

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第6430998号
(P6430998)

(45) 発行日 平成30年11月28日 (2018.11.28)

(24) 登録日 平成30年11月9日 (2018.11.9)

(51) Int.Cl.		F I			
C 1 2 N	15/11	(2006.01)	C 1 2 N	15/11	Z
C 1 2 Q	1/6837	(2018.01)	C 1 2 Q	1/6837	
C 1 2 Q	1/6874	(2018.01)	C 1 2 Q	1/6874	
G 0 6 F	19/22	(2011.01)	G 0 6 F	19/22	

請求項の数 14 外国語出願 (全 118 頁)

(21) 出願番号	特願2016-117074 (P2016-117074)	(73) 特許権者	513156537
(22) 出願日	平成28年6月13日 (2016.6.13)		ナテラ, インコーポレイテッド
(62) 分割の表示	特願2008-542450 (P2008-542450) の分割		アメリカ合衆国 カリフォルニア 940 70, サン カルロス, インダストリ アル ストリート 201, スイート 410
原出願日	平成18年11月22日 (2006.11.22)	(74) 代理人	100099759
(65) 公開番号	特開2016-184429 (P2016-184429A)		弁理士 青木 篤
(43) 公開日	平成28年10月20日 (2016.10.20)	(74) 代理人	100077517
審査請求日	平成28年6月13日 (2016.6.13)		弁理士 石田 敬
(31) 優先権主張番号	60/739,882	(74) 代理人	100087871
(32) 優先日	平成17年11月26日 (2005.11.26)		弁理士 福本 積
(33) 優先権主張国	米国 (US)	(74) 代理人	100087413
(31) 優先権主張番号	60/742,305		弁理士 古賀 哲次
(32) 優先日	平成17年12月6日 (2005.12.6)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 予測を行うための、遺伝子データを清浄化し、そして、そのデータを使用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】

標的個体の所与の染色体の所与のセグメント上の多数遺伝子座の測定値を用いて、標的個体のゲノム中の前記所与のセグメントの存在数を決定する方法であって、

(i) 標的個体のゲノムに存在する前記所与のセグメントの存在数に関する1以上の仮説のセットを創製し、

(ii) 前記所与のセグメント上の複数の遺伝子座における可能な対立遺伝子の一部又は全部についての遺伝子データの量を測定し、ここで、前記遺伝子データは、遺伝子型及び/又は表現型のデータであり、

(iii) 標的個体の遺伝子データの測定値、及び、関連個体の遺伝子データの測定値に基づいて、前記仮説の各々の相対的確率を決定し、ここで、前記仮説の各々の相対的確率の決定は、前記遺伝子座において予測される遺伝子データ及び観察される遺伝子データの統計学的な分布と、前記遺伝子座間の交差確率とに基づき、該仮説の各々の尤度を計算することにより行われ、ここで、前記関連個体は、前記標的個体の親から選択され、

(iv) 各仮説に関連する相対的確率を用いて、最も高い確率を有する仮説を、標的個体の現実の遺伝物質の最も可能性の高い前記所与のセグメントの存在数として決定することを含む方法。

【請求項2】

標的個体のゲノムに存在する染色体のセグメントの存在数の決定が、染色体異常のスクリーニングを通じて行われ、前記異常は、モノソミー、片親二染色体、トリソミー、異数

性、アンバランスなトランスロケーション、およびその組合せからなる群より選択される、請求項 1 記載の方法。

【請求項 3】

各仮説の相対的確率の決定が、マッチドフィルタリングの概念を用いて行われる、請求項 1 又は 2 に記載の方法。

【請求項 4】

各仮説の相対的確率の決定が、対立遺伝子要求を行わない定量的技術を用いて行われ、ここで、各遺伝子座の測定値の平均および標準偏差は、既知、未知、または均一のいずれかである、請求項 1 又は 2 に記載の方法。

【請求項 5】

各仮説の相対的確率の決定が、対立遺伝子要求を用いる定性的技術を用いて行われる、請求項 1 又は 2 に記載の方法。

【請求項 6】

各仮説の相対的確率の決定が、参照配列の公知の対立遺伝子と、対立遺伝子の定量的測定値とを用いて行われる、請求項 1 又は 2 に記載の方法。

【請求項 7】

標的個体が、成人ヒト、若年ヒト、ヒト胎児、ヒト胚、非ヒト成体、非ヒト若年体、非ヒト胎児、および非ヒト胚よりなる群から選択される、請求項 1 ~ 6 の何れか一項に記載の方法。

【請求項 8】

標的個体の遺伝子データが、ポリメラーゼ鎖反応 (PCR)、リガーゼ媒介 PCR、縮重オリゴヌクレオチドプライマー PCR、多重置換増幅、対立遺伝子 - 特異的増幅およびその組合せよりなる群から選択されるツールおよび / または技術を用いて増幅される、請求項 1 ~ 7 の何れか一項に記載の方法。

【請求項 9】

標的個体の遺伝子データが、分子逆転プローブ (MIP)、ゲノタイピングマイクロアレイ、Taqman SNPゲノタイピングアッセイ、Illuminaゲノタイピングシステム、他のゲノタイピングアッセイ、蛍光イン - サイチュハイブリダイゼーション (FISH)、およびその組合せよりなる群から選択されるツールおよび / または技術を用いて測定される、請求項 1 ~ 8 の何れか一項に記載の方法。

【請求項 10】

標的個体の遺伝子データが、標的個体のバルクジプロイド組織、標的個体から取られる 1 以上のジプロイド細胞、標的個体から取られる 1 以上の胚盤胞、標的個体上で見出された細胞外遺伝物質、母性血液で見出された標的個体からの細胞外遺伝物質、母性血液中の標的個体由来の細胞、標的個体に由来することが既知の遺伝物質、およびその組合せよりなる群から選択される物質を分析することにより測定される、請求項 1 ~ 8 の何れか一項に記載の方法。

【請求項 11】

標的個体の染色体または染色体セグメントの数の決定が、体外受精における胚選択の目的に用いられる、請求項 1 ~ 10 の何れか一項に記載の方法。

【請求項 12】

標的個体の染色体または染色体セグメントの数の決定が、出生前遺伝子診断の目的に用いられる、請求項 1 ~ 10 の何れか一項に記載の方法。

【請求項 13】

請求項 1 ~ 12 の何れか一項に記載の方法を実現するように構成された、コンピューターにより実施されるシステム。

【請求項 14】

標的個体の遺伝子データの不完全な知識と、標的個体に遺伝的に関連する 1 人以上の関連個体の遺伝子データの知識とに基づいて、標的個体の遺伝子データと、標的個体のゲノムに存在する染色体または染色体セグメントの存在数とを決定する方法であって、

10

20

30

40

50

ここで、前記遺伝子データは、遺伝子型及び/又は表現型のデータであり、ここで、前記関連個体は、前記標的個体の親から選択され、

(i) 関連個体のどの染色体のどのセグメントが、標的個体のゲノムに見出されるセグメントに対応するかについて、1以上の仮説のセットを創製し、

(ii) 標的個体のゲノムに存在する所与の染色体セグメントの数に関する1以上の仮説のセットを創製し、

(iii) 前記所与のセグメント上の複数の遺伝子座における可能な対立遺伝子の各々について遺伝子データの量を測定し、

(iv) 標的個体の遺伝子データの測定値の測定値と、関連個体の遺伝子データの測定値とに基づいて、前記仮説の各々の相対的確率を決定し、ここで、前記仮説の各々の相対的確率の決定は、前記複数の遺伝子座において予測される遺伝子データ及び観察される遺伝子データの統計学的な分布と、前記複数の遺伝子座間の交差確率とに基づき、該仮説の各々の尤度を計算することにより行われ、

(v) 各仮説に関連する相対的確率を用いて、最も高い確率を有する仮説を、標的個体の現実の遺伝物質の最も可能性が高い前記所与の染色体及びセグメントの存在数として決定する

ことを含む方法。

【発明の詳細な説明】

【技術分野】

【0001】

(関連出願への相互参照)

本願は、米国特許法の下で、以下の米国仮特許出願の利益を主張する：2005年11月26日出願の第60/739,882号；2005年12月6日出願の第60/742,305号；2005年12月9日出願の第60/754,396号；2006年2月21日出願の第60/774,976号；2006年4月4日出願の第60/789,506号；2006年6月30日出願の第60/817,741号；2006年7月31日出願の第11/496,982号；および2006年9月22日出願の第60/846,610号；これらの開示は、その全体が本明細書中に参考として援用される。

【0002】

(技術分野)

本発明は、一般には、医療的に予測される目的のための遺伝子データを獲得し、操作し、および用いる分野、具体的には、不完全に測定された遺伝子データを遺伝的に関連する個体の公知の遺伝子データを用いることによってより正確とし、それにより、種々の表現型結果をもたらす遺伝子不規則性のより効果的な同定を可能とするシステムに関する。また、本発明は、一般に、遺伝子、表現型および臨床的情報を分析し、管理し、それに作用させ、およびその情報を用いて、医療的決定の表現型結果を予測する分野に関する。さらに詳しくは、本発明は、対象の群からの一体化され、確認された遺伝子および表現型データを用いて、特定の対象に関して良好な決定を行う方法およびシステムに関する。

【背景技術】

【0003】

(関連技術の背景)

出生前および着床前遺伝子診断

出生前診断の現行の方法は医師および親に対して成長する胎児における異常を警告することができる。出生前診断がなければ、50人の赤ん坊の内1人は深刻な身体または精神的ハンディキャップを備えたまま誕生し、30人の内1人のように多くの者は先天的奇形のいくつかの形態を有するであろう。あいにくと、標準的な方法は侵襲性テストを必要とし、流産の大きか1%の危険性を有している。これらの方法は羊水穿刺、絨毛膜絨毛パイオプシーおよび胎児血液サンプリングを含む。これらの内、羊水穿刺は最も普通の手法であり；2003年において、それは全ての妊娠のほぼ3%で行われていたが、その使用頻度は過去15年にわたって減少してきた。出生前診断の主な欠点は、限定された活動のコ

10

20

30

40

50

ースを仮定すれば、一旦異常が検出されれば、それは非常に深刻な欠陥についてテストするには価値がありかつ倫理的であるに過ぎない。結果として、出生前診断は、典型的には、高い危険性の妊娠の場合に試みられるに過ぎず、そこでは、潜在的異常の深刻性と組合わされた上昇した欠陥の確率が危険性を凌ぐ。これらの危険性を緩和する出生前診断の方法に対する要望が存在する。

【 0 0 0 4 】

最近、無細胞胎児DNAおよび無傷胎児細胞が母体血液循環に入ることができるのが発見された。結果として、これらの細胞の分析は、早期の非侵襲性出生前遺伝子診断(NIPGD)を可能とすることができる。NIPGDを用いることにおける鍵となる挑戦は、母体血液から胎児の細胞または核酸を同定し、それを抽出する仕事である。母体血液における胎児細胞の濃度は胎児の妊娠の段階および状態に依存するが、見積もりは母体血液1ミリリットル毎に1ないし40の胎児細胞、または100,000母体有核細胞当たり1未満の胎児細胞の範囲である。現在の技術は母親の血液から少量の胎児細胞を単離することができるが、胎児細胞をいずれかの量の純度まで豊富化するのには非常に困難である。この関係での最も効果的な技術はモノクローナル抗体の使用を含むが、胎児細胞を単離するのに用いられる他の技術は密度遠心、成人赤血球の選択的溶解、およびFACSを含む。胎児DNA単離は、胎児-特異的DNA配列と共にプライマーを用いるPCR増幅を用いて示されてきた。各胚SNPの分子の10がこれらの技術を通じて利用可能なのに過ぎないので、高い忠実度での胎児組織の下のタイピングは現在可能ではない。

【 0 0 0 5 】

正常なヒトはジプロイド細胞毎に23染色体の2つの組を有し、1つのコピーは各親に由来する。異数性、余分なまたは失われた染色体を持つ細胞、および片親二染色体、一方の親に由来する2つの所与の染色体を持つ細胞は、着床の失敗、流産および遺伝病の大きなパーセンテージの原因であると考えられる。個体におけるある種の細胞のみが異数性である場合、該個体はモザイク現象を呈するといわれる。染色体異常の検出は成功した妊娠の確率の増大に加えて、とりわけ、ダウン症候群、クラインフェルター症候群およびターナー症候群のような疾患を持つ個体または胚を同定することができる。染色体異常についてのテストは母親の年齢のように特に重要であり；35歳および40歳の間では胚の40%および50%の間が異常であり、40歳を超えると、胚の半分を超えて異常であると見積もられる。

【 0 0 0 6 】

異数性およびモザイク現象の予測で用いられる伝統的な方法である核型分析は、他のより高いスループットのよりコスト的に有利な方法に対する途を開く。最近多大な注目を集めてきた1つの方法はフローサイトメトリー(FC)および蛍光イン・サイチュハイブリダイゼーション(FISH)であり、これを用いて、いずれかの相の細胞周期において異数性を検出できる。この方法の1つの利点は、それが核型分析よりも安価であるが、コストは、一般に、少し選択された染色体をテストするのでかなり十分である点である(通常、染色体13、18、21、X、Y；時々は8、9、15、16、17、22)；加えて、FISHは低いレベルの特異性を有する。15細胞を分析するのにFISHを用い、95%信頼性を持って19%のモザイク現象を検出することができる。テストの信頼性はモザイク現象のレベルが低くなるにつれ、および分析する細胞の数が減少するにつれかなり低くなる。テストが、対立の細胞を分析する場合、15%と高い擬陽性率を有すると見積もられている。より高いスループット、より低いコスト、およびより大きな精度を有する方法に対する多大な要望が存在する。

【 0 0 0 7 】

遺伝病の古典的な出生前診断に対する代替法としての着床前遺伝子診断(PGD)の使用に向けて多くの研究がなされてきた。ほとんどのPGDは、今日、異数性のような高レベルの染色体異常、および成功した着床およびテイク-ホームベイビーである主な結果を伴うバランスしたトランスロケーションに焦点を当てている。着床前段階における胚のより広範なゲノタイピングのための方法に対する要望が存在する。既知の病気に関連する対

10

20

30

40

50

立遺伝子の数は、現在、OMIMによると389であり、常に上昇している。その結果、病気表現型に関連する多数の胚SNPを分析するのは益々重要となりつつある。出生前診断よりも優れた着床前遺伝子診断の明瞭な進歩は、それが、一旦望ましくない表現型が検出されたならば、作用の可能な選択に関して倫理的論争のいくつかを回避する点にある。

【0008】

ゲノタイピング

単一の細胞を単離するための多くの技術が存在する。FACSマシンは種々の適用を有し；1つの重要な適用は、サイズ、形状および総じてのDNA含有量に基づいて細胞間を区別することである。FACSマシンは、単一細胞をいずれかの所望の容器に分類するように設定することができる。多くの異なるグループが、出生前遺伝子診断、組換え実験、および染色体不均衡の分析を含めた、多数の適用のために単一細胞DNA分析を用いてきた。単一 - 精子ゲノタイピングは、従前、精子試料の法医学分析で用いて（混合試料から生起する問題を減少させ）、および単一 - 細胞組換え実験のために用いられてきた。

10

【0009】

ヒト胚からの単一細胞の単離は、高度に技術的であるが、今日、体外受精クリニックにおいてルーチン的である。今日まで、出生前診断のほとんど大部分は、蛍光イン・サイチュハイブリダイゼーション(FISH)を用いており、これは、(ダウン症候群、またはトリソミー21のような)大きな染色体異常を決定することができ、およびPCR/電気泳動を用いてきており、これは少量のSNPまたは他の対立遺伝子の要求を決定することができる。極体および胚盤胞は共に成功して単離されてきた。胚の一体性を危うくすることなく単一の胚盤胞を単離するのは非常に重要である。最も普通の技術は、3日胚(6または8細胞段階)から単一の胚盤胞を取り出すことである。胚を特殊な細胞培養基(カルシウムおよびマグネシウムを欠如する標準培養基)に移し、酸性溶液、レーザー、または機械的なドリリングを用いて穴を透明帯に導入する。技術者は、次いで、バイオプシーペットを用いて、単一の目に見える核を取り出す。臨床的実験は、この目的は着床の成功を減少させないことを示している。というのは、この段階において、胚細胞は未分化だからである。

20

【0010】

全ゲノム増幅(WGA)に対して利用できる3つの主な方法がある：連結 - 媒介PCR(LM-PCR)、縮重オリゴヌクレオチドプライマーPCR(DOP-PCR)、および多数置換増幅(MDA)。LM-PCRにおいては、アダプターと呼ばれる短いDNA配列をDNAの平滑末端に連結する。これらのアダプターは普遍的増幅配列を含有し、これはPCRによってDNAを増幅するのに用いられる。DOP-PCRにおいては、普遍的増幅配列をやはり含有するランダムプライマーを第一ラウンドのアニーリングおよびPCRで用いる。次いで、第二ラウンドのPCRを用いて、普遍的プライマー配列をさらに持つ配列を増幅する。最後に、MDAはphi-29ポリメラーゼを用い、これは、DNAを複製する高度にプロセッシング可能な非特異的酵素であり、単一 - 細胞分析で用いられてきた。これらの方法のうち、DOP-PCRは、単一コピーの染色体を含めた、少量のDNAから多量のDNAを信頼性よく生産する。他方、MDAは最も速い方法であり、数時間以内にDNAの100折り畳み増幅を生産する。単一細胞からの増幅材料に対する主な制限は(1)極端に薄いDNA濃度または極端に小さな容量の反応混合物を用いる必要性、および(2)全ゲノムを横切って蛋白質からDNAを信頼性よく解離させる困難性である。それにもかかわらず、単一 - 細胞全ゲノム増幅は、何年もの間種々の適用に対して成功して用いられてきた。

