記の項目を提供する。

(項目1)

プロセッサによって実行される、2つ以上のクラスにデータセットを分類するコンピュータ実装方法であって、前記方法は、

(a) トレーニングデータセットおよびトレーニングクラスセットを受信するステップであって、前記トレーニングクラスセットは、前記トレーニングデータセット中の要素の各々と関連付けられるクラスを識別する、ステップと、

(b) テストデータセットを受信するステップと、

(c) 前記トレーニングデータセットおよび前記トレーニングクラスセットに機械学習技法を適用することによって、前記トレーニングデータセットについての第1の分類器を生成するステップと、

(d) 前記第1の分類器に従って、前記テストデータセット中の要素を分類することによって、第1のテストクラスセットを生成するステップと、

(e) 複数の反復の各々について、

(i) 前記トレーニングクラスセットおよび前記テストクラスセットのうちの少なくとも1つに基づいて、前記トレーニングデータセットを変換するステップと、

(i i) 前記テストデータセットを変換するステップと、

(i i i) 前記変換されたトレーニングデータセットおよび前記トレーニングクラスセットに基づいて、第2の分類器に従って前記変換されたテストデータセット中の前記要素を分類することによって、第2のテストクラスセットを生成するステップと、

(i v) 前記第1のテストクラスセットと前記第2のテストクラスセットとが異なる場合、前記第2のクラスセットを前記第1のクラスセットとして記憶し、前記変換されたテストデータセットを前記テストデータセットとして記憶し、ステップ(i)に戻るステップと

を含む、方法。

(項目2)

前記第1のテストクラスセットと前記第2のテストクラスセットとが異なる場合に、前記第2のクラスセットを出力するステップをさらに含む、項目1に記載の方法。

(項目3)

前記トレーニングデータセットの要素は、疾患を有する患者についての、前記疾患に耐性がある患者についての、または、前記疾患がない患者についての遺伝子発現データを表す、項目1~2のいずれかに記載の方法。

(項目4)

前記トレーニングデータセットおよび前記テストデータセットは、集約データセット中のサンプルを前記トレーニングデータセットまたは前記テストデータセットにランダムに割り当てることによって、生成される、項目1~3のいずれかに記載の方法。

(項目5)

ステップ(i)、ステップ(i i)、または、ステップ(i)とステップ(i i)との両

10

20

30

40

50

方の前記変換は、前記データセットの重心に基づいて前記データセットの要素を調整することによって、バイアス補正技法を行うステップを含む、項目 1 ~ 4 のいずれかに記載の方法。

(項目 6)

前記バイアス補正技法は、前記データセットの各要素から前記重心の成分を差し引くステップを含む、項目 5 に記載の方法。

(項目 7)

ステップ (i)、ステップ (i i)、または、ステップ (i) とステップ (i i) との両方における前記変換は、回転、シアー、線形変換、または、非線形変換を適用するステップを含む、項目 1 ~ 6 のいずれかに記載の方法。

10

(項目 8)

前記複数の反復の各々について、前記第 1 のテストクラスセットを前記第 2 のテストクラスセットと比較するステップをさらに含む、項目 1 ~ 7 のいずれかに記載の方法。

(項目 9)

前記複数の反復の各々について、前記変換されたトレーニングデータセットおよび前記トレーニングデータセットに機械学習技法を適用することによって、前記変換されたトレーニングデータセットについての前記第 2 の分類器を生成するステップをさらに含む、項目 1 ~ 8 のいずれかに記載の方法。

(項目 10)

ステップ (i i) での前記変換は、ステップ (i) の同一の変換を適用することによって行われる、項目 1 ~ 9 のいずれかに記載の方法。

20

(項目 11)

前記第 2 のテストクラスセットを表示デバイス、印刷デバイス、または、記憶デバイスに提供するステップをさらに含む、項目 1 ~ 10 のいずれかに記載の方法。

(項目 12)

前記第 1 のテストクラスセットおよび前記第 2 のテストクラスセットは、前記第 1 のテストクラスセットの任意の要素が前記第 2 のテストクラスセットの対応する要素と異なる場合に、異なる、項目 1 ~ 11 のいずれかに記載の方法。

(項目 13)

エラー率に基づいて、前記第 2 の分類器の性能測定基準を計算するステップをさらに含む、項目 1 ~ 12 のいずれかに記載の方法。

30

(項目 14)

コンピュータ可読命令を備えるコンピュータプログラム製品であって、前記コンピュータ可読命令は、少なくとも 1 つのプロセッサを備えるコンピュータ化システムにおいて実行される場合、前記少なくとも 1 つのプロセッサに項目 1 ~ 13 のいずれかに記載の方法の 1 つ以上のステップを実行させる、コンピュータプログラム製品。

(項目 15)

非一時的なコンピュータ可読命令を伴って構成される少なくとも 1 つのプロセッサを備えるコンピュータ化システムであって、前記非一時的なコンピュータ可読命令は、実行される場合、前記プロセッサに項目 1 ~ 13 のいずれかに記載の方法を実行させる非一時的なコンピュータ可読命令を伴って構成される少なくとも 1 つのプロセッサを備える、コンピュータ化システム。

40

【0009】

上記で説明される方法の特定の実施形態において、本方法はさらに、第 1 のテストクラスセットと第 2 のテストクラスセットとが異なる場合に、第 2 のクラスセットを出力するステップを含む。特定すると、上記で説明されるような反復は、第 1 のテストクラスセットおよび第 2 のテストクラスセットが収束し、かつ、予測された分類の間に差異がなくなるまで、繰り返され得る。上記で説明される方法の特定の実施形態において、トレーニングデータセットの要素は、疾患を有する患者について、疾患に耐性がある患者について、または、疾患がない患者についての遺伝子発現データを表す。トレーニングクラスセ

50

ットの要素は、トレーニングデータセット中のデータサンプルについての既知のクラス識別子に対応し得る。例えば、クラス識別子は、「疾患陽性」、「疾患免疫性」、または、「疾患なし」等のカテゴリを含み得る。

【0010】

上記で説明される方法の特定の実施形態において、トレーニングデータセットおよびテストデータセットは、集約データセット中のサンプルをトレーニングデータセットまたはテストデータセットにランダムに割り当てることによって、生成される。集約データセットをトレーニングデータセットとテストデータセットとにランダムに分割することが、クラスを予測してロバストな遺伝子シグネチャを生成するために望ましくあり得る。さらに、集約データセットのサンプルは、分割の前に破棄され得るか、または、トレーニングデータセットあるいはテストデータセットのサンプルは、分割後に破棄され得る。上記で説明される方法の特定の実施形態において、トレーニングデータセットを変換するステップ、テストデータセットを変換するステップ、または、トレーニングデータセットを変換するステップとテストデータセットを変換するステップとの両方は、データセットの重心に基づいてデータセットの要素を調整することによって、バイアス補正技法を行うステップを含む。変換は、トレーニングクラスセットに基づいて変換を定義し得る変換関数に従って行われる。上記で説明される方法の特定の実施形態において、バイアス補正技法は、データセットの各要素から重心の成分を差し引くステップを含む。例えば、バイアス補正技法の結果は、データセットにおいて表される各クラスの重心を考慮することによって、トレーニングデータセット、テストデータセット、または、トレーニングデータセットおよびテストデータセットの両方の各要素が、「再び中心に置かれる」ことであり得る。上記で説明される方法の特定の実施形態において、トレーニングデータセットを変換するステップ、テストデータセットを変換するステップ、または、トレーニングデータセットを変換するステップとテストデータセットを変換するステップとの両方は、回転、シアア (shear)、線形変換、または、非線形変換を適用するステップを含む。