30

40

【0011】

これらの関連でDNA増幅を用いるのに多数の困難がある。PCRによる単一 - 細胞DNA(または少数の細胞からの、またはより少量のDNAからのDNA)の増幅は、該ケースの5ないし10%において報告されているように完全に失敗しかねない。これは、しばしば、DNAの汚染、細胞の喪失、そのDNA、またはPCR反応の間におけるDNAの接近性である。増幅およびマイクロアレイ分析による胚DNAの測定で生じ得る誤差の

50

他の源は、特定のヌクレオチドがPCRの間に誤ってコピーされるDNAポリメラーゼによって導入される転写誤差、およびアレイ上での不完全なハイブリダイゼーションによるマイクロアレイのリーディング誤差を含む。しかしながら、最大の問題は、ヘテロ接合性細胞における2つの対立遺伝子のうちの一方を増幅できないことと定義される対立遺伝子ドロップアウト(ADO)のままである。ADOは増幅の40%を超えるまで影響しかならず、既に引き起こされたPGD誤診断を引き起こしてきた。ADOは特に優性病の症例において健康の論争となり、ここで、増幅できないことは侵された胚の着床に導きかねない。(ヘテロ接合体における)各マーカー当たり1を超えるプライマーの組に対する必要性はPCRプロセスを複雑とする。従って、より信頼性があるPCRアッセイがADO起源の理解に基づいて開発されつつある。単一細胞増幅のための反応容器は実験中である。アンプリコンのサイズ、DNA分解の量、凍結および解凍およびPCRプログラムおよび条件は、各々、ADOの速度に影響する。

10

【0012】

しかしながら、全てのそれらの技術は、単一細胞における増幅で利用可能なDNAの微量に依存する。このプロセスにはしばしば汚染が伴う。適当な滅菌条件およびマイクロサテライトサイジングは、汚染DNAの確率を排除することができる。というのは、出生前対立遺伝子においてのみ検出されるマイクロサテライト分析は汚染を排除するからである。対立細胞レベルまで分子診断プロトコルを信頼性よく導入する研究は、最近、マイクロサテライトマーカーの第一ラウンド多重PCR、続いての、リアルタイムPCRおよびマイクロサテライトサイジングを用いて追求されて、汚染の機会を排除してきた。多重PCRは単一細胞DNA分析における非常に重要な要件である単一反応における多数断片の増幅を可能とする。慣用的なPCRはPGDで用いられた最初の方法であるが、蛍光イン・サイチュハイブリダイゼーション(FISH)は今日普通である。乱れていない細胞および組織構築物内での拡散の検出を可能とするのはデリケートなビジュアルアッセイである。それは、先ず、分析すべき細胞の固定に依拠する。その結果、試料の固定および貯蔵条件の最適化が、特に、単一細胞懸濁液で求められる。

20

【0013】

単一細胞レベルでの多数の病気の診断を可能とする最新の技術は相間染色体変換、比較ゲノムハイブリダイゼーション(CGH)、蛍光PCR、および全ゲノム増幅を含む。これらの技術の全てによって得られたデータの信頼性は、DNA調製の質に依拠する。PGDは高価でもあり、その結果、ミニ配列決定のような安価なアプローチに対する要望が存在する。ほとんどの突然変異検出技術とは異なり、ミニ配列決定は低いADO率での非常に小さなDNA断片の分析を可能とする。増幅およびPGDについての単一細胞DNAを調製する良好な方法が従って求められており、研究されている。より新規なマイクロアレイおよび比較ゲノムハイブリダイゼーション技術は、依然として結局は、分析されるDNAの質に依拠する。

30

【0014】

いくつかの技術が、少数の細胞、単一細胞(例えば、胚盤胞)、少数の染色体のDNAについての、またはDNAの断片からの多数SNPを測定するために開発されている。ポリメラーゼ鎖反応(PCR)、続いてのマイクロアレイゲノタイピング分析を用いる技術がある。いくつかのPCRベースの技術は、多数置換増幅(MDA)、および単一对のプライマーでのPCRを用いて増幅することができる多数のタグドオリゴヌクレオチドを用いてゲノタイピングを行う分子逆転プロンプ(MIPS)のような全ゲノム増幅(WGA)技術を含む。非PCRベースの技術の例は蛍光イン・サイチュハイブリダイゼーション(FISH)である。該技術は、対立遺伝子ドロップアウト、不完全なハイブリダイゼーション、および汚染のような効果のインパクトを亢進するであろう限定された量の遺伝物質によりひどく誤差の傾向があることが明らかである。

40

【0015】

ゲノタイピングデータを供する多くの技術が存在する。TaqmanはApplied Biosystemsによって生産され、分配されるユニークなゲノタイピング技術で

50

ある。Taqmanはポリメラーゼ鎖反応(PCR)を用いて、注目する配列を増幅する。PCRサイクリングの間に、対立遺伝子特異的な従たる溝バインダー(MGB)は増幅された配列にハイブリダイズする。ポリメラーゼ酵素によるストランド合成はMGBプロープに連結されたレポーター色素を放出し、次いで、Taqman光学リーダーは色素を検出する。このように、Taqmanは定量的対立遺伝子区別を達成する。アレイベースのゲノタイピング技術と比較して、Taqmanは反応当たりかなり高価であり、(～\$0.40/反応)、およびスループットは比較的低い(実行当たり384遺伝子型)。反応当たり1ngのDNAが必要とされるに過ぎないが、Taqmanによる数千の遺伝子型はマイクログラム量のDNAを必要とし、従って、Taqmanは必ずしもマイクロアレイよりも少ないDNAを用いない。しかしながら、IVF遺伝子型ワークフローに関しては、Taqmanは最も容易に適用できる技術である。これはアッセイの高い信頼性および、最も重要なことには、アッセイのスピードおよび容易性のためである(実行当たりほぼ3時間、および最小の分子生物学工程)。また、(500k Affymetrixアレイのような)多くのアレイ技術とは異なり、Taqmanは高度に慣用化でき、これは、IVF市場で重要である。さらに、Taqmanは高度に定量的であり、従って、異数性はこの技術単独で検出できよう。

10

【0016】

Illuminaは、最近、高-スループットゲノタイピングにおけるリーダーとして出現した。Affymetrixとは異なり、Illuminaゲノタイピングアレイはハイブリダイゼーションに専ら依拠しない。その代わりに、Illumina技術が対立遺伝子-特異的DNA延長工程を用い、これは、元の配列の決定について、ハイブリダイゼーション単独よりもかなり感受性であって、特異的である。従って、これらの対立遺伝子の全てはPCRによって多重的に増幅され、次いで、これらの産物はビーズアレイにハイブリダイズされる。これらのアレイでのビーズはユニークな「アドレス」タグを含有し、天然配列を含有せず、従って、このハイブリダイゼーションは高度に特異的であって、感受性である。次いで、対立遺伝子がヘッドアレイの定量的スキャンニングによって呼ばれる。Illumina Golden Gateアッセイシステムは1536までの遺伝子座を同時に遺伝子型分けし、従って、スループットはAaqmanよりも良好であるが、Affymetrix 500kアレイほどは高くない。Illumina遺伝子型のコストはTaqmanよりも低い、Affymetrixアレイよりも高い。また、Illuminaプラットフォームは500k Affymetrixアレイと同程度完全となるまでには長くを必要とし(72時間まで)、これはIVFゲノタイピングでは問題である。従って、Illuminaはかなり良好なコールレートを有し、アッセイが定量的であり、従って、異数性がこの技術で検出可能である。Illumina技術が500k AffymetrixアレイよりもSNPの選択においてかなりフレキシブルである。

20

30

【0017】

一定時間において250,000SNPまでの測定を可能とする最高スループット技術の内の1つはAffymetrix GeneChip 500Kゲノタイピングアレイである。この技術はPCRをやはり用い、続いて、ハイブリダイゼーションによる分析、および水晶表面における異なる位置で化学的に合成されたDNAプロープに対する増幅されたDNA配列の検出を用いる。これらのアレイの不利は低いフレキシビリティおよびより低い感度である。「完全なマッチ」および「ミスマッチプロープ」のような選択性を増加させることができる修飾されたアプローチがあるか、これらはアレイ当たりのSNPコールの数を犠牲にしてそれを行う。

40

【0018】

パイロ配列決定、または合成による配列決定もまたゲノタイピングおよびSNP分析で用いることもできる。パイロ配列決定に対する主な利点は、極端に速いターンアラウンドおよび曖昧でないSNPコールを含むが、アッセイは、現在、高-スループット平行分析に導かれている。PCR、続いての、ゲル電気泳動は、着床前診断においてほとんどの成

50

功に適合したかなり単純な技術である。この技術において、研究者はネステッドPCRを用いて、注目する短い配列を増幅する。次いで、彼らは特殊なゲル上でこれらのDNA試料を実行して、PCR産物を可視化する。異なる塩基は異なる分子量を有し、従って、どれくらい速く産物がゲル中を泳動するかに基づいて塩基含有量を決定することができる。この技術は低-スループットであり、現行技術を用いる科学者による主題の分析を必要とするが、スピードの利点を有する(1ないし2時間のPCR、1時間のゲル電気泳動)。この理由で、それは、セラセミア、神経線維腫症2型、白血球接着欠乏症I型、アロポー-シーメンス病、鎌状細胞貧血、網膜芽細胞腫、ペリツェーウス-メルツバッヒャー病、ドゥシェーヌ筋ジストロフィー、およびクラリノ症候群を含めた、膨大な病気についての出生前ゲノタイピングで従前用いられてきた。

10

【0019】

非常に高い忠実度でもって少量の遺伝物質を遺伝子型分けするために開発されたもう1つの有望な技術は、Affymetrix's Genflexアレイのような分子逆転プローブ(MIP)である。この技術は、平行して多数のSNPを測定する能力を有し; 平行して測定された10,000を超えるSNPsが証明されている。少量の遺伝物質については、この技術についてのコールレートは概略95%において確立されており、なされたコールの精度は99%を超えることが確立されている。これまで、該技術は所与のSNPについて150分子と小さなゲノムデータの量について実行されてきた。しかしながら、該技術は、着床前遺伝子診断について要求されるように、単一細胞、またはDNAの単一ストランドからのゲノムデータで証明されてきた。

20

【0020】

MIP技術は、その2つの端部が、それらがDNAの直ちに隣接する標的配列にハイブリダイズする場合に連結によって接合できる線状オリゴヌクレオチドであるパドロックプローブを用いる。プローブがゲノムDNAにハイブリダイズされた後に、ギャップを満たす酵素をアッセイに加え、これは4つのヌクレオチドの内1つをギャップに加えることができる。もし加えられたヌクレオチド(A, C, T, G)が測定下でSNPに対して相補的であるならば、それはDNAにハイブリダイズし、連結によってパドロックプローブの端部を接合するであろう。次いで、管状産物、または閉じたパドロックプローブをエキソヌクレオリシスによって線状プローブから区別される。エキソヌクレアーゼは、線状プローブを分解し、環状プローブを残すことによって、千倍以上だけ、閉じた-vs-閉じていないプローブの相対的濃度を変化させるであろう。次いで、残ったプローブをもう1つの酵素によって切断部位において開き、DNAから取り出し、PCRによって増幅する。各プローブは20塩基タグよりなる異なるタグ配列が付され(16,000が作り出されている)、例えば、Affymetrix Genflexタグアレイによって検出することができる。特定のギャップを満たす酵素が加えられた反応からのタグドプローブからの存在は、関連SNP上での相補的アミノ酸の存在を示す。

30

【0021】

MIPSの分子生物学利点は:(1)単一反応における多重ゲノタイピング、(2)遺伝子型「コール」はギャップを満たし連結することによって起こるが、ハイブリダイゼーションによっては起こらない、および(3)ユニバーサルタグのアレイへのハイブリダイゼーションは、ほとんどのアレイハイブリダイゼーションに固有な偽陽性を減少させることを含む。伝統的な500k、TaqManおよび他のゲノタイピングアレイにおいて、全ゲノタイプ試料はアレイにハイブリダイズされ、これは種々の完全なマッチおよびミスマッチプローブを含有し、アルゴリズムはミスマッチおよび完全なマッチプローブの強度に基づく遺伝子型を要求するようである。しかしながら、DNA試料の複雑性、およびアレイ上での膨大な数のプローブのため、ハイブリダイゼーションは固有にノイズがある。他方、MIPは、より長く、従って、より特異的であり、従って、プローブを環状化するのに頑強な連結工程を用いる多重プローブを用いる(すなわち、アレイ上にはない)。対立遺伝子ドロップアウトは(貧弱な実行プローブのため)高いであろうが、バックグラウンドは(特異性のため)このアッセイにおいてはかなり低い。

40

50

【 0 0 2 2 】

この技術を単一細胞（または少数の細胞）からのゲノムデータで用いる場合、それは、PCRベースのアプローチのように、一体性の争いに悩んでいる。例えば、パドロックブロープがゲノムDNAにハイブリダイズできないことは、対立遺伝子ドロップアウトを引き起こすであろう。これは体外受精の関係で悪くなるであろう。というのは、ハイブリダイゼーション反応の効率は低く、かつそれは相対的に速く進行して、限定された時間内に胚を遺伝子型分けする必要があるからである。ハイブリダイゼーション反応は販売業者が推奨するレベルよりも十分低く減少でき、マイクロ-流動技術を用いて、ハイブリダイゼーション反応を加速することもできる。ハイブリダイゼーション反応のための時間を減少させることに対するこのアプローチは減少したデータの質を引き起こすであろう。

10

【 0 0 2 3 】

予測ゲノミックス

一旦遺伝子データが測定されれば、次の工程が予測目的でデータを用いることである。多くの研究が予測ゲノミックスにおいてなされ、これは、表現型予測を遺伝子型に基づいてなすことができるように、蛋白質、RNAおよびDNAの正確な機能を理解することを試みる。カノニカル技術は単一-ヌクレオチド多形(SNP)の機能に焦点を当てるが、より進歩した方法は多因子表現型特徴を担うようにされつつある。これらの方法は、遺伝子および表現型予測の組、および測定された結果の組の間の数学的関係を決定するように試みる、直線回帰および非直線神経ネットワークのような技術を含む。また、遺伝子データに典型的なように、結果の数に対して多くの潜在的プレディクターが存在し、データが過少決定される場合でさえパラメーターの重要な組を解決することができるように、さらなる制限を回帰パラメーターに適應するまばらなデータ組を収容するように設計されたRidge回帰、log回帰および段階的選択のような回帰分析技術の組もある。他の技術は、未決定データ組から情報を抽出するために主な成分分析を適用する。決定ツリーおよび偶発性の表のような他の技術は、それらの独立した変数に基づいて主題を細分化して、主題を、表現型結果が同様であるカテゴリーまたはピンに入れるための戦略を用いる。論理的回帰といわれる最近の技術は、カテゴリー的に独立した変数の間の異なる論理的相互関係についてサーチして、遺伝子データに関連する多数の独立変数の間の相互作用に依存する変数をモデル化する方法を記載している。用いる方法に拘わらず、予測の質は、予測をなすのに用いる遺伝子データの質に自然に高度に依存する。

20

30

【 0 0 2 4 】

DNA配列決定のコストは迅速に低下しており、近い将来において、個人の利益のための個々のゲノム配列決定はより普通になるであろう。個人的遺伝子データの知識は、広範な表現型予測が個人に対してなされるのを可能とするであろう。正確な表現型予測をなすためには、関係を問わず、高い質の遺伝子データが非常に重要である。出生前または着床前遺伝子診断の場合には、複雑化因子は入手可能な遺伝物質の相対的少量である。限定された遺伝物質をゲノタイピングで用いる場合に、測定された遺伝子データの性質に固有にノイズがあると仮定すれば、一次データの忠実度を増大させ、それをクリーンとできる方法に対する多大な要望が存在する。

【 0 0 2 5 】

臨床的決定がなされる現行の方法は、存在する情報の最良な可能な使用を行わない。医療的、生化学的および情報技術の進歩としては、増大した量のデータが作り出され、アカデミックおよび臨床的実験の関係においての個々の患者について双方を貯蔵する。分析で利用可能な遺伝子、表現型および臨床的情報の量における最近の急増に従い、臨床的に関連する相関関係を見出して、人々がより長く、より健康でかつよりエンジョイできる人生を送るのを助けるのに多大の努力が払われてきた。従前には臨床家および研究者は彼らの分析を少量の明らかな潜在的因子に焦点を当て、データの局所的貯蔵を用いるが、他の剤のスコアによって測定されたデータを活用することができ、および所与の遺伝子型または表現型に相関する従前に疑われていない因子を同定することができるより複雑なモデルを用いる潜在的利点がより明瞭になりつつある。この状況は、一旦個人的な遺伝子データが