【0011】

上記で説明される方法の特定の実施形態において、本方法はさらに、複数の反復の各々について、第1のテストクラスセットと第2のテストクラスセットとを比較するステップを含む。比較の結果として、第1のテストクラスセットおよび第2のテストクラスセットは、第1のテストクラスセットの任意の単一の要素が第2のテストクラスセットの対応する要素とは異なる場合に、異なると言われ得る。概して、第1のテストクラスセット中の少なくとも所定の数の要素が第2のテストクラスセット中の対応する要素と異なる場合に、第1のテストクラスセットと第2のテストクラスセットとが異なると言われ得るように、閾値が設定され得る。

【0012】

上記で説明される方法の特定の実施形態において、本方法はさらに、複数の反復の各々について、変換されたトレーニングデータセットおよびトレーニングデータセットに機械学習技法を適用することによって、変換されたトレーニングデータセットについての第2の分類器を生成するステップを含む。上記で説明される方法の特定の実施形態において、テストデータセットの変換は、トレーニングデータセットを変換するステップの変換と同一の変換を伴う。上記で説明される方法の特定の実施形態において、本方法はさらに、表示デバイス、印刷デバイス、または、記憶デバイスに第2のテストクラスセットを提供するステップを含む。上記で説明される方法の特定の実施形態において、本方法はさらに、エラー率に基づいて、第2の分類器の性能測定基準を計算するステップを含む。特定の実施形態において、限定されないが、線形判別分析 (LDA)、ロジスティック回帰、サポートベクトルマシン、ナイーブベイズ分類器等の線形分類器が好ましい。

【0013】

本発明のコンピュータシステムは、上記で説明されるような方法の種々の実施形態を実装するための手段を備える。例えば、コンピュータプログラム製品が説明され、本製品は、少なくとも1つのプロセッサを備えるコンピュータ化システムにおいて実行される場合

10

20

30

40

50

、上記で説明される方法のうちのいずれかの1つ以上のステップをプロセッサに実行させるコンピュータ可読命令を備える。別の例において、コンピュータ化システムが説明され、本システムは、実行される場合、上記で説明される方法のうちのいずれかをプロセッサに実行させる非一時的なコンピュータ可読命令を伴って構成されるプロセッサを備える。本明細書で説明されるコンピュータプログラム製品およびコンピュータ化方法は、1つ以上のプロセッサを各々が含む1つ以上のコンピューティングデバイスを有するコンピュータ化システムにおいて実装され得る。概して、本明細書で説明されるコンピュータ化システムは、本明細書で説明されるコンピュータ化方法のうちの1つ以上を実行するようにハードウェア、ファームウェア、および、ソフトウェアを伴って構成されるコンピュータ、マイクロプロセッサ、論理デバイス、または、他のデバイスもしくはプロセッサ等の、プロセッサまたはデバイスを含む1つ以上のエンジンを備え得る。これらのエンジンのうちのいずれか1つ以上は、いずれか1つ以上の他のエンジンから物理的に分離可能であり得るか、または、共通のまたは異なる回路基板上の別個のプロセッサ等の、複数の物理的に分離可能な構成要素を含み得る。本発明のコンピュータシステムは、上記で説明されるような方法およびその種々の実施形態を実装するための手段を備える。エンジンは、随時、相互接続され得、さらに、随時、摂動データベース、測定可能値データベース、実験データのデータベース、および、文献データベースを含む1つ以上のデータベースに接続され得る。本明細書で説明されるコンピュータ化システムは、ネットワークインターフェースを通して通信する1つ以上のプロセッサおよびエンジンを有する分散型コンピュータ化システムを含み得る。そのような実装は、複数の通信システムにわたる分散型計算のために適切であり得る。

10

20

【図面の簡単な説明】

【0014】

本開示のさらなる特徴、その性質、および、種々の利点は、類似参照文字が全体を通して類似部分を指す添付図面と関連して検討される下記の詳細な説明を考慮すると明白になる。

【0015】

【図1】図1は、1つ以上のバイオマーカシグネチャを識別するための例示的なシステムを描写する。

【図2】図2は、データセット中の要素の分類を図示する。

30

【図3 - 1】図3は、データセットを分類するための例示的なプロセスの流れ図である。

【図3 - 2】図3は、データセットを分類するための例示的なプロセスの流れ図である。

【図4】図4は、図1のシステムの構成要素のうちのいずれか等のコンピューティングデバイスのブロック図である。

【図5】図5は、トレーニングデータセット中の遺伝子シグネチャのヒートマップである。

【発明を実施するための形態】

【0016】

本明細書で説明されるシステムおよび方法の全体的な理解を提供するために、ここで、遺伝子バイオマーカシグネチャを識別するためのシステムおよび方法を含む特定の例証の実施形態が、説明される。しかしながら、本明細書で説明されるシステム、コンピュータプログラム製品、および、方法は、任意のデータ分類適用等の他の好適な適用のために適合させられかつ修正され得、そのような他の追加および修正は、その範囲から逸脱しないことが、当業者によって理解される。概して、本明細書で説明されるコンピュータ化システムは、本明細書で説明されるコンピュータ化方法のうちの1つ以上を実行するようにハードウェア、ファームウェア、および、ソフトウェアを伴って構成されるコンピュータ、マイクロプロセッサ、または、論理デバイス等の1つ以上のエンジン、プロセッサ、または、デバイスを備え得る。

40

【0017】

図1は、本明細書で開示される分類技法が実装され得る、1つ以上のバイオマーカシグ

50

ネチャを識別するための例示的なシステム100を描写する。システム100は、バイオマーカジェネレータ102と、バイオマーカコンソリデータ104とを含む。システム100はさらに、バイオマーカジェネレータ102およびバイオマーカコンソリデータ104の動作の特定の局面を制御するための中央制御装置(CCU)101を含む。動作中に、遺伝子発現データ等のデータが、バイオマーカジェネレータ102で受信される。バイオマーカジェネレータ102は、複数の候補バイオマーカおよび対応するエラー率を生成するようにデータを処理する。バイオマーカコンソリデータ104は、これらの候補バイオマーカおよびエラー率を受信し、最適な性能尺度およびサイズを有する好適なバイオマーカを選択する。

【0018】

バイオマーカジェネレータ102は、データを処理して一組の候補バイオマーカおよび候補エラー率を生成するためのいくつかの構成要素を含む。特定すると、バイオマーカジェネレータは、データをトレーニングデータセットとテストデータセットとに分割するためのデータ前処理エンジン110を含む。バイオマーカジェネレータ102は、トレーニングデータセットおよびテストデータセットを受信してテストデータセットの要素を2つ以上のクラス(例えば、罹患および非罹患、感染しやすい、および、免疫がある、および、罹患等)のうちの1つに分類するための分類エンジン114を含む。バイオマーカジェネレータ102は、データ前処理エンジン110によって選択されるテストデータに適用される場合の分類器の性能を決定するための分類器性能監視エンジン116を含む。分類器性能監視エンジン116は、分類器(例えば、分類にとって最も重要であるデータセットの要素の成分)に基づいて候補バイオマーカを識別し、1つ以上の候補バイオマーカについて、候補エラー率を含み得る性能尺度を生成する。バイオマーカジェネレータ102はさらに、1つ以上の候補バイオマーカおよび候補性能尺度を記憶するためのバイオマーカ記憶部118を含む。

【0019】

バイオマーカジェネレータは、自動的に制御またはユーザ操作され得るCCU 101によって制御され得る。特定の実施形態において、バイオマーカジェネレータ102は、データをトレーニングデータセットとテストデータセットとにランダムに分割する度に、複数の候補バイオマーカを生成するように動作し得る。そのような複数の候補バイオマーカを生成するために、バイオマーカジェネレータ102の動作は、複数回、反復され得る。CCU 101は、所望の数の候補バイオマーカを含む1つ以上のシステム反復パラメータを受信し得、それらは、次に、バイオマーカジェネレータ102の動作が反復され得る回数を決定するように使用され得る。CCU 101はまた、バイオマーカ中の構成要素の数(例えば、バイオマーカ遺伝子シグネチャ中の遺伝子の数)を表し得る所望のバイオマーカサイズを含む他のシステムパラメータを受信し得る。バイオマーカサイズ情報は、トレーニングデータから候補バイオマーカを生成するために分類器性能監視エンジン116によって使用され得る。バイオマーカジェネレータ102の動作、特定すると分類エンジン114の動作は、図2~4への参照によってさらに詳細に説明される。

【0020】

バイオマーカジェネレータ102は、1つ以上の候補バイオマーカおよび候補エラー率を生成し、それらは、ロバストなバイオマーカを生成するためにバイオマーカコンソリデータ104によって使用される。バイオマーカコンソリデータ104は、複数の候補バイオマーカを受信して複数の候補バイオマーカにわたって最も頻繁に発生する遺伝子を有する新しいバイオマーカシグネチャを生成するバイオマーカコンセンサスエンジン128を含む。バイオマーカコンソリデータ104は、複数の候補バイオマーカにわたって全体的なエラー率を決定するためのエラー計算エンジン130を含む。バイオマーカジェネレータ102と同様に、バイオマーカコンソリデータ104もまた、自動的に制御またはユーザ操作され得るCCU 101によって制御され得る。CCU 101は、最小バイオマーカサイズについての好適な閾値を受信および/または決定し得、バイオマーカジェネレータ102およびバイオマーカコンソリデータ104の両方を動作させる反復の数を決定

10

20

30

40

50

するように、この情報を使用し得る。1つの実施形態において、各反復中に、CCU 101は、バイオマーカサイズを1つ減少させ、閾値が達せられるまでバイオマーカジェネレータ102およびバイオマーカコンソリデータ104の両方を反復する。そのような実施形態において、バイオマーカコンセンサスエンジン128は、各反復について、新しいバイオマーカシグネチャおよび新しい全体的なエラー率を出力する。したがって、バイオマーカコンセンサスエンジン128は、閾値から最大バイオマーカサイズまで様々である異なるサイズを各々が有する一組の新しいバイオマーカシグネチャ(複数)を出力する。バイオマーカコンソリデータ104はさらに、これらの新しいバイオマーカシグネチャの各々の性能尺度またはエラー率を検討して出力のために最適なバイオマーカを選択するバイオマーカ選択エンジン126を含む。バイオマーカコンソリデータ104およびそれぞれのエンジンの動作は、図2~4への参照によってさらに詳細に説明される。

10

【0021】

図3は、データセットを分類するための例示的なプロセスの流れ図である。ステップ302で、分類エンジン114は、トレーニングデータおよびテストデータを受信する。下記で説明されるように、分類エンジン114は、1つ以上の分類器を開発するためにトレーニングデータを使用し、次いで、1つ以上の分類器をテストデータに適用する。図3で図示されるように、トレーニングデータは、トレーニングデータセットT0.train304と、トレーニングクラスセットcl.train306とを含む。トレーニングデータセットT0.train304中の各要素は、データサンプル(例えば、特定の患者からの発現データのベクトル)を表し、トレーニングクラスセットcl.train306中の既知のクラス識別子に対応する。例えば、3クラスシナリオにおいて、トレーニングデータセットT0.train304中の第1の要素は、特定の疾患を有する患者についての遺伝子発現データを表し得、トレーニングクラスセットcl.train306中の第1の要素「疾患陽性」に対応し得、トレーニングデータセットT0.train304中の第2の要素は、特定の疾患に耐性または免疫がある患者についての遺伝子発現データを表し得、トレーニングクラスセットcl.train306中の第2の要素「疾患免疫性」に対応し得、トレーニングデータセットT0.train304中の第3の要素は、特定の疾患がない患者についての遺伝子発現データを表し得、トレーニングクラスセットcl.train306中の第3の要素「疾患なし」に対応し得る。ステップ302で受信されるテストデータは、テストデータセットT0.test308を含み、そのテストデータセットT0.testは、トレーニングデータセットT0.train304中のデータサンプルと同一の基礎的な種類のデータを表すが、例えば、異なる患者または異なる実験から採取されたサンプルを表し得る。任意選択で、分類エンジン114はまた、分類器がテストデータセットT0.test308に適用される場合に分類エンジン114によって生成される分類器の性能を評価するために使用され得る、テストデータセット中のデータサンプルについての既知のクラス識別子を含むテストクラスセットcl.test310を受信する。いくつかの実装において、テストデータセットT0.test308中のデータサンプルについてのいかなる既知のクラスも利用可能ではなく、したがって、テストクラスセットcl.test310は、分類エンジン114に提供されない。

20

30

40

【0022】

概して、ステップ302で受信されるデータは、サンプル中の複数の異なる遺伝子の発現値等の、分類が引き出され得る任意の実験データまたは別様に得られたデータ、および/または、任意の生物学的に意味のある被分析物のレベル等の種々の表現型の特性を表し得る。特定の実施形態において、データセットは、疾患状態についてのおよび対照状態についての発現レベルデータを含み得る。本明細書で使用される場合、「遺伝子発現レベル」という用語は、遺伝子によってコード化される分子(例えば、RNAまたはポリペプチド)の量を指し得る。mRNA分子の発現レベルは、mRNAの量(mRNAをコード化する遺伝子の転写活性によって決定される)、および、mRNAの安定性(mRNAの半減期によって決定される)を含み得る。遺伝子発現レベルはまた、遺伝子によってコード

50

化される所与のアミノ酸配列に対応するポリペプチドの量を含み得る。したがって、遺伝子の発現レベルは、遺伝子から転写される mRNA の量、遺伝子によってコード化されるポリペプチドの量、または、それら両方に対応することができる。遺伝子の発現レベルはさらに、遺伝子産物の異なる形態の発現レベルによってカテゴライズされ得る。例えば、遺伝子によってコード化される RNA 分子は、差次的に発現させられたスプライスバリエーション (differentially expressed splice variant)、異なる開始または終結部位を有する転写産物、および/または、他の特異的に処理された形態を含み得る。遺伝子によってコード化されるポリペプチドは、ポリペプチドの開裂および/または修飾形態を含み得る。ポリペプチドは、リン酸化、脂質化、プレニル化、硫酸化、水酸化、アセチル化、リボシル化、ファルネシル化、炭水化物の追加、および、同等物によって修飾されることができ、さらに、所与の種類の修飾を有するポリペプチドの複数の形態が、存在し得る。例えば、ポリペプチドは、複数の部位においてリン酸化され、異なるレベルの特異的にリン酸化されたタンパク質を発現し得る。