40

50

病気の原因および治療、および対象の他の素因を理解するにおいてより抽象的な役割を占めれば、かなりより複雑になるであろう。次の10年以内に、臨床試験のために、または個人化された治療およびまたは薬物割当ての目的のために、患者の全ゲノムをスキャンし、ならびに膨大な表現型データ点を収集するのが可能であろう。

【0026】

利用可能なデータの量が膨大となり、それが依然として迅速に増大するにつれ、問題の最も重要な点は、最も適当な関係が発見し、かつそれを用いて人々に役に立つのを可能とする設計および実行する良好な方法となった。分析するのに利用可能な変数の数が増大するにつれ、天文学的数の潜在的関係を会得でき、先見的にそれらのいずれかを除外しない方法を開発するのがより重要となった。同時に、それらの研究を同一プロトコルで実行しなかつた場合でさえ、多数の研究の知見を総合し、それを利用することができる方法を開発するのが重要である。また、所与の分析において用いるために最適な方法を正しく同定することができるシステムを開発するために、研究されてきた非常に多数の予測モデルを仮定すれば、それは益々重要になりつつある。

10

【0027】

HIVの関係におけるバイオインフォマティクス

HIVは三千万を超える人々が現在HIVに罹って生きているヒトにおいてHIVは広域病と考えられ、毎年二百万を超える死亡がHIVに帰せられている。HIVの主な特徴の1つはその速い複製サイクル、および逆転写酵素の高い誤差率および組換え原性の結果としてのその高い遺伝子可変性である。その結果、HIVウイルスの種々の株は異なるレベルの異なる薬物に対する耐性を示し、最適な治療養生法は感染性株の同一性およびその特別な罹患性を考慮することができる。

20

【0028】

今日まで認可されたART薬物は11のRTI：7のヌクレオシド、1つのヌクレオチド、および3つの非ヌクレオシド；7つのPI；および1つの融合/エンター阻害剤のリストよりなる。世界中でのART薬物が現在広く行きわたっていることを仮定すれば、ウイルスの耐性株の出現は、耐性に対する低い遺伝子バリア、および貧弱な薬物固執双方のため不可避的である。その結果、どのようにして突然変異したウイルスが抗-レトロウイルス療法に应答するかを予測する技術は益々重要となっている。というのは、それらはサルベージ療法についての結果に影響するだろうからである。ウイルス遺伝子配列決定の迅速に現象しているコスト-予備的に調製された配列については5ドルと低い容量価格は、よりコストがかかりかつ関連するイン-ビトロ表現型測定よりはむしろ、ウイルス遺伝子配列データに基づく薬物の選択を魅力的なオプションとする。しかしながら、配列データの使用はウイルス遺伝子突然変異の出現に基づく、ウイルス薬物应答の正確な予測を必要とする。ウイルス突然変異の多くの異なる組合せは、全ての遺伝子補因子およびそれらの相互作用を含むモデルを設計し、限定されたデータでもってモデルを訓練するのを困難とする。後者の問題は、薬物養生法の多くの異なる組合せが、変数、すなわち、ベースライン臨床状態、処置履歴、臨床的結果および遺伝子配列を含有するいずれかの特定の養生法について十分に大きなデータ組を収集するのを困難とする場合に、イン-ビボ薬物应答をモデル化する関係が悪化した。

30

40

【0029】

抗ウイルス薬物に対する耐性は、RTまたはプロテアーゼ配列内の1つの突然変異、または複数の突然変異の組合せの結果であり得る。RT酵素は560コドンの鍵となる組によってコードされ；プロテアーゼ酵素は99のコドンによってコードされる。アミノ酸を改変する突然変異のみをコードすることによって各アミノ酸遺伝子座は19の可能な突然変異を有し；従って、RT酵素について野生型とは異なる合計10,640の可能な突然変異、およびプロテアーゼ酵素についての1,981の可能な突然変異がある。単純な直線モデルを用い、データで総合した各突然変異(全ての突然変異が起こるのではない)が特定の重み付け、または直線回帰パラメーターと関連させる場合、数千のパラメーターが存在し得る。もし数百人の患者の試料のみが各薬物で利用できるならば、問題は過剰決定

50

的であるか、またはHadamardの意味において不適切である。というのは、独立した方程式よりも評価するより多くのパラメータがあるからである。不適切な問題のためにモデル構築する問題に適用することができる多くの技術が存在する。これらが先見的専門知識を観察と組み合わせ、専門家のルールに基づくシステム、ならびにi)リッジ回帰、ii)主要成分分析、iii)決定ツリー、iv)段系的選択技術、v)神経ネットワーク、vi)最小絶対収縮および選択オペレーター(LASSO)およびvii)Support Vector Machines(SVM)を含めた統計的方法を作り出すことを含む。

【0030】

3つの主な産業 - 標準専門家システムを典型的に用いて、ART薬物へのHIVウイルスの罹患性：ANRS-AC11システム、Regaシステム、およびStanford HIVdbシステムを予測する。新しいアルゴリズムがこれらの専門家システムに対して評価されるのは文献において通常である。しかしながら、これらの専門家システムのいずれも、表現型応答の直接的予測を行うように設計されていないが、むしろ、異なる薬物をそれにより比較することができる数値スコアを供し、または感受性、中程度および耐性のような区別されるグループに薬物を分類するように設計されている。加えて、段階的選択でもって訓練された直線回帰モデルのような統計学的アルゴリズムは、表現型結果の予測において専門家システムを実質的に凌ぐことが明瞭に確立されている。結果として、統計学的技術の組のみが、文献に最近開示された方法を最良に実行することを含む詳細な記載中の新規な方法と比較される。

【0031】

サルベージARTの臨床的結果の予測に対する現在のアプローチは、薬物養生法および遺伝子突然変異の多くの異なる順列と組み合わせた、ほとんどは、統計学的に有意な結果のデータの欠如のため、良好な予測パワーを示さない。この分野は多数の不均一なデータ組の一体化、および薬物応答予測の増強の双方のための緊急の要望を有する。

【0032】

癌の関係でのバイオインフォマティクス

見積って80,000の年次臨床試験のうち、2,100は癌薬物のためである。癌療法のための危険性および利点をバランスさせることは、表現型および遺伝子型情報の組合せ使用についての臨床的先駆者を表す。過去数十年において化学療法で大きな進歩があったにもかかわらず、腫瘍学者は彼らの癌患者を、癌細胞について正常な細胞に対してしばしば毒性である原始的全身薬物で依然として治療している。かくして、化学の最大毒性用量および治療用量の間に微妙な線がある。さらに、用量 - 制限毒性は、他の患者ではなくある患者においてよりひどく、治療運動をより高くまたはより低くシフトさせ得る。例えば、乳癌治療で用いられるアントラサイクリンは有害な心血管事象を引き起こしかねない。現在、もし患者が心臓病に対して低い危険性であると決定できても、治療ウィンドウをより大きな用量のアントラサイクリン療法を可能とするようにシフトできたとしても、全ての患者はあたかも心血管毒性の危険性があるように治療される。

【0033】

各患者についての化学療法の利点および危険性をバランスさせるために、副作用のプロフィール、および医薬介入の治療的有効性を予測することができる。癌療法は、しばしば、ユニークな宿主および腫瘍遺伝子型についての不適切な調整のため失敗する。単一の多形は、稀には、薬物応答において有意な変動を引き起こし；むしろ、マニフォールド多形の結果ユニークな生体分子組成物をもたらす、臨床的結果の予測を困難とする。「ファルマコゲネティクス」は、広く、遺伝子変異が薬物に対する患者の応答に影響する方法と定義される。例えば、肝臓酵素における天然の変異は薬物代謝に影響する。癌化学療法の将来は標的化医薬であり、これは、癌を、多数の遺伝子的、分子的、細胞的、および生化学的異常を含む病気プロセスとして理解する必要がある。酵素 - 特異的薬物の出現に伴い、腫瘍が特異的にまたは正常な組織よりも高いレベルで分子標的を発現することを確実にするために注意することができる。腫瘍細胞および健康な細胞の間の相互作用を考慮する

10

20

30

40

50

ことができる。というのは、患者の正常な細胞および酵素は腫瘍薬物の曝露を制限でき、または有害な事象をよりありそうにしかねないからである。

【 0 0 3 4 】

バイオインフォマティクスは癌治療に大変革を起こさせ、仕立てられた治療が利点を最大化し、有害な事象を最小化するのを可能とする。応答を予測するのに用いられる機能的マーカーはコンピュータアルゴリズムによって分析することができる。乳癌、結腸癌、肺癌および前立腺癌は4つの最も普通の癌である。これらの癌に対する2つの治療の例は乳癌を治療するのに用いられるタモキシフェン、および結腸癌患者において用いられるイリノテカンである。タモキシフェンまたはイリノテカンも、各々、乳癌または結腸癌を治療するのに必要でなく、または十分でない。癌および癌の治療は、患者の副作用のプロファイルおよび腫瘍応答に従って、療法の改正および、しばしば、組合せ療法を必要とする動的なプロセスである。もし癌治療を決定的なツリーとイメージして、他の療法の前、後またはそれと共にいずれかの1つの治療を与え、またはそれを差し控えるならば、このツリーは決定決断点のサブセットを含み、そこではツリーの多く(すなわち、他の治療)はブラックボックスと考えることができる。それにも拘わらず、医師を最も効果的な治療に部分的にガイドするためのデータを有することは有益であり、より多くのデータを集めるに従い、このデータに基づいて治療の決定を行うための効果的な方法は数千人の癌患者において平均余命および生活の質を有意に改善することができよう。

10

【 0 0 3 5 】

結腸または大腸は胃腸(GI)管の最後の6-フットのセクションである。合衆国癌協会は、結直腸癌の145,000の症例が2005年において診断され、56,000人が結果として死亡するであろうと見積もっている。結直腸癌はグレード、または細胞の異常、および段階について評価され、これは腫瘍のサイズ、リンパ節の関与、および遠い転移の存在または不存在に細分化される。結直腸癌の95%は、結腸のルーメンをライニングする遺伝子的突然変異体上皮細胞から発生する腺癌である。症例の80ないし90%において、外科的処置単独が看護の標準であるが、転移の存在は化学療法を必要とする。転移性結直腸癌に対する多くの一次療法の1つは5-フルオロウラシル、ロイコボリン、およびイリノテカンの養生法である。

20

【 0 0 3 6 】

イリノテカンは、スーパーコイルドDNAの絡みを解いて、DNA複製が分裂細胞において進行するようにし、細胞をアポトーシスに対して感受性とするトポイソメラーゼを阻害するカンプトテシンアナログである。イリノテカンは生物学的経路において明確な役割を有さず、従って、臨床的結果は予測するのが困難である。用量-限定的毒性はひどい(グレードIIIないしIV)下痢および骨髄抑制を含む、その双方は直ちに医療的注意を必要とする。イリノテカンはウリジンニリン酸グルコロノシルトランスフェラーゼイソ形態1a1(UGT1A1)によって活性な代謝産物であるSN-38に代謝される。UGT1A1における多形はGIのひどさ、および骨髄副作用と関連する。

30

【 0 0 3 7 】

先行技術

本明細書中において、本発明の分野に関連する先行技術の組をリストする。この先行技術はいずれも、本発明の新規なエレメントを含まず、または断じてそれに言及しない。特許文献1において、Hartleyらは、作製された組換え部位および組換え蛋白質を用いてDNA分子のセグメントを移動させ、または交換する組換えクローニング方法を記載する。特許文献2において、Parrottらは、生体活性脂質のレベルについて体外受精培養の培地検体を分析して、当該特徴を決定することによって、総じての胚の健康、着床性、および出産予定日まで成功して発生する増大した尤度を含めた体外受精胚の種々の生物学的特徴を決定する方法を提供する。特許文献3において、Threadgillらは、複数の単離された親細胞における部位-特異的有糸分裂組換えに関連するイン・ビトロフェノタイピングおよび遺伝子マッピングで言うようなホモ接合性細胞ライブラリーを調製する方法を記載する。特許文献4において、Stewartらは、血清において直接

40

50

的に、またはIVF/ET手法の一部として患者から抽出された顆粒膜黄体細胞を培養することによって間接的にレラキシンを測定することによって成功する体外受精（IVF）の確率を決定する方法を記載する。特許文献5において、Cookeらは、女性患者からの生物学的試料中の11-ヒドロキシステロイドデヒドロゲナーゼのレベルを測定することによってIVFの結果を予測する方法を提供する。特許文献6において、Larderrらは、神経ネットワークを用いて、療法剤に対する病気の抵抗性を予測する方法を記載する。特許文献7において、Vingerhoetsらは、所与のHIV株のインテグラーゼ遺伝子型を、関連表現型と共にHIVインテグラーゼ遺伝子型の公知のデータベースと単純に比較して、マッチング遺伝子型を見出す方法を記載する。特許文献8において、Dentonらは、個人のハプロタイプを一般的集団におけるハプロタイプの公知のデータベースと比較して、治療に対する臨床的応答を予測する方法を記載する。特許文献9において、Schadtらは、遺伝子マーカーのマップを構築し、個人の遺伝子および特性を分析して遺伝子-特性遺伝子座データを与え、次いで、これを遺伝子的に相互作用する経路を同定するための方法としてクラスター化し、これを多変数分析を用いて確認する方法を記載する。特許文献10において、Veltriらは、パラメーターとしてバイオマーカーのコレクションを利用して、前立腺癌の再発の危険性を評価する神経ネットワークの使用を含む方法を記載する。特許文献11において、Mascarenhasは、患者についての生化学的プロフィールを確立し、テストコフォルトのメンバーにおいて応答性を測定し、次いで、患者の生化学的プロフィールのパラメーターを個々にテストして、薬物応答性の尺度との相関性を見出すことによって薬物応答性を予測する方法を記載する。

【特許文献1】米国特許第6,720,140号明細書

【特許文献2】米国特許第6,489,135号明細書

【特許文献3】米国特許出願公開第2004/0033596号明細書

【特許文献4】米国特許第5,994,148号明細書

【特許文献5】米国特許第5,635,366号明細書

【特許文献6】米国特許第7,058,616号明細書

【特許文献7】米国特許第6,958,211号明細書

【特許文献8】米国特許第7,058,517号明細書

【特許文献9】米国特許第7,035,739号明細書

【特許文献10】米国特許第6,025,128号明細書

【特許文献11】米国特許第5,824,467号明細書

【発明の概要】

【課題を解決するための手段】

【0038】

（発明の要旨）

開示するシステムは、情報の源として二次的遺伝子データを用い、またその遺伝子データを用いて、表現型および臨床的予測をする、不完全またはノイズがある遺伝子データの清浄化を可能とする。開示はヒト対象からの遺伝子データに焦点を当てているが、開示する方法は関連する範囲において生物の範囲の遺伝子データに適用されることは注意すべきである。遺伝子データを清浄化するために記載する技術は、体外受精の間の着床前診断、羊水穿刺と組み合わせた出生前診断、絨毛膜バイオプシー、および胎児血液サンプリング、および非侵襲性出生前診断との関係で最も関連し、ここで、少量の胎児遺伝物質は母体血液から単離される。診断は遺伝病、欠点または異常の増大した尤度、ならびに臨床的およびライフスタイルの決定を促進するための個体についての表現型予測を行うことに焦点をあてることができる。本発明は、先に議論された先行技術の欠点に取り組む。表現型および臨床的予測を行うための本明細書中に記載された技術は、着床前診断、出生前診断との関係、または医療的疾患、または罹患性を持つ個人を含めた、多数の関係で関連する。本明細書中に開示される技術のある実施形態は、個体についての遺伝子、表現型および/または臨床的情報の組を仮定し、個体についての表現型結果または表現型罹患性の性格な予測を行うためのシステムを記載する。1つの態様において、遺伝子データに典型

的なように、測定された結果の数と比較して多くの潜在的予測が存在する場合に表現型を正確に予測することができる線形および非線形回帰モデルを形成するための技術が開示され；本発明のもう1つの態様において、該モデルは分割表に基づき、パブリックドメインで入手可能な情報から形成される。なおもう1つの発明において、システムが記載され、ここで、多数のモデルが関連データセットで訓練され、関連予測を行うのに最も正確なそのモデルを用いる。

【0039】

本発明の1つの態様において、方法は、減数分裂のメカニズムの知識、および胚DNAの不完全な測定と共に、母親および父親の遺伝子データの不完全な知識を用いて、高度な信頼性をもって鍵となるSNPの位置において胚DNAをイン・シリコにて再構築する。親データは、貧弱に測定されたSNPのみならず、挿入、欠失、およびSNP、または全く測定されなかったDNAの全領域の再構築を可能とすることに注意するのは重要である。

10

【0040】

開示された方法は体外受精との関係で適応でき、ここで、着床についてコードされる各胚からのゲノタイピングで利用できる。開示された方法は、少数の胎児細胞、または胎児DNAの断片のみが母親の血液から単離されている非侵襲性出生前診断(NIPD)の關係に等しく適応できる。開示された方法は、羊水穿刺の場合、および胎児の血液が直接的にサンプリングされる他の方法において等しく適応可能である。開示された方法は、限定された量の遺伝子データが標的個人から入手でき、およびさらなる遺伝子データが標的に遺伝的に関連する個体から入手できるいずれの場合においてもより一般的に適用可能である。