10

【0023】

特定の実施形態において、細胞または組織における遺伝子発現レベルは、遺伝子発現プロファイルによって表され得る。遺伝子発現プロファイルは、細胞または組織等の検体における遺伝子の発現レベルの特徴的な表現を指し得る。個体からの検体における遺伝子発現プロファイルの決定は、個体の遺伝子発現状態を表す。遺伝子発現プロファイルは、メッセンジャー RNA またはポリペプチドの発現、あるいは、細胞中または組織中の 1 つ以上の遺伝子によってコード化されるそれらの形態を反映する。発現プロファイルは、概して、異なる細胞または組織の間で異なる発現パターンを示す生体分子 (核酸、タンパク質、炭水化物) のプロファイルを指し得る。遺伝子発現プロファイルを表すデータサンプルは、発現レベルのベクトルとして記憶され得、ベクトルにおける各入力は、特定の生体分子または他の生物学的実体に対応する。

20

【0024】

特定の実施形態において、データセットは、サンプル中の複数の異なる遺伝子の遺伝子発現値を表す要素を含み得る。他の実施形態において、データセットは、質量分析によって検出されるピークを表す要素を含み得る。概して、各データセットは、複数の生物学的状態クラスのうちの一つに各々が対応するデータサンプル (複数) を含み得る。例えば、生物学的状態クラスは、サンプルのソース (すなわち、サンプルが取得される患者) における疾患の有無、病期、疾患のリスク、疾患の再発の可能性、1 つ以上の遺伝子座における共有遺伝子型 (例えば、共通 HLA ハプロタイプ、遺伝子における突然変異、メチル化等の遺伝子の修飾等)、作用物質 (例えば、毒性物質または潜在的に毒性の物質、環境汚染物質、候補薬剤等) または条件 (温度、pH 等) への曝露、人口学的特性 (年齢、性別、体重、家族歴、既往歴等)、作用物質への耐性、作用物質への感受性 (例えば、薬剤への反応性)、および、同等物を含むことができるが、それらに限定されない。

30

【0025】

データセットは、最終的な分類器選択における収集バイアスを低減するように、互いから独立し得る。例えば、それらは、複数のソースから収集されることができ、異なる除外または包含の基準を使用して異なる時間に異なる場所から収集され得、すなわち、データセットは、生物学的状態クラスを定義する特性外の特性を考慮する場合に、比較的ヘテロジニアスであり得る。ヘテロジェニシティ (heterogeneity) に寄与する要因は、性別、年齢、民族性による生物学的変動、摂食、運動、睡眠の挙動による個体的変動、および、血液処理のための臨床プロトコルによるサンプル取り扱い変動を含むが、それらに限定されない。しかしながら、生物学的状態クラスは、1 つ以上の共通特性を備え得る (例えば、サンプルソースは、疾患および同一の性別、または、1 つ以上の他の共通の人口学的特性を有する個体を表し得る)。特定の実施形態において、複数のソースからのデータセットは、異なる時間および/または異なる条件下における患者の同一の集団からのサンプルの収集によって生成される。

40

【0026】

50

特定の実施形態において、複数のデータセットは、複数の異なる臨床試験場から取得され、各データセットは、各個別試験場で取得される複数の患者サンプルを備える。サンプル種類は、血液、血清、血漿、乳頭吸引物、尿、涙、唾液、髄液、リンパ液、細胞および/または組織溶解物、レーザ顕微解剖組織または細胞サンプル、(例えば、パラフィンブロック中の、または、凍結された)埋め込み細胞または組織、(例えば、剖検からの)新鮮なまたは保存用のサンプルを含むが、それらに限定されない。サンプルは、例えば、インビトロで細胞または組織培養から得ることができる。代替として、サンプルは、生体から、または、単細胞生物等の生物の集団から得ることができる。1つの例において、特定の癌についてのバイオマーカーを識別する場合、2つのテスト場で独立したグループによって選択される対照から、血液サンプルが収集され、それによって、独立した独立したデータセットが開発されるサンプルを提供し得る。

10

【0027】

いくつかの実装において、トレーニングセットおよびテストセットは、バルクデータを受信してそのバルクデータをトレーニングデータセットとテストデータセットとに分割するデータ前処理エンジン110(図1)によって生成される。特定の実施形態において、データ前処理エンジン110は、データをこれら2つのグループにランダムに分割する。データをランダムに分割することが、クラスを予測してロバストな遺伝子シグネチャを生成するために望ましくあり得る。他の実施形態において、データ前処理エンジン110は、データの種類または標識に基づいて、データを2つ以上のグループに分割する。概して、データは、本開示の範囲から逸脱することなく、所望に応じた任意の好適な方法で、トレーニングデータセットおよびテストデータセットに分割されることができる。トレーニングデータセットおよびテストデータセットは、任意の好適なサイズを有し得、同一のまたは異なるサイズであり得る。特定の実施形態において、データ前処理エンジン110は、データをトレーニングデータセットとテストデータセットとに分割することの前に、1つ以上のデータを破棄し得る。特定の実施形態において、データ前処理エンジン110は、任意のさらなる処理の前に、トレーニングデータセットおよび/またはテストデータセットから1つ以上のデータを破棄し得る。

20

【0028】

ステップ311において、分類エンジン114は、カウンタ変数*i*を1に等しく設定する。ステップ312において、分類エンジン114は、トレーニングデータセットT0.train_304およびトレーニングクラスセットcl.train_306に基づいて、第1の分類器rf_314を生成する。図2は、データセット中の要素の分類を図示する。分類エンジン114は、サポートベクトルマシン技法、線形判別分析技法、ランダムフォレスト技法、k最近傍技法、部分最小二乗技法(部分最小二乗および線形判別分析特徴を組み合わせる技法を含む)、ロジスティック回帰技法、ニューラルネットワークベースの技法、決定木ベースの技法、および、(例えば、「Diagnosis of multiple cancer types by shrunken centroids of gene expression」PNAS, v. 99, n. 10, 2002で、Tibshirani、Hastle、Narasimhan、および、Chuによって説明されるような)収縮重心技法(shrunken centroid technique)を含むが、それらに限定されないいずれか1つ以上の既知の機械学習アルゴリズムをステップ312で使用し得る。いくつかのそのような技法は、線形判別分析、サポートベクトルマシン、ランダムフォレスト(Breiman, Machine Learning, 45(1):5-32(2001))、k最近傍(Bishop, Neural Networks for Pattern Recognition, ed. O.U. Press, 1995)、部分最小二乗判別分析、および、PAMR(Tibshirani et al., Proc Natl Acad Sci USA, 99(10):6567-6572(2002))に対応する、lda、svm、randomForest、knn、pls_lda、および、pamrを含むRプログラミング言語用パッケージとして利用可能である。分類エンジ

30

40

50

ン 1 1 4 は、ステップ 3 1 2 で、第 1 の分類器 $r f \quad 3 1 4$ をメモリに記憶し得る。

【 0 0 2 9 】

ステップ 3 1 6 において、分類エンジン 1 1 4 は、第 1 の分類器 $r f \quad 3 1 4$ (ステップ 3 1 2 で生成される) をテストデータセット $T 0 . t e s t \quad 3 0 8$ に適用することによって、一組の予測されたテスト分類 $p r e d c l . t e s t \quad 3 1 8$ を生成する。分類エンジン 1 1 4 は、ステップ 3 1 6 で、予測された分類 $p r e d c l . t e s t \quad 3 1 8$ をメモリに記憶し得る。