20

【0041】

本発明の1つの態様において、再構築された胎児または胚ゲノムデータを用いて、細胞が異数性であるか、すなわち、少数の、または2を超える特定の染色体が細胞に存在するかを検出することができる。この疾患の普通の例はトリソリン-21であり、これはダウン症候群を生起させる。再構築されたデータを用いて、所与の染色体の2つが存在し、その双方が1つの親に由来する疾患である片親二染色体についても検出することができる。これは、DNAの潜在的状態についての仮説の組を創製し、いずれの1つが測定されたデータを仮定して真実である最高の確率を有するかを見るためにテストすることによってなされる。異数性をスクリーニングするための高スループットゲノタイピングデータの使用は、各胚からの単一の胚盤胞が多数病気-関連遺伝子座を測定し、ならびに染色体異常についてスクリーニングする双方で用いられるのを可能とするのに注意されたし。

30

【0042】

本発明のもう1つの態様において、複数の遺伝子座に存在する増幅されたまたは増幅されていない遺伝物質の量の直接的測定を用いて、異数性、または片親二染色体について検出することができる。この方法の背後にある考えは、単に、増幅の間に存在する遺伝物質の量は初期試料における遺伝子情報の量に比例し、多数の遺伝子座においてこれらのレベルを測定することは統計学的に有意な結果を与えることである。染色体異常についてスクリーニングする方法は、遺伝子データを清浄化するための本明細書中に記載された関連方法と組合せて用いることができる。

40

【0043】

本発明のもう1つの態様において、開示された方法は、外来性遺伝物質によって生じたデータを同定することにより外来性DNAまたはRNAに汚染されている個体の遺伝物質を清浄化できる。汚染DNAによって生じた偽シグナルは、異数性によって生じた染色体-幅特異的シグナルを検出できる方法と同様に認識することができる。

【0044】

本発明のもう1つの態様において、標的細胞が単離され、これらの細胞に含有される遺伝子データが増幅され、以下の技術：PCR-ベースの増幅技術、PCR-ベースの測定技術、または分子逆転プローブに基づく検出技術、またはGeneChipまたはTa q

50

Manシステムのようなマイクロアレイのうちの1以上の組合せを用いて多数SNPの測定を行う。次いで、この遺伝子データを本明細書中に記載されたシステムで用いる。

【0045】

本発明のもう1つの態様において、双方の親からのジプロイドおよびハプロイドデータを用いて、個体の遺伝子データを清浄化できる。別法として、親からのハプロイドデータは、もし親のジプロイドおよびハプロイドデータを測定することができれば、シミュレートすることができる。もう1つの態様において、個体に対する公知の遺伝子関連のいずれかの個人からの遺伝子データを用いて、親、兄弟姉妹、祖父母、子孫、従兄弟、叔父、叔母などを含めた、個体のデータを清浄化することができる。

【0046】

本発明のもう1つの態様において、標的および/または関連個体の遺伝子データはイン・シリコにて部分的にまたは全体的に知ることができ、いくつかの直接的測定の必要性を軽減する。遺伝子データの部分は、隠れたMarkovモデルを利用するインフォーマティックアプローチによってイン・シリコにて作り出すことができる。

【0047】

本発明の1つの態様において、SNPの決定における信頼性を見積もることが可能である。

【0048】

本明細書中に記載された技術は、1つの、または少数の細胞における遺伝物質の測定、ならびに非侵襲性出生前診断(NIPD)との関係で母親の血液から単離することができるもののようなより少量のDNAについての測定の双方に関連することに注意されたし。また、この方法はイン・シリコでの、すなわち、遺伝物質から直接的に測定されないゲノムデータに等しく適応することができる。

【0049】

本発明の1つの態様において、OMIM(男性におけるオンラインメンデル遺伝)データを介するように刊行物を介して、およびHapMapプロジェクトおよびヒトゲノムプロジェクトの他の態様から入手可能なデータを用いて入手可能なデータから構築することができる分割表に基づいてモデルを作り出すための技術が提供される。この技術のある実施形態は、モデルの予測的精度を改良するために、遺伝子の間の関連についての、および遺伝子および病気の間に関連についての出現する公のデータを用いる。

【0050】

なおもう1つの態様において、最良のモデルを、特定の患者で利用できるデータで見出すことができる技術を開示する。この態様において、多くの異なるモデリング技術と共に、変数の多くの異なる組合せを調べることができ、他の対象からのテストデータと共に交差・確認に基づいて個々の対象についての最良の予測を生じるであろうその組合せを選択することができる。

【0051】

いくつかの場合において、個体についての表現型の結果または表現型の感受性の正確な予測を行うにおいて最良のものを生じさせることができるモデルを、凸最適化技術を用いて訓練して、データの特定の組についての全体的に最適なパラメータを見出すのが保証されるように、プレディクターの連続的サブセット選択を行う。この特徴は、モデルが複雑であり得、遺伝子突然変異または遺伝子発現レベルのような多くの潜在的プレディクターを含有することができる場合に特に有利である。さらに、いくつかの例においては、それらが単純な方法でデータを説明するように、凸最適化技術を用いて、モデルを希薄とすることができる。この特徴は、モデルにおける潜在的プレディクターの数が、訓練データにおける測定された結果の数と比較して大きい場合でさえ、訓練されたモデルが正確に一般化されるのを可能とする。同様な技術は学問的雑誌に公表されている(Rabinowitz, M.ら, 2006, "Accurate Prediction of HIV-1 drug response from the reverse transcriptase and protease amino acid sequence

10

20

30

40

50

s using sparse models created by convex optimization.” Bioinformatics 22(5):541-9) 。この論文からの情報は背景および文脈のために本書類に含めてあることに注意されたし。

【0052】

本明細書中に開示されたある説明的実施形態はヒト対象からの遺伝子データに焦点を当て、癌またはHIVにかかった人々についての、またはアルツハイマー病または心筋梗塞のような病気に対する彼らの罹患性を理解したい人々についての特別な実施形態を提供するが、開示された方法は多数の異なる関係の範囲において生物の範囲の遺伝子データに適用されるのに注意すべきである。表現型予測および薬物応答予測について本明細書中に記載された技術は、種々の癌、遺伝子病、細菌、真菌またはウイルス感染の治療との関係で、並びに臨床的およびライフスタイルの決定を促進するために個体について表現型予測を行うにおいて関連し得る。さらに、該システムを用いて、遺伝子データ、具体的にはIVFとの関係で胚(着床前)の、または羊水穿刺を含めた非侵襲性または侵襲性出生前診断との関係で胎児のSNP(単一ヌクレオチド多形)データを仮定し、特定の表現型結果の尤度を決定することができる。

10

【0053】

1つの実施形態において、予測的モデルを、標準化された計算可能なフォーマットで貯蔵されている特定の個人についての遺伝子データに適用することができる。個人は、彼らに関連する特定の論点を記載することができ、あるいはシステムは、いずれの表現型罹患性がその個体に関連するかを自動的に決定することができる。新しい研究データが病気-遺伝子関連、治療、またはライフスタイルの嗜好性について入手できるようになるので、個体には、集合されたゲノムおよび臨床データから開発された予測的モデルに基づいて、彼らの決定および嗜好性についてのこの情報のインパクトを知らせることができる。別法として、該システムは新しい研究データを用いて、個体についての疑われていない危険性をここに検出することができ、その個体にはこの情報のインパクトを知らせることができる。

20

【0054】

もう1つの実施形態において、遺伝子データ、表現型データおよび関連診断テストを含めた臨床記録のデータベースから一体化されたデータについて訓練された結果予測モデルを用いて臨床家のために増強された報告を作成することができる。このシステムは、限定されるものではないが、HIV、癌、アルツハイマー病および心臓病を含めた、病気および/または病気素因を持つ個体についての増強された報告の創生を提供できる。この増強された報告は治療する医師に、いずれの病気-管理または予防的処置が与えられた個体についてより適当であるか、またはあまり適当でないであろうことを示すであろう。報告は、集合された対象データについて訓練されたモデルを用いるその個体についての鍵となる結果についての予測および信頼性限界を含むであろう。

30

【0055】

もう1つの実施形態によると、特定の個体についてのデータを用いて、分割表に基づき、パブリックドメインで入手可能な情報から形成されたモデルを用いて該個体についての予測を行い、該データは該個体の遺伝子データ、該個体の表現型データ、および個体の臨床データ、およびその組合せよりなる群から取られ、ここで、該予測は該個体の表現型、表現型罹患性および可能な臨床的結果を含む群から取られたトピックに関し、およびここで、該情報は、遺伝子型-表現型関連についての情報、ある遺伝子対立遺伝子の頻度についての情報、遺伝子対立遺伝子内のある関連の頻度についての情報、遺伝子対立遺伝子のある実施形態を仮定したある表現型の1以上の状態の確率についての情報、ある表現型の状態を仮定した遺伝子対立遺伝子のある組合せの確率についての情報、およびその組合せを含む群から取られるシステムおよび方法が開示される。

40

【0056】

なおもう1つの実施形態によると、それにより、特定の個体についてのデータを用いて

50

、最良の精度を示すモデルを利用できるように集合データについて訓練された種々の数学的モデルを用い該個体についての予測を行うことができ、ここで、該個体のデータは該個体の遺伝子データ、該個体の表現型データ、および該個体の臨床的データよりなる群から取られ、およびここで、該予測は該個体の表現型、表現型罹患性、可能な臨床的結果、およびその組合せから取られるトピックスに関連するシステムおよび方法が提供される。ある実施形態において、該方法は、多数のモデルおよび多数のチューニングパラメータを用いて、データの所与の組において異なる独立した変数および従属した変数の組合せの多くまたはすべてを調べることができ、次いで、最良の表現型予測を行う目的でテストデータにて最高の相関係数を達成した独立した変数および従属した変数およびその組合せ、そのモデルおよびそれらのチューニングパラメータを選択する。

10

【0057】

もう1つの実施形態によると、本明細書中に開示された方法のいずれも予測を用いて、該個体に関連する1以上のトピックスに関連する特定の個体についての報告を作成することができ、ここで、該トピックスはライフスタイルの決定、ダイエットの嗜好性、ホルモンサプリメント、病気についての可能な治療養生法、病原体に対する可能な治療養生法、薬物介入、およびその組合せを含む群から取られ、およびここで、該予測は該個体の遺伝子メイキャップ、該個体の表現型特徴、該個体の臨床的履歴およびその組合せに関連するデータに基づく。

【0058】

他の実施形態によると、本明細書中に開示された方法のいずれも予測を用いて、医師または臨床家のような特定の個人の代理人のための報告を作成することができ、ここで、該予測は該個体に関連する情報を供することによって該代理人を助けることができ、およびここで、該情報の主題はライフスタイルの決定、ダイエットの嗜好性、ホルモンサプリメント、病気についての可能な治療養生法、病原体についての可能な治療養生法、薬物介入、他の治療的介入、およびその組合せを含むトピックスの群から取られ、およびここで、該予測は該個体の遺伝子メイキャップ、該個体の表現型特徴、該個体の臨床的履歴およびその組合せに関するデータに基づく。

20

【0059】

もう1つの実施形態によると、本明細書中に開示された方法のいずれも予測を用いて、癌にかかった特定の個体に利益を与えることができ、およびここに該予測は、その個体および該個体の特定の癌に関連する情報を供することによって臨床家を助けることができ、およびここで、該情報の主題は治療養生法、ライフスタイルの決定、およびダイエットの嗜好性、薬物介入、他の治療的介入およびその組合せを含むトピックスの群から取られ、およびここで、該予測は該個体の遺伝子メイキャップ、該個体の表現型特徴、該個体の臨床的履歴、およびその組合せに関するデータに基づく。

30

【0060】

1つの実施形態によると、本明細書中に開示された方法のいずれも、病原体に罹った特定の個体に利益を与えるために用いることができ、およびここで、該予測は、その個体、および該個体を感染する特定の病原体に関連する情報を供することによって臨床家を助けることができ、ここで、該病原体は細菌、ウイルス、微生物、アメーバ、真菌および他の寄生虫よりなる群から選択されるクラスのものであり、およびここで、該情報の主題は治療養生法、ライフスタイルの決定、およびダイエットの嗜好性、薬物介入、他の治療的介入、およびその組合せを含むトピックスの群から取られ、およびここで、該予測は該個体の遺伝子メイキャップ、該個体の表現型特徴、該個体の臨床的履歴、およびその組合せに関するデータに基づく。

40

【0061】

もう1つの実施形態によると、本明細書中に開示された方法のいずれも、具体的な個体についての予測、新しい知識、およびデータを用いることができる。というのは、その知識およびデータは入手可能となるからであり、これを用いて、該個体に関連するトピックスについての、情報報告を自動的にまたは要求に応じて作成することができ、ここで、ト

50

ピックスはライフスタイルの決定、ダイエットの嗜好性、ホルモンサプリメント、病気についての可能な治療養生法、病原体についての可能な治療養生法、薬物の介入、他の治療的介入、およびその組合せを含む群から取られ、およびここで、新しい知識およびデータは性質において医療的であり、およびここで、該予測は、該個体の遺伝子のベークアップ、該個体の表現型特徴、該個体の臨床的履歴、およびその組合せに関するデータに基づく。

【 0 0 6 2 】

もう1つの実施形態によると、本明細書中に開示された方法のいずれも、特定の胚からの遺伝子データを用いる予測を用いることができ、該予測を用いて、該胚のある表現型に対する予測された感受性に基づくIVFの関係で胚の選択を助けることができる。

10

【 0 0 6 3 】

1つの実施形態によると、本明細書中に開示された方法のいずれも、特定の胎児からの遺伝子データを用いる予測を用いることができ、該予測を用いて、余命、乾癬の確率、または数学的能力の特定のレベルの確立のような、潜在的子孫についての特別な表現型の結果を見積もることができる。

【 0 0 6 4 】

この開示の利点を仮定すれば、他の態様、特徴および実施形態は本明細書中に開示された方法およびシステムの1以上を実施することができるのは当業者によって認識されるであろう。

例えば、本願発明は以下の項目を提供する。

20

(項目 1)

(i) 関連個体からのいずれの染色体のいずれのセグメントが標的個体ゲノムで見出されるセグメントに対応するかに関する1以上の仮説のセットを創製し、

(i i) 該標的個体遺伝子データの測定、および該関連個体遺伝子データの測定を仮定して該仮説の各々の確率を決定し、次いで、

(i i i) 各仮説に関連する確率を用いて、該標的個体の現実の遺伝物質の最もありそうな状態を決定する：

ことを含む、該標的個体の遺伝子データの不完全な知識、および該標的に遺伝的に関連する1以上の個体の遺伝子データの知識に基づいて該標的個体の遺伝子データを決定する方法。

30

(項目 2)

前記方法が、前記標的の遺伝子データの測定、および親の遺伝子データを仮定した特定の測定の尤度の決定に基づいて、親染色体のいずれの領域が、標的個体に寄与した配偶子の形成に寄与した最大尤度を有するかを、決定することを含む、項目1記載の方法。

(項目 3)

親の少なくとも1つのハプロタイプが、親のジブロイド試料から測定された遺伝子データ、およびジブロイド試料から測定されたいずれの対立遺伝子がいずれのハプロタイプに属するかを決定するのに用いられる親からのハプロイド試料から測定された遺伝子データを用いることによって決定されている、項目1記載の方法。

(項目 4)

前記遺伝的に関連する個体からの遺伝子データが、ジブロイド母性試料、ハプロイド母性試料、ジブロイド父性試料およびハプロイド父性試料からの遺伝子データを含む群から選択される、項目1記載の方法。

40

(項目 5)

清浄化された胚遺伝子データにおける個々のSNP要求の各々について信頼性が計算される、項目1記載の方法。

(項目 6)

前記遺伝的に関連する個体からの遺伝子データが、ジブロイド母性細胞、ジブロイド父性細胞、ハプロイド父性細胞、母性祖父からのジブロイド細胞、および母性祖父からのハプロイド細胞からの遺伝子データを含む群から選択される、項目1記載の方法。

50

(項目7)

前記前記遺伝的に関連する個体からの遺伝子データが、ジプロイド母性細胞、ジプロイド父性細胞、および問題となる表現型のキャリアーであることが知られた関連個体からのジプロイド細胞からの遺伝子データを含む群から選択される、項目1記載の方法。

(項目8)

遺伝的に関連する個体が、父親、母親、息子、娘、兄弟、姉妹、祖父、祖母、叔父、叔母、甥、姪、孫息子、孫娘、従兄弟、クローン、前記標的に対する公知の遺伝的関係を持つ他の個体、およびその組合せよりなる群から選択される、項目1記載の方法。

(項目9)

前記標的個体が、成人ヒト、若年ヒト、ヒト胎児、ヒト胚、非ヒト成体、非ヒト若年体、非ヒト胎児、および非ヒト胚よりなる群から選択される、項目1記載の方法。

10

(項目10)

前記個体の遺伝子データの1以上が、ポリメラーゼ鎖反応(PCR)、リガンド媒介PCR、縮重オリゴヌクレオチドプライマーPCR、多重置換増幅、対立遺伝子-特異的増幅技術、およびその組合せよりなる群から選択されるツールおよび/または技術を用いて増幅される、項目1記載の方法。

(項目11)

前記個体の遺伝子データの1以上が、分子逆転プローブ(MIP)、ゲノタイピングマイクロアレイ、Taqman SNPゲノタイピングアッセイ、Illuminaゲノタイピングシステム、他のゲノタイピングアッセイ、蛍光イン-サイチュハイブリダイゼーション(FISH)、およびその組合せを含む群から選択されるツールおよび/または技術を用いて測定される、項目1記載の方法。