【 0 0 3 0 】

ステップ 3 2 0 において、分類エンジン 1 1 4 は、トレーニングデータセット $T 0 . t r a i n \quad 3 0 4$ を変換する。この変換は、トレーニングクラスセット $c l . t r a i n \quad 3 0 6$ に基づいてトレーニングデータセット $T 0 . t r a i n \quad 3 0 4$ を変換する変換関数 $c o r r e c t e d D a t a$ に従って進む。ステップ 3 1 0 の変換の結果は、分類エンジン 1 1 4 がメモリに記憶し得る変換されたトレーニングデータセット $T 0 . t r a i n . 2 \quad 3 2 2$ である。いくつかの実装において、ステップ 3 2 0 で分類エンジン 1 1 4 によって行われる変換は、バイアス補正技法を含む。例えば、変換は、全体として採取されるデータセットの重心、または、データセットにおいて表される各クラスの重心に関して、トレーニングデータセット $T 0 . t r a i n \quad 3 0 4$ の要素を調整することによって、トレーニングデータセット $T 0 . t r a i n \quad 3 0 4$ を「再び中心に置いて」もよい。

【 0 0 3 1 】

1 つの特定の再中心化技法は、異なるグループの重心の中心に基づいて、トレーニングデータセット $T 0 . t r a i n \quad 3 0 4$ の要素を中心に置くことを伴う。トレーニングデータセット $T 0 . t r a i n \quad 3 0 4$ 中に n 個のデータサンプルがあり、かつ、各データサンプルが p 個の入力 (例えば、 p 個の異なる遺伝子についての発現レベルを表す) を有するベクトルである場合、 x_{ij} にデータサンプル j の i 番目の入力を表させる。トレーニングクラスセット $c l . t r a i n \quad 3 0 8$ が K 個の異なるクラスを表す場合、クラス k における n_k 個のサンプルの指数を C_k に表させる。分類エンジン 1 1 4 は、クラス k の重心の i 番目の成分を下記のように計算し得、

【数 1】

$$\bar{x}_{ik} = \sum_{j \in C_k} \frac{x_{ij}}{n_k} \quad (1)$$

かつ、クラス重心の中心の i 番目の成分を下記のように計算し得る。

【数 2】

$$\bar{x}_i^c = \sum_{k=1}^K \frac{\bar{x}_{ik}}{K} \quad (2)$$

【 0 0 3 2 】

分類エンジン 1 1 4 はまた、全体的な重心の i 番目の成分を下記のように計算し得る。

【数 3】

$$\bar{x}_i = \sum_{j=1}^n \frac{x_{ij}}{n} \quad (3)$$

【 0 0 3 3 】

次いで、分類エンジン 1 1 4 は、下記によって求められる差を加えることによって、ト

10

20

30

40

50

レーニングデータセット $T0.train$ 304 の各要素の中の i 番目の入力を調整することを含む変換を行ってもよい。

【数 4】

$$\Delta = -\bar{x}_i^c \quad (4)$$

【0034】

いくつかの実装において、ステップ 320 で行われる変換は、方程式 1 ~ 4 への参照によって上記で説明されるもの以外の偏移 (shift)、回転、シアー、これらの変換の組み合わせ、または、任意の他の線形あるいは非線形の変換を含む。

10

【0035】

ステップ 324 において、分類エンジン 114 は、テストデータセット $T0.test$ 308 を変換する。テストデータセット $T0.test$ 308 に適用される変換、correctedData は、ステップ 320 でトレーニングデータセット $T0.train$ 304 に適用される同一の種類の変換であるが、 $T0.train$ 304 および $predcl.train$ 314 の代わりに、引数 $T0.test$ 308 および $predcl.test$ 318 に関して適用される。例えば、トレーニングデータセット $T0.train$ 304 の要素が、トレーニングデータセット $T0.train$ 304 のクラスの重心に関して計算されるような方程式 4 によって求められる の値によって、ステップ 320 で調整される場合には、テストデータセット $T0.test$ 308 の要素は、テストデータセット $T0.test$ 308 のクラスの重心に関して計算されるような方程式 4 によって求められる の値によって、ステップ 324 で調整される。ステップ 324 の変換の結果は、分類エンジン 114 がメモリに記憶し得る変換されたテストデータセット $T0.test.2$ 326 である。

20

【0036】

ステップ 327 において、分類エンジン 114 は、反復カウンタ i の値が 1 に等しいかどうかを決定する。そうである場合、分類エンジン 114 は、分類エンジン 114 が、第 2 の分類器 $rf2$ 329 を生成するために、変換されたトレーニングデータセット $T0.train.2$ 322 およびトレーニングクラスセット $cl.train$ 306 を使用するステップ 328 を続けて実行する。ステップ 332 およびステップ 336 への参照によって上記で説明されるように、任意の機械学習技法が、ステップ 328 で分類器を生成するために適用され得る。第 2 の分類器 $rf2$ 329 は、第 1 の分類器 rf 314 (例えば、両方の SVM 分類器) と同一の種類であり得るか、または、異なる種類であり得る。

30

【0037】

ステップ 331 において、分類エンジン 114 は、反復カウンタ i をインクリメントし、次いで、分類エンジン 114 が第 2 の分類器 $rf2$ 329 を (ステップ 324 で分類エンジン 114 によって生成されるような) 変換されたテストデータセット $T0.test.2$ 326 に適用するステップ 333 を続けて実行する。ステップ 333 の出力は、変換されたデータセット $T0.test.2$ 326 のための一組の予測された分類 $predcl.test.2$ 330 である。分類エンジン 114 は、表示デバイス、印刷デバイス、記憶デバイス、ネットワークにわたって分類エンジン 114 と通信している別のデバイス、または、システム 100 の内部あるいは外部の任意の他のデバイスに、予測された分類を出力し得る。

40

【0038】

ステップ 332 において、分類エンジン 114 は、(ステップ 316 で生成されるような) 予測された分類セット $predcl.test$ 318 の分類と (ステップ 328 で生成されるような) 予測された分類セット $predcl.test.2$ 330 の分類との間に何らかの差異があるかどうかを決定する。予測された分類のセットが一致する (すなわち、テストデータセット $T0.test$ 308 中の各データサンプルについて、そ

50

のデータサンプルについての予測されたクラスが、2つの予測された分類セットの間で同一である)場合には、分類エンジン114は、ステップ338へ進み、予測された分類セットpredcl.test.2 330(同等に、予測された分類セットpredcl.test 318)をテストデータセットT0.test 308の最終的な分類として出力する。

【0039】

分類エンジン114が分類データセットpredcl.test 318と分類データセットpredcl.test.2 330との間の差異を識別する場合、分類エンジン114は、ステップ334へ進み、テストデータセットT0.test 308の以前に記憶された値を、(ステップ324の変換によって生成されるような)変換されたテストデータセットT0.test.2 326の値と置換する。結果として、テストデータセットT0.test 308は、変換されたテストデータセットT0.test.2 326の値を有する。分類エンジン114は、ステップ336へ進み、(ステップ316で生成されるような)予測された分類セットpredcl.test 318の以前に記憶された値を、(ステップ328で生成されるような)予測された分類セットpredcl.test.2 330の値と置換する。結果として、予測された分類セットpredcl.test 318は、予測された分類セットpredcl.test.2 330の値を有する。

10

【0040】

テストデータセットT0.test 308の値が変換されたテストデータセットT0.test.2 326の値で更新され、かつ、予測された分類セットpredcl.test 318が予測された分類セットpredcl.test.2 330の値で更新されると、分類エンジン114は、ステップ324に戻って新しい変換を行い、分類エンジン114が(ステップ332で)予測された分類の間に差異がないことを決定するまで、このプロセスを反復する。

20

【0041】

分類器性能監視エンジン116は、好適な性能測定基準を使用して、図3のプロセスの終わりに、分類エンジン114によって生成される最終的な分類の性能を分析し得る。特定の実施形態において、性能測定基準は、エラー率を含み得る。性能測定基準はまた、試行された予測の総数によって除算された正しい予測の数を含み得る。性能測定基準は、本開示の範囲から逸脱することなく、任意の好適な尺度であり得る。

30

【0042】

本主題の実装は、本明細書で説明されるような1つ以上の特徴と、1つ以上の機械(例えば、コンピュータ、ロボット)に本明細書で説明される動作を実現させるように動作可能な機械可読媒体を備える物品とを備えるシステム、方法、および、コンピュータプログラム製品を含むことができるが、それらに限定されない。本明細書で説明される方法は、単一のコンピューティングシステムまたは複数のコンピューティングシステムに存在する1つ以上のプロセッサまたはエンジンによって実装されることができる。そのような複数のコンピューティングシステムは、接続されることができ、複数のコンピューティングシステムのうちの1つ以上の間の直接接続を介したネットワーク(例えば、インターネット、無線広域ネットワーク、ローカルエリアネットワーク、広域ネットワーク、有線ネットワーク、または、同等物)を経由した接続を含むが、それに限定されない1つ以上の接続を介して、データおよび/またはコマンド、あるいは、他の命令または同等物を交換することができる。

40

【0043】

図4は、図1~3への参照によって説明されるプロセスを行うための回路を含む図1のシステム100の構成要素のうちのいずれか等の、コンピューティングデバイスのブロック図である。システム100の構成要素の各々は、1つ以上のコンピューティングデバイス400上に実装され得る。特定の局面において、複数の上記の構成要素およびデータベースは、1つのコンピューティングデバイス400内に含まれ得る。特定の实装において

50

、構成要素およびデータベースは、いくつかのコンピューティングデバイス400にわたって実装され得る。

【0044】

コンピューティングデバイス400は、少なくとも1つの通信インターフェースユニットと、入力/出力コントローラ410と、システムメモリと、1つ以上のデータ記憶デバイスとを備える。システムメモリは、少なくとも1つのランダムアクセスメモリ(RAM 402)と、少なくとも1つの読み取り専用メモリ(ROM 404)とを含む。これらの要素は全て、中央処理ユニット(CPU 406)と通信し、コンピューティングデバイス400の動作を促進する。コンピューティングデバイス400は、多くの異なる方法で構成され得る。例えば、コンピューティングデバイス400は、従来のスタンドアロンコンピュータであり得るか、または、代替として、コンピューティングデバイス400の機能は、複数のコンピュータシステムおよびアーキテクチャにわたって分散され得る。コンピューティングデバイス400は、データ分割、区別、分類、スコア化、ランク付け、および、記憶の動作のうちのいくつかまたは全てを行うように構成され得る。図4において、コンピューティングデバイス400は、ネットワークまたはローカルネットワークを介して、他のサーバまたはシステムにリンクされる。

10

【0045】

コンピューティングデバイス400は、分散されたアーキテクチャにおいて構成され得、データベースおよびプロセッサは、別個のユニットまたは場所において格納される。いくつかのそのようなユニットは、一次処理機能を行い、最低限でも、一般コントローラまたはプロセッサおよびシステムメモリを含む。そのような局面において、これらのユニットの各々は、通信インターフェースユニット408を介して、他のサーバ、クライアント、または、ユーザコンピュータ、および、他の関連デバイスとの一次通信リンクとしての役割を果たす通信ハブまたはポート(図示せず)に取り付けられる。通信ハブまたはポートは、それ自体が最小処理能力を有し、主に、通信ルータとしての役割を果たし得る。種々の通信プロトコルは、限定されないが、Ethernet(登録商標)、SAP、SAS(登録商標)、ATP、Bluetooth(登録商標)、GSM(登録商標)、および、TCP/IPを含むシステムの一部であり得る。

20

【0046】

CPU 406は、1つ以上の従来のマイクロプロセッサ等のプロセッサ、および、CPU 406から作業負荷をオフロードするための数値演算コプロセッサ等の1つ以上の補助コプロセッサを備える。CPU 406は、通信インターフェースユニット408および入力/出力コントローラ410と通信し、それらを通して、CPU 406は、他のサーバ、ユーザ端末、または、デバイス等の他のデバイスと通信する。通信インターフェースユニット408および入力/出力コントローラ410は、例えば、他のプロセッサ、サーバ、または、クライアント端末と同時に通信するための複数の通信チャネルを含み得る。相互に通信しているデバイスは、継続的に相互に伝送している必要はない。反対に、そのようなデバイスは、必要に応じて相互に伝送する必要が少なく、実際には、ほとんどの時間、データを交換することを控え得、いくつかのステップが行われることを要求することにより、デバイス間の通信リンクを確立し得る。

30

40

【0047】

CPU 406はまた、データ記憶デバイスと通信する。データ記憶デバイスは、磁気、光学、または、半導体のメモリの適切な組み合わせを備え得、例えば、RAM 402、ROM 404、フラッシュドライブ、コンパクトディスクまたはハードディスクあるいはドライブ等の光学ディスクを含み得る。CPU 406およびデータ記憶デバイスは、各々、例えば、単一のコンピュータまたは他のコンピューティングデバイス内に全体的に位置し得るか、または、USBポート、シリアルポートケーブル、同軸ケーブル、Ethernet(登録商標)型ケーブル、電話回線、無線周波数送受信機、または、他の類似の無線もしくは有線の媒体、あるいは、前述のものの組み合わせ等の通信媒体によって、相互に接続され得る。例えば、CPU 406は、通信インターフェースユニット40

50

8を介して、データ記憶デバイスに接続され得る。CPU 406は、1つ以上の特定の処理機能を行なうように構成され得る。

【0048】

データ記憶デバイスは、例えば、(i)コンピューティングデバイス400のためのオペレーティングシステム412、(ii)本明細書で説明されるシステムおよび方法に従って、特に、CPU 406に関して詳細に説明されるプロセスに従って、CPU 406に命令するように適合させられた1つ以上のアプリケーション414(例えば、コンピュータプログラムコードまたはコンピュータプログラム製品)、または、(iii)プログラムによって要求される情報を記憶するために利用され得る情報を記憶するように適合させられたデータベース(単数または複数)416を記憶し得る。いくつかの局面において、データベースは、実験データ、および、既刊文献モデルを記憶するデータベース(単数または複数)を含む。

10

【0049】

オペレーティングシステム412およびアプリケーション414は、例えば、圧縮、アンコンパイル、および、暗号化されたフォーマットにおいて記憶され得、コンピュータプログラムコードを含み得る。プログラムの命令は、ROM 404またはRAM 402から等、データ記憶デバイス以外のコンピュータ可読媒体から、プロセッサのメインメモリに読み込まれ得る。プログラムにおける命令のシーケンスの実行は、CPU 406に、本明細書で説明されるプロセスステップを行なわせるが、有線回路が、本発明のプロセスの実装のためのソフトウェア命令の代わりに、または、それと組み合わせて使用され得る。したがって、説明されるシステムおよび方法は、ハードウェアおよびソフトウェアの任意の特定の組み合わせに限定されない。