20

(項目12)

前記個体の遺伝子データの1以上が、該個体のバルクジプロイド組織、該個体から取られた1以上のジプロイド細胞、該個体から取られた1以上の胚盤胞、該個体の精液、該個体の卵、該個体で見出される細胞外遺伝物質、母性血液で見出される該個体からの細胞外遺伝物質、母性血漿で見出される該個体からの細胞外遺伝物質、母性血液で見出される該個体からの細胞、該個体に由来することが知られている遺伝物質、およびその組合せを含む群から選択される物質を分析することによって測定される、項目1記載の方法。

(項目13)

前記関連個体遺伝子データの1以上が、イン・シリコにて部分的にまたは全体的に知られているか、あるいは前記標的個体の遺伝子データを決定する以外の個人によって提供される、項目1記載の方法。

30

(項目14)

前記個体の1以上のハプロイド遺伝子データが、ジプロイドデータからハプロイドデータをシミュレートするコンピュータアルゴリズムによってイン・シリコにて部分的にまたは全体的に創製される、項目1記載の方法。

(項目15)

前記コンピュータアルゴリズムが隠れMarkovモデルを含む項目14記載の方法。

(項目16)

前記標的遺伝子データの決定が、体外受精の関係で胚選択を目的として用いられる、項目1記載の方法。

40

(項目17)

前記標的遺伝子データの決定が、出生前遺伝子診断の目的で用いられる、項目1記載の方法。

(項目18)

前記標的遺伝子データの決定が、統計学的モデルおよび/または専門家則を用いて表現型罹患性の予測を行う目的で用いられる、項目1記載の方法。

(項目19)

前記標的遺伝子データの決定が表現型予測を行う目的で用いられ、ここで、該表現型のい

50

くつかまたは全てを提示する尤度は、他の従前に知られた表現型情報によって影響される、項目 1 記載の方法。

(項目 20)

前記標的遺伝子データの決定が表現型予測を行う目的で用いられ、ここで、該予測は、該標的遺伝子データを、パブリックドメインで見出される公知の遺伝子マーカーと比較することによってなされる、項目 1 記載の方法。

(項目 21)

標的遺伝子データの決定が、臨床的決定を行う目的で用いられる、項目 1 記載の方法。

(項目 22)

標的遺伝子データの決定が、臨床的決定を行う目的で表現型マーカーと組合せて用いられる、項目 1 記載の方法。

10

(項目 23)

前記標的遺伝子データの決定が、1以上の病気に対する罹患性についてスクリーニングする目的で用いられ、ここで、家族の病歴が存在しない、項目 1 記載の方法。

(項目 24)

前記標的遺伝子データの決定が、1以上の表現型に対する罹患性についてスクリーニングする目的で用いられ、ここで、該表現型のいくつかまたは全てが多重遺伝子的である、項目 1 記載の方法。

(項目 25)

前記標的遺伝子データの知識が、汚染 DNA または RNA からの偽データを含むことが知られた、または含むことが疑われる、項目 1 記載の方法。

20

(項目 26)

前記個体の 1 以上の遺伝子データが、複数の SNP についての対立遺伝子要求、および各 SNP が知られている信頼性を含む、項目 1 記載の方法。

(項目 27)

前記標的個体の SNP 要求における信頼性が、該 SNP が正しく v s 正しくなく要求される確率のオッズ比を計算することによって決定される、項目 1 記載の方法。

(項目 28)

項目 1 記載の方法を達成するように構成されたシステム。

(項目 29)

項目 1 記載の方法を達成するように構成されたコンピュータ実施システム。

30

(項目 30)

(i) 標的個体のゲノムに存在する所与のセグメントの存在の数についての 1 以上の仮説のセットを創製し、

(ii) 該所与のセグメント上の複数の遺伝子座における可能な対立遺伝子のいくつかまたは全てについての遺伝子データの量を測定し、

(iii) 該標的個体の遺伝子データおよび、恐らくはまた、関連個体の遺伝子データの測定を仮定して該仮説の各々の相対的確率を決定し、次いで、

(iv) 各仮説に関連する相対的確率を用いて、該標的個体の現実の遺伝物質の最もありそうな状態を決定する；

40

ことを含む、該標的個体の所与の染色体の所与のセグメント上の多数遺伝子座の測定を用いて、該標的個体のゲノム中の所与のセグメントの存在の数を決定する方法。

(項目 31)

該標的ゲノムに存在する染色体のセグメントの存在の数の決定が、染色体異常についてスクリーニングする関係で行われ、この異常は、モノソミー、片親ジソミー、トリソミー、他の異数性、アンバランスなトランスロケーション、およびその組合せを含むリストから選択される、項目 30 記載の方法。

(項目 32)

各仮説の相対的確率の決定が、マッチドフィルタリングの概念を用いて行われる、項目 30 記載の方法。

50

(項目33)

各仮説の相対的確率の測定が、対立遺伝子要求を行わない定量的技術を用いてなされ、ここで、各遺伝子座の測定についての平均および標準偏差が既知、未知、または均一のいずれかである、項目30記載の方法。

(項目34)

各仮説の相対的確率の決定が、対立遺伝子要求を用いる定性的技術を用いてなされる、項目30記載の方法。

(項目35)

各仮説の相対的確率の決定が、参照配列の公知の対立遺伝子、および定量的対立遺伝子測定を用いることによってなされる、項目30記載の方法。

10

(項目36)

前記標的個体が、成人ヒト、若年ヒト、ヒト胎児、ヒト胚、非ヒト成体、非ヒト若年体、非ヒト胎児、および非ヒト胚よりなる群から選択される、項目30記載の方法。

(項目37)

前記標的個体の遺伝子データが、ポリメラーゼ鎖反応(PCR)、リガーゼ媒介PCR、縮重オリゴヌクレオチドプライマーPCR、多重置換増幅、対立遺伝子-特異的増幅およびその組合せを含む群から取られるツールおよび/または技術を用いて増幅される、項目30記載の方法。

(項目38)

前記標的個体の遺伝子データが、分子逆転プローブ(MIP)、ゲノタイピングマイクロアレイ、Taqman SNPゲノタイピングアッセイ、Illuminaゲノタイピングシステム、他のゲノタイピングアッセイ、蛍光イン-サイチュハイブリダイゼーション(FISH)、およびその組合せを含む群から選択されるツールおよび/または技術を用いて測定される、項目30記載の方法。

20

(項目39)

前記標的個体の遺伝子データが、該標的個体のバルクジプロイド組織、該標的個体から取られる1以上のジプロイド細胞、該標的個体から取られる1以上の胚盤胞、該標的個体上で見出された細胞外遺伝物質、母性血液で見出された該標的個体からの細胞外遺伝物質、母性血液で見出される該標的個体からの細胞、該標的個体に由来することが知られた遺伝物質、およびその組合せを含む群から取られる物質を分析することによって測定される、項目30記載の方法。

30

(項目40)

前記標的における染色体または染色体セグメントの数の決定が、体外受精の関係で胚選択を目的として用いられる、項目30記載の方法。

(項目41)

前記標的の染色体または染色体セグメントの数の決定が、出生前遺伝子診断の目的で用いられる、項目30記載の方法。

(項目42)

項目30記載の方法を達成するように構成されたシステム。

(項目43)

項目30記載の方法を達成するように構成されたコンピュータ実施システム。

40

(項目44)

(i) 関連個体からのいずれの染色体のいずれのセグメントが標的個体のゲノムで見出されるセグメントに対応するかについての1以上の仮説のセットを創製し、

(ii) 該標的のゲノムに存在する所与の染色体セグメントの数についての1以上の仮説のセットを創製し、

(iii) 該所与のセグメント上の複数の遺伝子座における可能な対立遺伝子の各々についてゲノムデータの量を測定し、

(iv) 該標的個体の遺伝子データの測定、および該関連個体の遺伝子データの測定を仮定して仮説の各々の相対的確率を決定し、次いで、

50

(v) 各仮説に関連する相対的確率を用いて、該標的個体の現実の遺伝物質の最もありそうな状態を決定する；

ことを含む、該標的個体の遺伝子データの不完全な知識、および該標的に遺伝的に関連する1以上の個体の遺伝子データの知識に基づいて、該標的個体の遺伝子データ、および該標的ゲノムに存在する染色体、または染色体のセグメントの存在の数を決定する方法。

(項目45)

(i) 遺伝子-病気関連についての公に入手可能な情報から形成された偶発事象表に基づいてモデルを構築し；次いで、

(ii) 該モデルを適用して、個体に関連するデータに対して操作することによって予測を行う；

ことを含む、該個体に関連する予測を行う方法。

(項目46)

多数の独立変数を使用する前記偶発事象表の精度が、結果データを用いて洗練することができ、ここで、独立変数のサブセットのみが測定される、項目45記載の方法。

(項目47)

多数の独立変数を使用する前記偶発事象表の精度が、前記独立変数の関連についてのデータを用いて洗練することができる、項目45記載の方法。

(項目48)

多数の独立変数を使用する前記偶発事象表の精度が、前記独立変数のある値の出現の頻度についてのデータを用いて洗練することができる、項目45記載の方法。

(項目49)

(i) 予測すべき結果が知られている個体の第二のセットからの集合データを用いて複数のモデルを創製し、それをテストし；

(ii) 第一の個体で利用可能なデータを仮定した予測を行うための種々のモデルの相対的精度を計算し；次いで、

(iii) 最も正確なものとして同定されるモデルを用いて、該第一の個体について予測を行う；

ことを含む、第一の個体に関する予測を行う方法。

(項目50)

前記個体に関連するデータのタイプは、該個体の遺伝子型データ、該個体の表現型データ、該個体の臨床データ、および該個体の実験室データよりなる群から選択されるデータを含む、項目45記載の方法。

(項目51)

前記個体に関連するデータのタイプが、該個体の遺伝子型データ、該個体の表現型データ、および該個体の臨床データ、ならびに該個体の実験室データよりなる群から選択されるデータを含む、項目49記載の方法。

(項目52)

前記データのタイプが、また、前記個体を感染させる病原体のデータよりなる、項目45記載の方法。

(項目53)

前記データのタイプが、また、前記個体を感染させる病原体のデータよりなる、項目49記載の方法。

(項目54)

前記予測が、前記個体の表現型、表現型罹患性、可能な臨床的結果、ライフスタイルの決定、身体の運動、ダイエットの嗜好性、ホルモンサプリメント、栄養サプリメント、病気のための治療、病原体のための処理、望まない疾患についての治療、医薬での治療、およびその組合せよりなる群から選択されるトピックに関する、項目45記載の方法。

(項目55)

前記予測が、前記個体の表現型、表現型罹患性、可能な臨床的結果、ライフスタイルの決定、身体の運動、精神的運動、ダイエット嗜好性、ホルモンサプリメント、栄養サブリメ

10

20

30

40

50

ント、病気についての治療、病原体についての処理、望ましくない疾患についての治療、
医薬での治療、およびその組合せよりなる群から選択されるトピックスに関する、項目 4
9 記載の方法。

(項目 5 6)

前記予測を用いて、前記個体のための、または該個体の代理人のための報告を作成する、
項目 4 5 記載の方法。

(項目 5 7)

前記予測を用いて、前記個体のための、または該個体の代理人のための報告を作成する、
項目 4 9 記載の方法。

(項目 5 8)

前記操作の行為が、新しいデータについて操作して、前記個体の予測を更新することを含
み、ここで、該データは新しい研究データ、または他の対象についての新しい集合データ
を含む群から選択される、項目 4 5 記載の方法。

(項目 5 9)

前記操作の行為が、新しいデータについて操作して、前記個体の予測を更新することを含
み、ここで、該データは新しい研究データまたは他の対象についての新しい集合データ
を含む群から選択される、項目 4 9 記載の方法。

(項目 6 0)

項目 4 5 記載の方法を達成するように構成されたシステム。

(項目 6 1)

項目 4 9 記載の方法を達成するように構成されたシステム。

【図面の簡単な説明】

【0065】

【図 1】 配偶子形成についての減数分裂における組換えの概念の説明図。

【図 2】 ヒト染色体 1 の 1 つの領域に沿っての組換えの可変速度の説明図。

【図 3】 異なる仮定に対する偽陰性および偽陽性の確率の決定。

【図 4】 混合された女性試料、全てのヘテロ遺伝子座からの結果。

【図 5】 混合された男性試料、全てのヘテロ遺伝子座からの結果。

【図 6】 女性試料についての Ct 測定とは異なる男性試料についての Ct 測定。

【図 7】 混合された女性試料からの結果； Taqman 単一色素。

【図 8】 混合された男性試料からの結果； Taqman 単一色素。

【図 9】 混合された男性試料についての反復測定の分布。

【図 10】 混合された女性試料からの結果； qPCR 尺度。

【図 11】 混合された男性試料からの結果； qPCR 尺度。

【図 12】 女性試料についての Ct 測定とは異なる男性試料についての Ct 測定。

【図 13】 第三の似ていない染色体での異数性の検出。

【図 14】 定常対立遺伝子ドロップアウト速度での 2 つの増幅分布の説明図。

【図 15】 アルファのガウス確率密度関数のグラフ。

【図 16】 入力データ、データベースデータ、アルゴリズムおよび出力の一般的な関係の
ダイアグラム。

【図 17】 $P(H|M)$ をどのように駆動するかを視覚的概観。

【図 18】 シミュレートされたデータについての清浄化アルゴリズムの有効性を示すのに
用いられるアルゴリズムを記載するフローチャートの視覚的表示。

【図 19】 IVF の間における胚の表現型予測の関係での、本明細書中に開示された方法
を達成するように構成されたシステムの説明図。

【図 20】 疎な解を生じる LASSO 傾向の説明図。 Ridge 回帰解は 2 つの円の接合
に存在し、LASSO 解は円および四角形の接合に存在する。

【図 21】 訓練およびテストデータの 10 の異なる 9 : 1 スプリットにわたって平均し、
次いで、各々、7 つの PI または 10 の RTI にわたって平均した、測定したおよび予測
した応答の相関係数 (% で表した R) の表。

10

20

30

40

50

【図 2 2】P I 応答を予測するためのプロテアーゼ酵素における突然変異に関連する L A S S O モデルパラメーターの値のグラフ表示。最大の絶対的大きさを持つ 4 0 のパラメーターのみを示す。

【図 2 3】N R T I 薬物応答を予測するための R T 酵素における突然変異に関連する L A S S O モデルパラメーターの値のグラフ表示。最大の絶対的大きさを持つ 4 0 のパラメーターのみを示す。

【図 2 4】N N R T I 薬物応答を予測するための R T 酵素における突然変異に関連する L A S S O モデルパラメーターの値のグラフ表示。最大の絶対的大きさを持つ 4 0 のパラメーターのみを示す。

【図 2 5】記載なし。

10

【図 2 6】記載なし。

【図 2 7】記載なし。

【図 2 8】記載なし。

【 0 0 6 6 】

表 1 : O M I M / N C B I に見出される病気遺伝子のまとめ。

表 2 : 異なる異数性検出技術のまとめ。

表 3 : 低度な共分離を持つ S N P を用いて記載された方法についての入力データの例。

表 4 : 高度な共分離を持つ S N P を用いて記載された方法についての入力データの例。

表 5 : 表 2 に示された入力データに代えての出力データの例。

表 6 : 表 4 に示された入力データに代えての出力データの例。

20

表 7 : 予備的シミュレーションの結果。

表 8 : 方法の全シミュレーションの結果。

表 9 : アルツハイマー病の開始への影響における A P O E および A C E における突然変異の役割を理解するための F a r r e r (2 0 0 5)、L a b e r t (1 9 9 8)、および A l v a r e z (1 9 9 9) の結果を表す 3 つの分割表。

表 1 0 : 表 7 の実験のメタ - 分析から生じた結果。

表 1 1 : 訓練およびテストデータの 1 0 の異なる 9 : 1 スプリットにわたって平均した、種々の方法についてのプロテアーゼ阻害剤 (P I) 薬物に対する測定されたおよび予測された応答の相関係数 (% で表した R) の表。結果の標準偏差 (S t d . t e v .) は灰色で示す ; 測定された薬物応答の数は最後の列に示す。

30

法 1 2 : 訓練およびテストデータ 1 0 の異なる 9 : 1 スプリットにわたって平均された、種々の方法についての逆転写酵素阻害剤 (R T I) 薬物に対する測定されたおよび予測された応答の相関係数 (% で表した R) の表。結果の標準偏差 (S t d . d e v .) は灰色で示す ; 測定された薬物応答の数は最後の列に示す。

表 1 3 : プロテアーゼ阻害剤 (P I) 薬物応答についてのプレディクターとしての最小絶対選択および収縮オペレーター (L A S S O) によって選択された非ゼロ重みを持つ突然変異の数と共に、種々の回帰方法についての訓練で用いられる試料の数、および突然変異の合計数。

表 1 4 : 逆転写酵素阻害剤 (R T I) 応答についてのプレディクターとしての L A S S O によって選択された非ゼロ重みを持つ突然変異の数と共に、種々の方法での訓練で用いられる試料の数、および突然変異の合計数。

40

表 1 5 : イリノテカン実験についての表現型データ。

【発明を実施するための形態】

【 0 0 6 7 】

(好ましい実施形態の詳細な説明)

システムの概念的概観

開示されたシステムの 1 つの目標は、遺伝子診断の目的の高度に正確なゲノムデータを提供することである。個体の遺伝子データが有意な量のノイズ、またはエラーを含有する場合、開示されたシステムは、関連個体の遺伝子データ、およびその第二の遺伝子データに含まれる情報の間の同様性を用いて、標的ゲノムにおけるノイズを清浄化する。これは

50

、染色体のいずれのセグメントが配偶子形成に関与し、およびどこで減数分裂の間に交差が起こったか、従って、第二のゲノムのいずれのセグメントが標的ゲノムのセクションに対してほとんど同一であると予測されるかを決定することによってなされる。ある状況においては、この方法を用いてノイズな塩基対測定を清浄化することができるが、それを用いて、測定されなかったDNAの個々の塩基対または全領域の同一性を推定することもできる。加えて、なされた各再構成要求について信頼性を計算することができる。高度に単純化された説明を最初に示し、非現実的な仮定をなして、本発明の概念を説明する。今日の技術に適用することができる詳細な統計学的アプローチを以後示す。

【 0 0 6 8 】

システムのもう1つの目標は、染色体の異常な数、染色体のセクション、および染色体の起源を検出することにある。異数性であり、アンバランスなトランスロケーション、片親二染色体、または他の正味の染色体異常を有する一般的試料において、複数の遺伝子座に存在する遺伝物質の量を用いて、試料の染色体状態を決定することができる。この方法に対して多数のアプローチが存在し、それらのうちいくつかをここに記載する。いくつかのアプローチにおいて、試料に存在する遺伝物質の量は、異数性を直接的に検出するのに十分である。他のアプローチにおいて、遺伝物質を清浄化する方法を用いて、染色体不均衡の検出の効率を増強させることができる。なされた各染色体要求に対して信頼性を計算することができる。

【 0 0 6 9 】

該システムのもう1つの目標は、遺伝子データに関連する変数の効果をモデル化するように設計された項目の広いアレイを開発することによって、遺伝子データから最も単純かつ触知可能な統計学的モデルを抽出する有効かつ効果的手段を提供することにある。より具体的には、遺伝子データに基づいて表現型または表現型感受性をモデル化するための現在利用可能な方法のほとんどまたは全ては以下の欠点を有する：(i)それらは凸最適化技術を用いず、かくして、所与の訓練データセットに対するモデルパラメーターについての局所的な最小解を見出すことは保証されない；(ii)それらはモデルの複雑性を最小化する技術を用いず、かくして、それらは、独立した変数の数に対して少数の結果が存在する場合に十分に一般化されるモデルを形成しない；(iii)それらは、正規分布したデータの単純化仮定をなすことなく、論理的回帰の関係でデータからの最も単純な触知のルールを抽出を可能とせず；(iv)それらは遺伝子-遺伝子関連、遺伝子-表現型関連および遺伝子-病気関連についての先見的情報を活用して、表現型または表現型感受性の最良の可能な予測をしない；(v)それらは1を超えるモデルを提供せず、かくして、訓練データに対する種々のモデルの交差-確証に基づいて最良の可能なデータを選択するための一般的なアプローチを提供しない。これらの欠点は、遺伝子および表現型情報に関連する多量のデータクラスの分析に基づいて結果を予測する関係で臨界的である。まとめると、現在利用可能な方法は個体が遺伝子型が所与の特定の表現型特徴の尤度についての、または親の遺伝子型特徴を仮定した子孫における特定の表現型特徴の尤度についての質問に答えるのに効果的に力を与えない。

【 0 0 7 0 】

以下に掲げる説明のいくつかは、本書類の著者によって従前に公表された仕事を含むことに注意されたい。それは背景情報として提供されて、本明細書中に開示された材料の理解を容易とし、および該材料に対するより大きな関係を与える。

【 0 0 7 1 】

3つのカテゴリーにおいて遺伝子型-表現型予測モデルを考慮することができる：i) 遺伝子欠陥または対立遺伝子は100%の確実性をもって病気表現型を引き起こすことが知られている；ii) 病気表現型の確率を増加させる遺伝子欠陥および対立遺伝子、ここで、プレディクターの数は表現型確率を分割表でモデル化できるのに十分に小さい；およびiii) 多次元線形または非線形回帰モデルを用いて表現型を予測するのに用いることができる遺伝子マーカーの複雑な組合せ。オンラインメンデル遺伝データベース(Online Mendelian Inheritance Database (OMIM))

10

20

30

40

50

)における現在知られている配列および病気表現型を持つ359の遺伝子(表1、列2参照)のうち、大部分はカテゴリー(i)に入り;残りは圧倒的にカテゴリー(ii)に入る。しかしながら、経時的に、多数の遺伝子型-表現型モデルがカテゴリー(iii)において生起していると予測され、ここで、多数の対立遺伝子または突然変異の相互作用は、特定の表現型の確率を見積もるためにモデル化される必要があるであろう。例えば、シナリオ(iii)は、確実に、今日、HIVウイルスの遺伝子データに基づいて抗-レトロウイルス療法に対するHIVウイルスの応答を予測する関係で当てはまる。

【0072】

シナリオ(i)については、経験則に基づいて表現型の発生を予測するのは通常直接的である。1つの態様において、シナリオ(ii)について表現型の正確な予測をなすのに用いることができる統計的技術が記載されている。もう1つの態様において、シナリオ(iii)について正確な予測を行うのに用いることができる統計学的技術が記載されている。もう1つの態様において、特定の表現型、集合データの特定の組、および特定の個々のデータについて最良のモデルを選択することができる方法が記載されている。

【0073】

本明細書中に開示された方法のある実施形態は、分割表を実行して、シナリオ(ii)において正確に予測を行う。これらの技術は遺伝子-遺伝子関連および遺伝子-病気関連についての先見の情報を利用して、表現型または表現型感受性の予測を改良する。これらの技術は、関連した独立変数の全てがサンプリングされるのではない従前の実験からのデータを活用するのを可能とする。それらが失われたデータを有するという理由でこれらの従前の結果を捨てる代わりに、概技術はHapMapプロジェクトおよびその他からのデータを活用して、関連する独立変数のサブセットのみが測定された従前の実験を用いる。このように、全ての関連する独立した変数が測定された対象からのデータを単純に集合させるよりはむしろ、予測モデルを全ての集合データに基づいて訓練することができる。

【0074】

本明細書中に記載されたある方法は凸最適化を用いて、シナリオ(iii)において正確な予測をなすのに用いることができる疎なモデルを創製する。遺伝子型-表現型モデリングの問題はしばしば過剰決定系であるか、または不適切である。というのは、潜在的プレディクター-遺伝子、蛋白質、突然変異およびそれらの相互作用-の数は、測定された結果の数に対して大きいからである。そのようなデータのセットは、依然として、Occam's Razorと同様な原理を発見することによって正確に一般化される疎なパラメーターモデルを訓練するのに用いることができる。多くの可能な理論が観察を説明することができる場合、最も単純なのは最も正しいらしいものである。この哲学は、先に議論したシナリオ(iii)において遺伝子型-表現型モデルの形成に関連する1つの態様において具体化される。遺伝子データへの適用について本明細書中に記載された技術は、過少判断されたまたは誤って条件付けされた遺伝子型-表現型データセットについて疎なパラメーターモデルを創製することを含む。疎なパラメーターセットの選択はOccam's Razorと同様な原理を発揮し、結果として、潜在的プレディクターの数が測定された結果の数に対して大きい場合でさえ、正確なモデルが開発されるのを可能とする。加えて、シナリオ(iii)において遺伝子型-表現型モデルを形成するための本明細書中に記載された技術のある実施形態は、所与の訓練データセットについてのモデルパラメーターに対する全体的最小解を見出すことが保証された凸最適化技術を用いる。

【0075】

集合データのセット、および個体についての入手可能なデータのセットを仮定すれば、その個体についての最良な表現型予測を行うために、いずれの予測アプローチが最も適当であるかは稀にしか明瞭でない。正確な表現型予測を行う傾向があるモデルのセットを記載することに加えて、本明細書中に開示された実施形態は、多数の方法をテストし、所与の表現型予測についての最適方法、集合データの所与のセット、および予測がなされるべき個体についての入手可能なデータの所与の組を選択するシステムを代表する。開示された方法およびシステムは、多重モデルおよび多重訓練パラメーターを用いるデータの所与

10

20

30

40

50

のセットにおける全ての異なる独立した変数および従属する変数の組合せを調べ、次いで、独立した変数、従属した変数、およびテストデータで測定された最良のモデリング精度を達成するチューニングパラメーターの組を選択する。シナリオ (i) に対応する場合には、専門家則を立案することができ；カテゴリー (i i) におけるような少数の独立した変数での他の場合には、分割表は最良の表現型予測を提供し；およびシナリオ (i i i) のような他の場合には、線形または非線形回帰技術を用いて、予測の最適な方法を提供することができる。本開示を読んだ後には、個体について予測をなすための最良のモデルを選択するアプローチをどのようにして用いて、本明細書中に開示されたものを超えて多くのモデリング技術から選択することができるかは当業者に明瞭であろうことを注記する。

【 0 0 7 6 】

10

技術のある実施形態はいくつかの関係で示されている。まず、それは、分割表、および遺伝子マーカーに基づいて、アルツハイマー病の予測に焦点を当てる多くの臨床的実験から一体化されたデータの不完全な組を用いてアルツハイマー病を発生する尤度を予測する関係で示されている。次に、該システムは、回帰分析、およびウイルスゲノムにおける遺伝子マーカーの知識を用いて1型ヒト免疫不全ウイルス (H I V - 1) の薬物応答をモデル化する関係で示されている。最後に、該システムは、各々、回帰分析、および個体についての双方の遺伝子マーカーの不完全なデータ、および癌に関連する実験質的および臨床的対象情報を用いる、乳癌および結腸癌の種々の症例の治療におけるタモキシフェンおよびイリノテカンの用法によって引き起こされる副作用の予測の点で示されている。

【 0 0 7 7 】

20

遺伝子型テストの減少する費用のため、信頼性よくウイルス薬物応答、癌薬物応答、および他の表現型応答または遺伝子データからの結果を予測する統計学的モデルは、それらが病気治療、ライフスタイルまたは嗜好性決定、または他の活動であるか否かを問わず適当な作用のコースの選択において重要なツールである。記載された最適化技術は、臨床的決定を増強させる目的で多くの遺伝子型 - 表現型モデリングの問題に応用を有するであろう。

【 0 0 7 8 】

システムの技術的記載

データの清浄化：単純化された例

図1は、親における配偶子の形成について減数分裂の間に起こる組換えのプロセスを説明する。個体の母親からの染色体101はオレンジ色（または灰色）で示す。個体の父親からの染色体102は白色で示す。減数分裂の前相Iの間の複糸期として知られたこの間隔の間に、4つの染色分体103のテトラドが目に見える。相同対の非姉妹染色分体の間の交差は組換え小節104として知られた地点で起こる。説明の目的で該例は単一の染色体、および3つの遺伝子の対立遺伝子の特徴付けると推定される3つの単一ヌクレオチド多形 (S N P) に焦点を当てる。この議論では、S N P は母性および父性染色体上で別々に測定できると仮定する。この概念は多くのS N P、多数のS N Pによって特徴付けられる多くの対立遺伝子、多くの染色体、および母性および父性染色体をゲノタイプング前には個々に単離することができない現行のゲノタイプング技術に適用することができる。

30

40

【 0 0 7 9 】

注目するS N Pの間における潜在的交差の地点に注意を払わなければならない。3つの母性遺伝子の対立遺伝子のセットは、S N P (S N P ₁ , S N P ₂ , S N P ₃) に対応する (a _{m1} , a _{m2} , a _{m3}) として記載することができる。3つの父性遺伝子の対立遺伝子のセットは (a _{p1} , a _{p2} , a _{p3}) として記載することができる。図1において形成された組換え小節をコードし、組換え染色分体の各対についてちょうど1つの組換えがあると仮定する。このプロセスで形成された配偶子のセットは遺伝子対立遺伝子： (a _{m1} , a _{m2} , a _{p3})、(a _{m1} , a _{p2} , a _{p3})、(a _{p1} , a _{m2} , a _{p3})、(a _{p1} , a _{p2} , a _{m3}) を有するであろう。染色分体の交差がない場合において、配偶子は対立遺伝子 (a _{m1} , a _{m2} , a _{m3})、(a _{p1} , a _{p2} , a _{p3}) を有するで

50

あろう。関連領域において交差の2つの地点がある場合において、配偶子是对立遺伝子 (a_{m1}, a_{p2}, a_{m3})、(a_{p1}, a_{m2}, a_{p3})を有するであろう。対立遺伝子のこれらの8つの異なる組合せを、その特定の親について、対立遺伝子の仮説セットという。

【0080】

胚DNAからの対立遺伝子の測定はノイズであろう。この議論の目的では、胚DNAからの単一染色体を取り、それが、その減数分裂を図1で説明する親に由来すると仮定する。この染色体上の対立遺伝子の測定は、もし胚染色体における測定された対立遺伝子が a_{m1} であれば $A_1 = 1$ であり、もし胚染色体における測定された対立遺伝子が a_{p1} であれば $A_1 = -1$ であって、もし測定された対立遺伝子 a_{m1} または a_{p1} でなければ $A_1 = 0$ であるインジケータ変数のベクトルの項: $A = [A_1 A_2 A_3]^T$ で記載することができる。推定親についての対立遺伝子の仮説セットに基づき、前記したすべての可能な配偶子に対応する8つのベクトルのセットを作り出すことができる。前記した対立遺伝子については、これらのベクトルは $a_1 = [1 \ 1 \ 1]^T$ 、 $a_2 = [1 \ 1 \ -1]^T$ 、 $a_3 = [1 \ -1 \ 1]^T$ 、 $a_4 = [1 \ -1 \ -1]^T$ 、 $a_5 = [-1 \ 1 \ 1]^T$ 、 $a_6 = [-1 \ 1 \ -1]^T$ 、 $a_7 = [-1 \ -1 \ 1]^T$ 、 $a_8 = [-1 \ -1 \ -1]^T$ となろう。システムのこの高度に単純化された適用において、胚のありそうな対立遺伝子は、仮説セットおよび測定されたベクトルの間の単純な相関分析を行うことによって決定することができる:

$$i^* = \arg \max_i A^T a_i, \quad i = 1 \dots 8 \quad (1)$$

一旦 i^* が見出されれば、仮説

【0081】

【数1】

a_{i^*}

が胚DNAにおける対立遺伝子の最もありそうなセットとして選択される。次いで、2つの異なる仮定、すなわち、胚染色体は母親または父親に由来するという仮定を立て、このプロセスを2回反復する。最大の相関

【0082】

【数2】

$A^T a_{i^*}$

を生じるその過程は正しいと仮定されるであろう。各場合において、母親または父親の各DNAの測定に基づき、対立遺伝子の仮説セットを用いる。開示された方法の典型的な実施形態においては、特定の病気表現型とのその関連のため重要であるSNPの間の多数のSNPを測定し - これらは表現型 - 関連SNPまたはPSNPといわれるであろうことに注意されたし。PSNPの間の非表現型 - 関連SNP (NSNP) は、個体間で実質的に異なる傾向があるRefSNPをNCBI dbSNPデータベースから選択することによって、(例えば、特殊化されたゲノタイピングアレイを開発するための) 先見的に選択することができる。別法として、PSNPの間のNSNPは親の特定の対について選択することができる。なぜならばそれらは親の間で異なるからである。PSNPの間のさらなるSNPの使用は、交差がPSNPの間で起こるか否かをより高いレベルの信頼性をもって決定することを可能とする。異なる「対立遺伝子」をこの注記において言及するが、これは単に便宜的なものであり; SNPは蛋白質をコードする遺伝子には関連しないであろうことに注意するのは重要である。

【0083】

現行の技術との関連でのシステム

もう1つのより複雑な実施形態において、特定の交差の確率を考慮して、対立遺伝子の事後確率を特定の測定を仮定して計算する。加えて、マイクロアレイに典型的なシナリオおよび他のゲノタイピング技術をアドレスし、ここで、ある時点で単一の染色体についてよりはむしろ染色体の対についてSNPを測定する。胚、父性および母性染色体について

の遺伝子座 i における遺伝子型の測定は、各々、SNP測定の対を表すランダム変数 ($e_{1,i}, e_{2,i}$)、($p_{1,i}, p_{2,i}$) および ($m_{1,i}, m_{2,i}$) によって特徴付けることができる。もしすべての測定が対としてなされるならば、母性および父性染色体における交差の存在を決定することができないので、該方法は修飾される：受精胚および父性および母性ジプロイド組織を遺伝子型分けするに加えて、各親からの1つのハプロイド細胞、すなわち、精子細胞および卵細胞も遺伝子型分けする。精子細胞の測定された対立遺伝子は $p_{1,i}, i = 1 \dots N$ によって表され、父性ジプロイド組織から測定された相補的対立遺伝子は $p_{2,i}$ によって表される。動揺に、卵細胞の測定された対立遺伝子は $m_{1,i}$ によって表され、母親のジプロイド細胞におけるそれらの相補体は $m_{2,i}$ によって表される。これらの測定は、どこで親染色体が測定された精子および卵細胞を生じるかにおいて交差したかについての情報を提供しない。しかしながら、卵または精子上の N 個の対立遺伝子の配列は少数の交差によって、または交差なしによって、親染色体から作り出されたと仮定することができる。これは開示されたアルゴリズムを適用するための十分な情報である。あるエラーの確率は、父性および母性SNPの要求に関連する。このエラーの確率の見積もりは、なされた測定 ($p_{1,i}, p_{2,i}$) および ($m_{1,i}, m_{2,i}$)、および用いる技術についてのシグナル - 対 - ノイズ比率に基づいて変化するであろう。これらのエラーの確率は、開示された方法に影響することなく各遺伝子座についてユニークに計算することができるが、父性および母性SNPを正しく要求する確立は、各々、 p_p および p_m において一定であると仮定することによってここでは代数は単純化される。