20

【0050】

好適なコンピュータプログラムコードは、本明細書で説明されるようなモデル化、スコア化、および、集約に関連する1つ以上の機能を果たすために提供され得る。プログラムはまた、オペレーティングシステム412、データベース管理システム、および、プロセッサが入力/出力コントローラ410を介してコンピュータ周辺デバイス(例えば、ビデオディスプレイ、キーボード、コンピュータマウス等)と連動することを可能にする「デバイスドライバ」等のプログラム要素を含み得る。

【0051】

コンピュータ可読命令を備えるコンピュータプログラム製品も、提供される。コンピュータ可読命令は、コンピュータシステム上にロードされて実行される場合、本方法または上記で説明される方法の1つ以上のステップに従って、コンピュータシステムを動作させる。本明細書で使用される場合、「コンピュータ可読媒体」という用語は、実行のために、コンピューティングデバイス400のプロセッサ(または、本明細書で説明されるデバイスの任意の他のプロセッサ)に命令を提供するかまたは提供に関与する任意の非一時的媒体を指す。そのような媒体は、不揮発性媒体および揮発性媒体を含むが、それらに限定されない多くの形態をとり得る。不揮発性媒体は、例えば、光学、磁気、または、光磁気のディスク、あるいは、フラッシュメモリ等の集積回路メモリを含む。揮発性媒体は、典型的にメインメモリを構成するダイナミックランダムアクセスメモリ(DRAM)を含む。コンピュータ可読媒体の共通の形態は、例えば、フロッピー(登録商標)ディスク、フレキシブルディスク、ハードディスク、磁気テープ、任意の他の磁気媒体、CD-ROM、DVD、任意の他の光学媒体、パンチカード、ペーパーテープ、孔のパターンを有する任意の他の物理的媒体、RAM、PROM、EPROM、または、EEPROM(電氣的に消去可能なプログラマブル読み取り専用メモリ)、FLASH-EEPROM、任意の他のメモリチップまたはカートリッジ、あるいは、コンピュータが読み取ることができる任意の他の非一時的媒体を含む。

30

40

【0052】

コンピュータ可読媒体の種々の形態は、実行のために、1つ以上の命令の1つ以上のシーケンスをCPU 406(または本明細書で説明されるデバイスの任意の他のプロセッ

50

サ)に搬送することに関与し得る。例えば、命令は、最初に、遠隔コンピュータ(図示せず)の磁気ディスク上にあり得る。遠隔コンピュータは、命令をその動的メモリ内にロードし、Ethernet(登録商標)接続、ケーブルライン、または、モデムを使用する電話回線をも経由して、命令を送信することができる。コンピューティングデバイス400(例えば、サーバ)にローカルの通信デバイスは、それぞれの通信ライン上でデータを受信し、プロセッサのためのシステムバス上にデータを置くことができる。システムバスは、データをメインメモリに搬送し、そこから、プロセッサは、命令を読み出して実行する。メインメモリによって受信される命令は、任意選択で、プロセッサによる実行の前または後のいずれかにおいて、メモリに記憶され得る。加えて、命令は、通信ポートを介して、種々のタイプの情報を搬送する無線通信またはデータストリームの例示的形態である電氣的、電磁的、または、光学的な信号として受信され得る。

10

【実施例】

【0053】

下記の公開データセットを、Gene Expression Omnibus(GEO)(<http://www.ncbi.nlm.nih.gov/geo/>)リポジトリからダウンロードする。

【表1】

- a. GSE10106 (www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE10106)
- b. GSE10135 (www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE10135)
- c. GSE11906 (www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE11906)
- d. GSE11952 (www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE11952)
- e. GSE13933 (www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE13933)
- f. GSE19407 (www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE19407)
- g. GSE19667 (www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE19667)
- h. GSE20257 (www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE20257)
- i. GSE5058 (www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE5058)
- j. GSE7832 (www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE7832)
- k. GSE8545 (www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE8545).

20

30

【0054】

トレーニングデータセットは、Affymetrixプラットフォーム(HGU-133+2)上にある。未加工データファイルを、R(R Development Core Team, 2007)の中のBioconductor(Gentleman, 2004)に属するaffyパッケージ(Gautier, 2004)のReadAffy機能によって読み取り、品質を、RNA分解プロット(affyパッケージのAffyRNAdeg機能を伴う)、NUSE、および、RLEプロット(機能affyPLM(Brettschneider, 2008)を伴う)を生成し、MA(RLE)値を計算し、品質管理チェック上の一組の閾値を下回るか、または、上記のデータセットの中で複製されるトレーニングデータセットからアレイを除外し、gcrmaアルゴリズム(Wu, 2004)を使用して品質管理チェックに合格するアレイを正規化することによって、管理する。トレーニングセットサンプル分類を、各データセットについてのGEOデータベースのシリーズマトリクスファイルから取得する。出力は、233個のサンプル(28個のCOPDサンプルおよび205個の対照サンプル)についての54675個のプロブセットを伴う遺伝子発現マトリクスから成る。均衡の取れたデータセットを製作するために、COPDサンプルは、同時係属中の米国仮特許出願第61/662812号で説明されるようなDual Ensemble方法を適用する前に、224個のCOPDサンプルを取得するための多重時間(multiple time)であった。205

40

50

人の対照および224人のCOPD患者を含む複合データセットを用いて、409個の遺伝子を有する遺伝子シグネチャを構築した。850個の二進値を、ランダムベクトルにおいて使用した。本方法で使用される分類方法は、下記のRパッケージ、すなわち、lda、svm、randomForest、knn、pls.lda、および、pamrを含んでいた。最大反復を、5000であるように設定した。マシューズ相関係数(MCC)、トレーニングデータセットにおける相互検証プロセスの精度は、それぞれ、0.743、0.87である。トレーニングデータセット中の遺伝子シグネチャのヒートマップを、図5に示す。図5のヒートマップにおいて、遺伝子発現値を、行ごとに中心に置いた。ヒートマップの色は、グレースケールでは明確に示されない場合もあるが、図5のデータは、対照データが左に示され、COPDデータが右側に示されていることを示す。テストデータセットは、16個の対照サンプルおよび24個のCOPDサンプルを含む民間供給業者(GeneLogic)から入手した未公開データセットである。本発明の変換不変方法を適用することなく、Dual Ensembleによって生成される遺伝子シグネチャは、合計40個のサンプルうちの29個のサンプルを正しく予測した。精度は0.725であり、MCCは0.527である。16個の対照サンプルにおいて、遺伝子シグネチャは、15個を対照として正しく予測したが、1個をCOPDとして誤って予測した。24個のCOPDサンプルの間で、遺伝子シグネチャは、14個をCOPDサンプルとして正しく予測したが、10個を対照として誤って予測した。

10

【0055】

しかしながら、変換不変方法が適用された場合、2つまたは複数のクラスを中心、および、100に設定された最大反復に従って偏移(shift)を伴った。同一の遺伝子シグネチャは、合計40個のサンプルのうちの30個のサンプルを正しく予測した。精度は0.75であり、MCCは0.533である。16個の対照サンプルにおいて、遺伝子シグネチャは、14個を対照として正しく予測したが、2個をCOPDとして誤って予測した。24個のCOPDサンプルの間で、遺伝子シグネチャは、16個をCOPDサンプルとして正しく予測したが、8個を対照として誤って予測した。