10

20

【0084】

測定は、測定 M という胚DNAで行われると仮定する。加えて、 A が今やセットであって、ベクトルではないように、表記法をわずかに修飾する： A とは、各親に由来する対立遺伝子の組合せ（またはセット）についての特定の仮説をいう。双方の親からの対立遺伝子 A のすべての可能な実施形態のセットを S_A として示す。目標は、測定 M を与えて、最大の事後確率をもって、対立遺伝子の組合せ（またはその仮説） $A \in S_A$ を決定することである：

$$A^* = \arg \max_{A \in S_A} P(A | M), \quad A \in S_A \quad (2)$$

条件付き確率の法則を用い、 $P(A | M) = P(M | A) P(A) / P(M)$ である。 $P(M)$ はすべての異なる A について共通するので、最適化サーチを：

30

$$A^* = \arg \max_{A \in S_A} P(M | A) P(A), \quad A \in S_A \quad (3)$$

として書き換えることができる。

【0085】

今や、 $P(M | A)$ の計算を考える。単一の遺伝子座 i で開始し、胚上のこの遺伝子座は親SNP $p_{t,1,i}$ および $m_{t,1,i}$ に由来すると仮定し、ここで、下付文字 t は、正しくても正しくなくてもよい行われた測定 $p_{1,i}$ および $m_{1,i}$ とは反対に、これらの親SNPの真の値を示すのに用いられる。胚SNPの真の値は ($e_{t,1,i}, e_{t,2,i}$) として示される。もし仮説 A が真であれば、($e_{t,1,i}, e_{t,2,i}$) = ($p_{t,1,i}, m_{t,1,i}$) または ($m_{t,1,i}, p_{t,1,i}$) である。測定 ($e_{1,i}, e_{2,i}$) のいずれが、いずれの親に由来するかを区別できないので、双方の順番を考慮しなければならず、従って、仮説セット $A = [(p_{t,1,i}, m_{t,1,i}), (m_{t,1,i}, p_{t,1,i})]$ となる。特定の測定 M の確率は、親SNPの真の値または基礎となる状態、すなわち、($p_{t,1,i}, p_{t,2,i}$) および ($m_{t,1,i}, m_{t,2,i}$) に依存する。4つのSNP、 $p_{t,1,i}, p_{t,2,i}, m_{t,1,i}, m_{t,2,i}$ が存在し、かつこれらの各々は4つのヌクレオチド塩基 A, C, T, G の値を取ることができるので、 4^4 または 256 の可能な状態が存在する。 $p_{t,1,i}, p_{t,2,i}, m_{t,1,i}, m_{t,2,i}$ であると仮定される1つの状態 s_1 についてアルゴリズムを説明する。この説明から、すべての 256 の可能な状態、 $s_k, k = 1 \dots 256$ にどのようにして該方法を適用するかは明瞭であろう。胚SNP ($e_{1,i}, e_{2,i}$) の測定 M を行い、結果 $e_{1,i} = p_{1,i}, e_{2,i} = m_{1,i}$ が

40

50

得られると仮定する。その仮説 A および状態 s_1 が真実であるとしたこの測定についての事前確率を計算する：

【 0 0 8 6 】

【 数 3 】

$$\begin{aligned}
 &P(e_{1,i} = p_{1,i}, e_{2,i} = m_{1,i} | A, s_1) = \\
 &P(e_{1,i} = p_{1,i}, e_{2,i} = m_{1,i} | A, s_1)P(e_{1,i} = p_{1,i} | e_{t,1,i} = p_{t,1,i})P(e_{2,i} = m_{1,i} | e_{t,2,i} = m_{t,1,i}) \\
 &+ P(e_{1,i} = m_{1,i}, e_{2,i} = p_{1,i} | A, s_1)P(e_{1,i} = p_{1,i} | e_{t,1,i} = m_{t,1,i}, p_{t,1,i} \neq m_{t,1,i})P(e_{2,i} = m_{1,i} | e_{t,2,i} = p_{t,2,i}, p_{t,2,i} \neq m_{t,1,i})
 \end{aligned}
 \tag{4}$$

第一項および第二項における最初の表現： $P(e_{1,i} = p_{1,i}, e_{2,i} = m_{1,i} | A, s_1) = P(e_{1,i} = m_{1,i}, e_{2,i} = p_{1,i} | A, s_1) = 0.5$ を考える。というのは、仮説 $A = [(p_{t,1,i}, m_{t,1,i}), (m_{t,1,i}, p_{t,1,i})]$ は胚 SNP についての 2 つの順序付けを等しくありそうとするからである。さて、第一項の第二の表現 $P(e_{1,i} = p_{1,i} | e_{t,1,i} = p_{t,1,i})$ を考え、これは、胚 SNP $e_{t,1,i}$ は現実には父性 SNP $p_{t,1,i}$ に由来すると仮定して $e_{1,i} = p_{1,i}$ を測定する確率である。父性 SNP、母性 SNP および胚 SNP を正しく測定する確率は p_p, p_m および p_e である。仮定 $(e_{t,1,i} = p_{t,1,i})$ を与えれば、測定 $(e_{1,i} = p_{1,i})$ は、胚および父性 SNP の双方が正しく測定されるか、あるいは双方は正しくなく測定され、それらは偶然に同一ヌクレオチド (A, C, T, または G) として正しくなく測定される、のいずれかを要求する。従って、 $P(e_{1,i} = p_{1,i} | e_{t,1,i} = p_{t,1,i}) = p_e p_p + (1 - p_e)(1 - p_p) / 3$ であり、ここで、単純性のために、4 つのヌクレオチドのすべてを正しくなく要求する確率は同等にありそうであると仮定される - 該アルゴリズムは、もう 1 つの特定のヌクレオチドについての測定を与えて特定のヌクレオチド (A, C, T, G) を要求する異なる確率を適合させるように容易に修飾することができる。同一アプローチを第一項中の 3 番目の表現に適用して、 $P(e_{2,i} = m_{1,i} | e_{t,2,i} = m_{t,1,i}) = p_e p_m + (1 - p_e)(1 - p_m) / 3$ を得ることができる。さて、第二項の 2 番目の表現を考える。 $P(e_{1,i} = p_{1,i} | e_{t,1,i} = m_{t,1,i}, m_{t,1,i} \neq p_{t,1,i})$ は、 $e_{1,i}$ または $p_{1,i}$ が正しくない測定であるか、または双方が正しくない測定であるかのいずれかを要求し、従って、測定された値は偶然に等しい： $P(e_{1,i} = p_{1,i} | e_{t,1,i} = m_{t,1,i}, m_{t,1,i} \neq p_{t,1,i}) = p_e(1 - p_p) / 3 + (1 - p_e)p_p / 3 + (1 - p_e)(1 - p_p)2 / 9$ 。同一の議論を第二項の最後の表現に適用して、 $P(e_{2,i} = m_{1,i} | e_{t,2,i} = p_{t,2,i}, m_{t,1,i} \neq p_{t,2,i}) = p_e(1 - p_m) / 3 + (1 - p_e)p_m / 3 + (1 - p_e)(1 - p_m)2 / 9$ を得ることができる。さて、これらの項のすべてを組合せ、単に代数を単純化するために、 $p_e = p_p = p_m = p$ と仮定して、

【 0 0 8 7 】

【 数 4 】

$$P(M(e_{1,i} = p_{1,i}, e_{2,i} = m_{1,i}) | A, s_1) = \frac{1}{162}(160p^4 - 160p^3 + 96p^2 - 28p + 13)
 \tag{5}$$

を計算することができる。計算は変化するが、本明細書中に記載されたものに対して同様な概念的アプローチをすべての 256 の可能な状態、 $s_k, k = 1 \dots 256$ で用いる。すべての 256 の状態 s_i について $P(e_{1,i} = p_{1,i}, e_{2,i} = m_{1,i} | A, s_i)$ を計算し、各 s_i の確率を合計し、 $P(e_{1,i} = p_{1,i}, e_{2,i} = m_{1,i} | A)$ を得る。言い換えれば：

【 0 0 8 8 】

【数5】

$$P(M | A) = \sum_{i=1..256} P(M | A, s_i) P(s_i)$$

(6)

である。各状態 s_i の確率 $P(s_i)$ を計算するために、別々の事象としての状態をなすすべての別々の対立遺伝子処理しなければならない。というのは、それらは別々の染色体上にあるからである、言い換えれば： $P(s_i) = P(p_{t,1,i}, p_{t,2,i}, m_{t,1,i}, m_{t,2,i}) = P(p_{t,1,i}) P(p_{t,2,i}) P(m_{t,1,i}) P(m_{t,2,i})$ である。ベイズ技術を適用して、個々の測定についての確率分布を見積もることができる。遺伝子座 i における母性または父性染色体上の対立遺伝子の各測定をコイン投げ実験として処理して、特定の値 (A, C, T , または G) であるこの対立遺伝子の確率を測定することができる。これらの測定を成人組織試料でなし、全く信頼性があるとして処理することができるが、対立遺伝子の対は各 SNP について測定し、いずれの対立遺伝子がいずれの染色体に由来かを決定するのは可能でない。 $w_{p,1,i} = P(p_{t,1,i})$ とし、これは、父親の染色体上の SNP i の確率が値 $p_{t,1,i}$ であることに対応する。以下の説明において、 $w_{p,1,i}$ の代わりに w を用いる。父親の染色体の SNP i で行った測定は収集データとして特徴付けられるものとする。 w についての確率分布 $p(w)$ を作り出し、データがベイズ理論： $p(w | D) = p(w) p(D | w) / p(D)$ に従って測定した後これを更新することができる。SNP i の n 個の対立遺伝子が観察され、 w に対応する特定の対立遺伝子が h 回出現する、言い換えれば、ヘッドは h 回観察されると仮定する。この観察の確率は二項分布によって特徴づけることができる。

10

20

【0089】

【数6】

$$p(D | w) = \binom{n}{h} w^h (1-w)^{n-h}$$

(7)

30

データを収集する前に、0 および 1 の間では均一である事前分布 $p(w)$ があると仮定する。ベイズ理論を適用することによって、直接的に、 $p(w | D)$ についての得られた分布は形式：

【0090】

【数7】

$$p(w | D) = \frac{1}{c} w^h (1-w)^{n-h} \quad \text{ここで} \quad c = \int_0^1 w^h (1-w)^{n-h} dw$$

(8)

のデータ分布であることを示し、 c は正規化定数である。しかしながら、次いで、ベイズ理論および新しい測定を適用することによって、 $p(w | D)$ を何回も更新し、それを、前記したデータ分布を有するように継続する。 $p(w)$ の見積もりは、新しい測定が収集されるごとに更新される。特定の SNP における異なる対立遺伝子の確立は人種および性別のグループ分けに依存するので、HapMap プロジェクトで用いたのと同じのグループ分けを用いて、異なる人種および異なる性別について異なる関数 $p(w)$ があることに注意されたし。 $P(s_i)$ の計算では、各染色体上の各対立遺伝子は見積もられた確率分布、すなわち、 $p_{p,1,i}(w_{p,1,i})$ 、 $p_{p,2,i}(w_{p,2,i})$ 、 $p_{m,1,i}(w_{m,1,i})$ および $p_{m,2,i}(w_{m,2,i})$ と関連するであろう。次いで、個々の分布の各々についての MAP 見積もりに従って $P(s_i)$ についての最大事後 (MAP) 見積もりを計算することができる。例えば、 $w_{p,1,i}^*$ は、 $p_{p,1,i}$

40

50

($w_{p,1,i}$) を最大化する議論であるとする。 $P(s_i)$ の M A P 見積もりは：
 $P(s_i)_{M A P} = w_{p,1,i} * w_{p,2,i} * w_{m,1,i} * w_{m,2,i}$
 (9)

に従って見出すことができる。各 w について確率分布が存在するので、M A P 見積もりを単に用いるよりはむしろ、確率分布に渡って積分することによって、いずれかの特定の信頼性レベルまで値 $P(s_i)$ の保存的見積もりを計算することもできる。例えば、これを行って、ある信頼性レベル内まで保存的に $P(M|A)$ を見積もることが可能である。保存的見積もりまたは M A P 見積もりを用いるかに拘わらず、 $P(s_i)$ の見積もりは $P(M|A)$ の計算のために継続的に洗練される。以下において、仮定された状態への言及をなくして、表記法を単純化し、状態 s_1 は詳細な計算のすべての説明のために仮定される。現実には、これらの計算は 256 の状態の各々について行われ、各々の確率に渡って合計することを銘記されたし。

10

【0091】

$P(M|A)$ を計算する方法は、今や、 M が胚上の S N P の N 個の対の測定のセット、 $M = [M_1, \dots, M_N]$ を表すと仮定し、多数の S N P 遺伝子座まで拡大される。また、 A は、いずれの親染色体がその S N P に貢献したかについての各 S N P に対する仮説のセットを表すと仮定する、 $A = [A_1, \dots, A_N]$ 。 $S A'$ が、 A とは異なる、またはセット A' に存在するというすべての他の可能な仮説のセットを表すものとする。 $P(M|A)$ および $P(M|A')$ を計算することができる：

【0092】

【数8】

$$P(M|A) = \prod_{i=1 \dots N} P(M_i | A_i), \quad P(M|A') = \sum_{A \in S_A'} P(A) \prod_{i=1 \dots N} P(M_i | A_i) \quad (10)$$

20

$P(A)$ の計算を考える。本質的には、これは、胚を形成する配偶子の形成において起こる特定の交差の尤度に基づく。特定の対立遺伝子セットの確率は2つの因子、すなわち、胚染色体が母親または父親に由来する確率、および交差の特定の組合せの確率に依存する。異数性をこうむらない胚染色体の清浄なセットについては、胚染色体が母親または父親に由来する事前確率は ~ 50% であり、その結果、すべての A について共通する。さて、組換え節の特定のセットの確率を考える。関連組換え部位 R の数は測定された S N P S : $R = N - 1$ の数に依存する。注目する P S N P の回りの N 個の N S N P を構成する DNA セグメントは比較的短いので、交差干渉は、同一染色体上の2つの交差が1つの領域で起こり得ることをかなりありそうもなくする。計算の効率の理由で、この方法は、唯一の交差が各関連染色体についての各領域で起こると仮定し、これは R 個の可能な部位で起こり得る。どのようにしてこの方法を拡大して、所与の領域に多数の交差がある確率を含めることができるかは当業者に明らかであろう。

30

【0093】

S N P の間の各領域における交差を P_r , $r = 1 \dots N - 1$ で示すものとする。一次的には、2つの S N P の間の領域 r における組換え節の確率は、(c モルガンで測定された) それらの S N P の間の遺伝子距離に比例する。しかしながら、多数の最近の研究は、2つの S N P 遺伝子座の間の組換えの確率の正確なモデリングを可能とした。精子の実験からの観察、および遺伝子変異のパターンは、組換えの率はキロベーススケールに渡って広く変化し、および多数の組換えは組換えホットスポットで起こり、連鎖非平衡を引き起こして、ブロック - 様構造を呈することを示す。ヒトゲノム上での組換え率についての N C B I データは、U C S C Genome Annotation Database を通じて公に入手可能である。

40

【0094】

種々のデータセットを単独で、または組合せて用いることができる。最も普通のデータセットの内の2つは H a p m a p プロジェクトおよび P e r l e g e n ヒトハプロタイププロジェクトからのものである。後者はより高い密度であり；前者はより高い質である。

50

HapMap相Iデータ、リリース16aに基づく、染色体1の位置1,038,423ないし4,467,775からの領域的組換え率については図2参照。これらの率は、パッケージLDHatで入手可能な可逆的ジャンプMarkov Chain Monte Carlo (MCMC)方法を用いて見積もられた。考えられる状態-空間は、ピース様定常組換え率マップの分布である。Markov鎖は、各セグメント201についての率に加えて、率変更点の数および位置の分布を探索する。これらの結果を用いて、SNP Sの間の各定常セグメントの長さの組換え率倍に渡って積分することによってP_rの見積りを得ることができる。ヌクレオチド202に渡っての累積組換え率を赤色で図2に示す。

【0095】

10

もし領域rおよびそうでなければ0で交差が起こったならばc_r = 1であるように、Cをインジケータ変数c_rのセットとする。もし交差が起こらないか、そうでなければ0であれば、c₀ = 1である。ただ1つの交差がN個のSNPの領域で起こり得ると仮定するので、セットCのただ1つのエレメントは非0である。よって、セットCによって表される交差の確率は：

【0096】

【数9】

$$P_c = \left(1 - \sum_{r=1 \dots N-1} P_r\right)^{c_0} \prod_{r=1} P_r^{c_r} \tag{11}$$