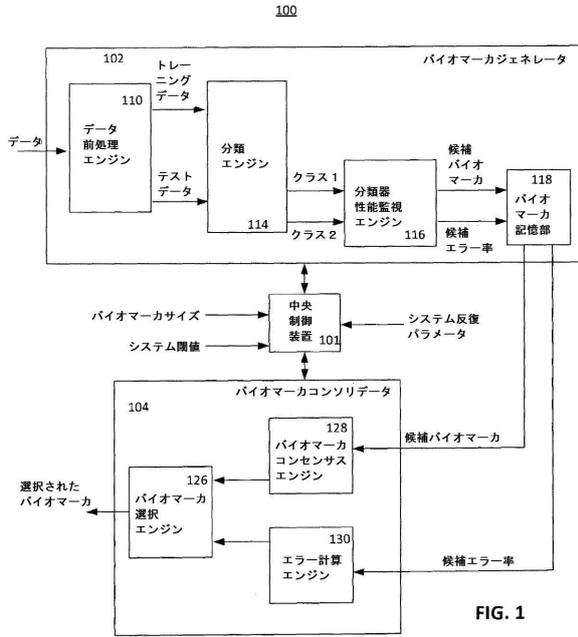
20

【0056】

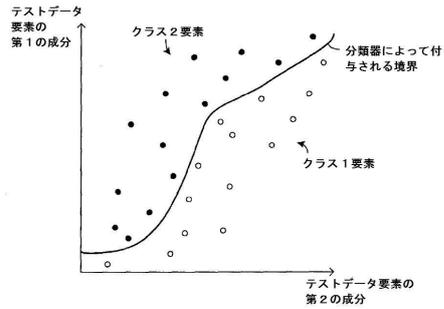
本発明の実装は、特定の例を参照して特定して示され、説明されているが、本開示の精神および範囲から逸脱することなく、形態および詳細における種々の変更がそれに行われ得ることが、当業者によって理解されるべきである。

30

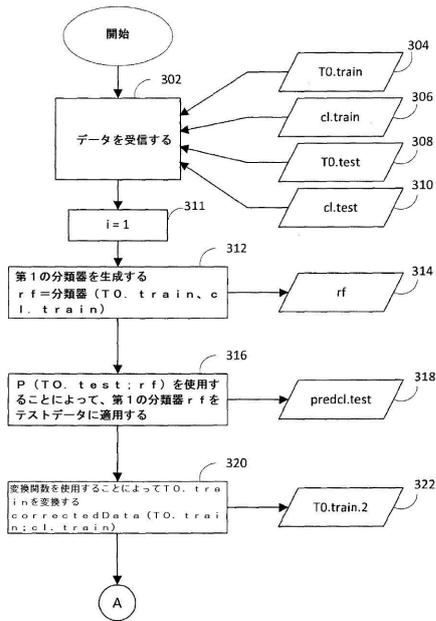
【図1】



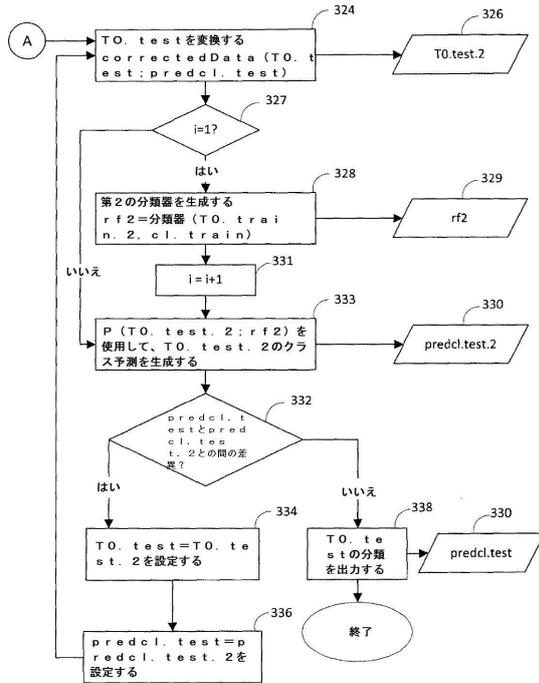
【図2】



【図3-1】



【図3-2】



【図4】

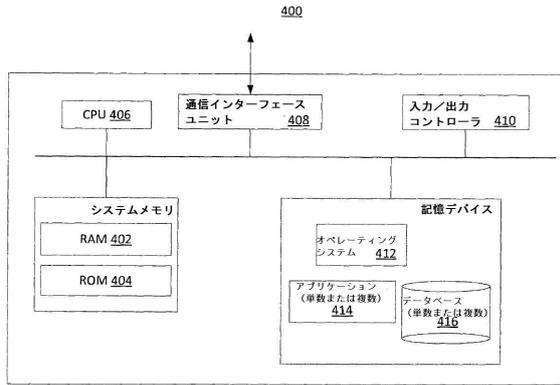


FIG. 4

【図5】

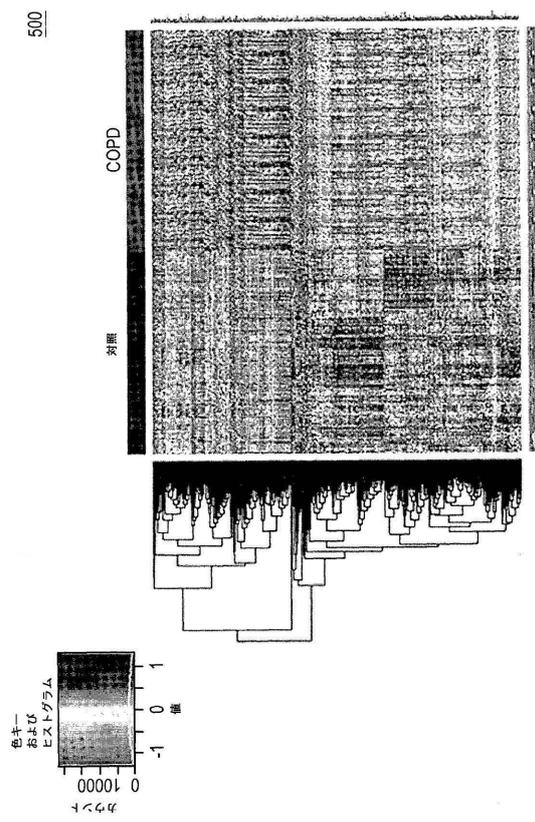


FIG. 5

フロントページの続き

(72)発明者 マルティン, フロリアン

スイス国 ツェーハー - 2034 プスー, シュマン ドゥ ロレー 1

(72)発明者 シアン, ヤン

スイス国 ツェーハー - 2000 ヌーシャテル, リュ ドゥ ロシェ 24

審査官 笠田 和宏

(56)参考文献 特表2012-501183(JP,A)

特開2009-075737(JP,A)

米国特許出願公開第2006/0074826(US,A1)

(58)調査した分野(Int.Cl., DB名)

IPC G06F 15/18

17/30

19/00 - 19/28

G06N 3/00 - 3/12

7/08 - 99/00

G06Q 10/00 - 99/00