20

であることが判明する。SNP 1 . . . Nについての仮説Aにおいて、関連する4つの潜在的交差がある。すなわち、i) (インジケータ変数のセットC_{pe}によって示される) 胚を形成した父性染色体、ii) 配列決定された精子を形成した父性染色体 (セットC_{ps})、iii) 胚を形成した母性染色体 (セットC_{me})、およびiv) 配列決定された卵を形成した母性染色体 (セットC_{ee})。2つのさらなる仮定はv) 第一の父性胚SNPがp_{t,1,1}またはp_{t,2,1}に由来するか、およびvi) 第一の母性胚SNPがm_{t,1,1}またはm_{t,2,1}に由来するかである。SNPの間の交差の確率は人種および性別の間で異なることが見出されるので、異なる交差確率は父性染色体についてはp_{p,r}として、および母性染色体についてはp_{m,r}として示されるであろう。従って、セットC_{pe}、C_{ps}、C_{me}、C_{ee}を包含する特定の仮説Aの確率は；

30

【0097】

【数10】

$$P(A) = \frac{1}{4} \left(1 - \sum_{r=1 \dots N-1} P_{p,r}\right)^{c_{pe0}} \prod_{r=1 \dots N-1} P_{p,r}^{c_{p,r}} \left(1 - \sum_{r=1 \dots N-1} P_{p,r}\right)^{c_{ps0}} \prod_{r=1 \dots N-1} P_{p,r}^{c_{p,r}} \left(1 - \sum_{r=1 \dots N-1} P_{m,r}\right)^{c_{me0}} \prod_{r=1 \dots N-1} P_{m,r}^{c_{m,r}} \left(1 - \sum_{r=1 \dots N-1} P_{m,r}\right)^{c_{ee0}} \prod_{r=1 \dots N-1} P_{m,r}^{c_{e,r}} \tag{12}$$

として表される。

【0098】

さて、P(A)およびP(M/A)を決定するための方程式に関しては、前記方程式3についてのA*を計算するのに必要な全ての要素は定義されている。よって、交差が起こった胚SNPの高度にエラー-傾向の測定から決定し、および高度な信頼性をもって胚測定を結果的に清浄化することが可能である。最良の仮説A*における信頼性の低度を決定することが残っている。これを決定するためには、オッズ比P(A* | M) / P(A* | M)を見出す必要がある。ツールは全てこの計算のために前記されている：

40

【0099】

【数11】

$$\frac{P(A^*|M)}{P(A^{*'}|M)} = \frac{P(A^*|M)}{1 - P(A^*|M)} = \frac{P(A^*)P(M|A^*)}{P(A^{*'})P(M|A^{*'})} = \frac{P(A^*)P(M|A^*)}{(1 - P(A^*))P(M|A^{*'})} = OR_{A^*} \tag{13}$$

次いで、A*における信頼性はP(A* | M) = OR_{A*} / (1 + OR_{A*})として与え

50

られる。この計算は特定の仮説 A^* における信頼性を示すが、SNPの特定の決定における信頼性を示さない。胚 $PSNP_n$ の決定における信頼性を計算するためには、このSNPの値を変化させない全ての仮説Aのセットを作り出す必要がある。このセットは $S_{A^*,n}$ として示され、これは、仮説 A^* によって予測されるように、同一の値を有する胚に $PSNP_n$ をもたらす全ての仮説に対応する。同様に、仮説 A^* によって予測される異なる値を有する $PSNP$ をもたらす全ての仮説に対応するセット $S_{A^*,n}$ を作り出す。さて、SNPが正しく要求される確率 - 対 - SNPが正しくなく要求される確率のオッズ比を計算することが可能である：

【 0 1 0 0 】
【 数 1 2 】

10

$$OR_{A^*,n} = \frac{\sum_{A \in S_{A^*,n}} P(A|M)}{\sum_{A \in S_{A^*,n}} P(A|M)} = \frac{\sum_{A \in S_{A^*,n}} P(A|M)}{1 - \sum_{A \in S_{A^*,n}} P(A|M)} = \frac{\sum_{A \in S_{A^*,n}} P(A)P(M|A)}{\sum_{A \in S_{A^*,n}} P(A)P(M|A)} \tag{14}$$

オッズ比 $OR_{A^*,n}$ に基づく胚SNPの特定の要求における信頼性は：

【 0 1 0 1 】
【 数 1 3 】

$$P(\text{正しく要求されたSNP}_n) = \sum_{A \in S_{A^*,n}} P(A|M) = \frac{OR_{A^*,n}}{1 + OR_{A^*,n}} \tag{15}$$

20

として計算することができる。

【 0 1 0 2 】

この技術を用いて、同一染色体の2つが同一の親からのものであり、他方、他の親からのその染色体のいずれも存在しない片親二染色体 (UPD) のような欠陥を検出することもできよう。親染色体における交差を推定しようと試みる際に、高い信頼性でもってデータを適切に説明する仮説はなく、もし複数のUPDを含む別の仮説が許容されるならば、それらはよりありそうであることが判明するであろう。

【 0 1 0 3 】

組換えラットにおける確実性の効果、およびSNP測定の信頼性のバウンディング

30

開示された方法は：特定のSNPの間の組換えの確立についての仮定；胚、精子、卵、父性および母性染色体についての各SNPの正しい測定の確率についての仮定；および異なる集団群内のある対立遺伝子の尤度についての仮定に依存する。これらの仮定の各々を考慮し：組換えのメカニズムは完全には理解され、モデル化されておらず、交差確率は、個人の遺伝子型に基づいて変化することが確立されている。さらに、組換え率が測定される技術は実質的可変性を示す。例えば、可逆的 - ジャンプ Markov Chain Monte Carlo (MCMC) 方法を実行するパッケージ LDAat は、仮定のセットを作成し、組換えのメカニズムおよび特徴付けについてのユーザーの入力のセットを必要とする。これらの仮定は、種々の実験によって得られた異なる結果によって証明されているように、SNPの間の予測された組換え率に影響し得る。

40

【 0 1 0 4 】

前記リストの全ての仮定のうち、組換え率についての仮定は方程式 15 に対して最もインパクトを有するであろうと予測される。前記した計算は、SNP_S、P_rの間の交差に対する確率の最良の見積もりに基づくべきである。その後、(正しくはSNP_nと呼ばれる) 信頼性尺度Pを低下させる方向において、例えば、組換え率についての95%信頼性範囲における値を用いてP_rで用いることができる。95%信頼性範囲は、組換え率の種々の実験によって生じた信頼性データに由来することができ、これは、異なる方法を用いて異なる群からの公表されたデータ間の不一致のレベルを見ることによって確認することができる。

【 0 1 0 5 】

50

同様に、95%信頼性範囲を、各SNPが正しく要求される確率の見積もりで用いることができる： p_p 、 p_m 、 p_e 。これらの数は、測定技術の信頼性についての経験的なデータと組み合わせた、ゲノタイピングアッセイ出力ファイルに含まれた現実の測定されたアレイ強度に基づいて計算することができる。これらのパラメータ p_p 、 p_m および p_e が確立されないNSNPは無視することができることを注記する。例えば、ジブroid親データは信頼性よく測定されるので、親のハクroid細胞、および親のジブroid組織の関連SNPについての対立遺伝子のいずれにも対応しない胚についてのNSNP測定を無視することができる。

【0106】

最後に、計算 $P(s_i)$ を生起する異なる集団内のある対立遺伝子の尤度についての仮定を考える。これらの仮定もまた開示された方法に対して大きなインパクトを有しないであろう。というのは、親ジブroidデータの測定は信頼性があり、すなわち、親試料からの状態 s_i の直接的測定は、典型的には、高い信頼性を持つデータをもたらすからである。それにも拘わらず、方程式8に記載された各 w についての確率分布を用いて、各状態 $P(s_i)$ の確率についての信頼性範囲を計算することが可能である。前記したように、(正しくはSNP $_n$ と呼ばれる)信頼性尺度 P を低下させる保存的方向における各 $P(s_i)$ についての95%信頼性範囲を計算することができる。

【0107】

(正しくはSNP $_n$ と呼ばれる) P の決定は、どのようにして多くのNSNPが各PSNPについて測定される必要があるかについての決定を知らせて、所望のレベルの信頼性を達成するであろう。

【0108】

開示された方法の概念を実施する、すなわち、親のDNAの測定、1以上の胚のDNAの測定、および減数分裂のプロセスの事前知識を組合せて、胚SNPの良好な見積もりを得る異なるアプローチがあることを注記する。事前知識の異なるサブセットが知られており、または知られておらず、または大きなまたは小さな低度の確実性でもって知られている場合に、どのようにして同様な方法を適用することができるかは当業者に明らかであろう。例えば、多数の胚の測定を用いて、特定の胚のSNPを要求する確実性を改良し、または親からの失われたデータを供給することができる。注目するPSNPを測定技術によって測定する必要がないことを注記する。たとえ測定システムによってPSNPが決定されなくても、それは、依然として、開示された方法によって高度な信頼性でもって再構築できる。

【0109】

一旦減数分裂の間に起こった交差の点が決定され、標的ゲノムの領域が親DNAの関連領域にマッピングされれば、注目する個体のSNPの同一性のみならず、測定における対立遺伝子ドロップアウトまたは他のエラーによる測定された標的ゲノムで失われているであろうDNAの全領域を推定することが可能である。または、親DNAにおける挿入および欠失を測定し、開示された方法を用いて、それらが標的DNAに存在すると推定することも可能である。

【0110】

種々の技術を用いて前記して開示アルゴリズムの計算の複雑性を改善することができる。例えば、母親および父親の間で異なるNSNPを選択することができるにすぎないか、または圧倒的に選択することができる。もう1つの考慮は、PSNPの近くに間隔が設けられたNSNPを用いて、注目するNSNPおよびPSNPの間で起こる交差のチャンスを最小化することであろう。また、多数のPSNPの適用範囲を最大化するために染色体に沿って間隔を設けたNSNPを用いることもできる。もう1つの考慮は、最初に少数のNSNPのみを用いて、大まかにどこで交差が起こったかを、限定された程度の確率のみでもって決定することであろう。次いで、さらなるNSNPを用いて、交差モデルを洗練し、正しくPSNPを要求する確率を増加させることができる。考慮する交差組合せの数は、 N がSNPの数であって、 C が最大数の交差である N^C として概略評価する。結果と

10

20

30

40

50

して、 $C = 4$ については、Pentium（登録商標）-IVプロセッサに対して計算可能に御しやすくしつつ、各PSNPについて概略 $N = 100$ を供給することが可能である。前記したアプローチ、および増大した計算効率についての他のアプローチを用い、 $N > 100$ 、 $C > 4$ を容易に供給することができる。1つのそのようなアプローチを以下に記載する。

【0111】

基本となる概念を変化させることなく、胚データ、親データ、および用いるアルゴリズムの特定のセットに基づいて、PSNPについての要求を行い、PSNPが正しく決定された確率の見積もりを生じる多くの他のアプローチがあることを注記する。この確率は個人の決定をなすのに、およびIVFまたはNIPGDの関係で信頼性のよいサービスを実行するのに用いることができる。

10

【0112】

遺伝子データ清浄化アルゴリズムに対する帰納的解

直線的に範囲を定めるアルゴリズムに関連する本発明のもう1つの実施形態をここに記載する。計算パワーの限定された性質を仮定すると、計算の長さは開示された方法の使用において重要な因子であり得る。計算を実行する場合、必要とされる計算の数がSNPの数と共に指数関数的に上昇するある値を計算しなければならないいずれのアルゴリズムも扱いにくくなり得る。SNPの数と共に直線的に増加する多数の計算を含む解は、常に、SNPの数が大きくなるにつれて時間の観点から好ましいであろう。以下に、このアプローチを記載する。

20

【0113】

全ての可能な仮説を考慮する単純なアプローチは、SNPの数が指数関数である実行時間と戦わなければならない。前記したように、 k 個のSNPについての測定された胚、父親および母親染色体の測定のコレクションであると仮定する。すなわち、 $M = \{M_1, \dots, M_k\}$ であり、ここで、 $M_i = (e_{1i}, e_{2i}, p_{1i}, p_{2i}, m_{1i}, m_{2i})$ である。前記したように、仮説空間は $S_H = \{H^1, \dots, H^q\} = \{\text{全ての仮説のセット}\}$ であり、ここで、各仮説はフォーマット $H^j = \{H^j_1, \dots, H^j_k\}$ のものであり、ここで、 H^j_i はフォーマット $H^j_i = (p_i^*, m_i^*)$ のスニップ i についての「ミニ」仮説であり、ここで、 $p_i^* = \{p_{1i}, p_{2i}\}$ および $m_i^* = \{m_{1i}, m_{2i}\}$ である。4つの異なる「ミニ」仮説 H^j_i 、特に：

30

$H^j_{i1} : (e_{1i}, e_{2i}) = \{(p_{1i}, m_{1i}) \text{ または } (m_{1i}, p_{1i})\}$

$H^j_{i2} : (e_{1i}, e_{2i}) = \{(p_{1i}, m_{2i}) \text{ または } (m_{2i}, p_{1i})\}$

$H^j_{i3} : (e_{1i}, e_{2i}) = \{(p_{2i}, m_{1i}) \text{ または } (m_{1i}, p_{2i})\}$

$H^j_{i4} : (e_{1i}, e_{2i}) = \{(p_{2i}, m_{2i}) \text{ または } (m_{2i}, p_{2i})\}$

がある。目標は、最もありそうな仮説 H^* を：

【0114】

【数14】

$$H^* = \arg \max_{H \in S_H} P(H|M) = \arg \max_{H \in S_H} F(M, H)$$

として選択することであり、ここで、関数 $F(M, H) = P(H|M)$ である。

40

【0115】

空間 S^H において 4^k の異なる仮説がある。全空間 S^H を専ら調べることによって最良の仮説を見出す試みによって、必要なアルゴリズムは $k \cdot O(\exp(k))$ における指数関数オーダーのものであり、ここで、 k は関連するSNPの数である。大きな k 、 $k > 5$ さえについても、これはかなり遅く、非現実的である。従って、一定時間内にサイズ $(k-1)$ の問題の関数としてサイズ k の問題を解く帰納的解に頼るのがより現実的である。本明細書中に示された解は $k \cdot O(k)$ における直線オーダーのものである。

【0116】

SNPの数において直線的な帰納的解

$F(M, H) = P(H|M) = P(M|H) \cdot P(H) / P(M)$ で始める。次いで、

50

$\text{argmax}_H F(M, H) = \text{argmax}_H P(M | H) * P(H)$ であり、目標は直線の時間内に $P(M | H) * P(H)$ を解くことである。 $M_{(s, k)} = \text{SNP}_{s \text{ ないし } k}$ での測定、 $H_{(s, k)} = \text{SNP}_{s \text{ ないし } k}$ についての仮説とし、表現方法 $M_{(k, k)} = M_k$ 、 $H_{(k, k)} = H_k = \text{SNP}_k$ についての測定および仮説を単純化する。先に示したように：

【 0 1 1 7 】
【 数 1 5 】

$$P(M_{(1,k)} | H_{(1,k)}) = \prod_{i=1}^k P(M_i | H_i) = P(M_k | H_k) * \prod_{i=1}^{k-1} P(M_i | H_i) = P(M_k | H_k) * P(M_{(1,k-1)} | H_{(1,k-1)})$$

10

である。また、

【 0 1 1 8 】
【 数 1 6 】

$$P(H_{(1,k)}) = 1/4 * \prod_{i=2}^k PF(H_{i-1}, H_i) = PF(H_{k-1}, H_k) * 1/4 * \prod_{i=2}^{k-1} PF(H_{i-1}, H_i) = PF(H_{k-1}, H_k) * P(H_{(1,k-1)})$$

であり、ここで、

【 0 1 1 9 】
【 数 1 7 】

$$PF(H_{i-1}, H_i) = \begin{cases} 1 - PC(H_{i-1}, H_i) & H_{i-1} = H_i \\ PC(H_{i-1}, H_i) & H_{i-1} \neq H_i \end{cases}$$

20

であり、 $PC(H_{i-1}, H_j) = H_{i-1}, H_i$ の間の交差の確率である。

【 0 1 2 0 】

最後に、 k 個の SNP については：

$$\begin{aligned} F(M, H) &= P(M | H) * P(H) = P(M_{(1, k)} | H_{(1, k)}) * P(H_{(1, k)}) \\ &= P(M_{(1, k-1)} | H_{(1, k-1)}) * P(H_{(1, k-1)}) * P(M_k | H_k) * PF(H_{k-1} | H_k) \end{aligned}$$

であり、従って、短くすると、

30

$$F(M, H) = F(M_{(1, k)}, H_{(1, k)}) = F(M_{(1, k-1)}, H_{(1, k-1)}) * P(M_k | H_k) * PF(H_{k-1}, H_k)$$

であり、すなわち、 k 個の SNP についての F の計算を $k - 1$ 個の SNP についての F の計算に変えることができる。

【 0 1 2 1 】

$H = (H_1, \dots, H_k)$ については、 k 個の SNP についての仮説：

【 0 1 2 2 】
【 数 1 8 】

$$\max_H F(M, H) = \max_{(H_{(1,k-1)}, H_k)} F(M, (H_{(1,k-1)}, H_k)) = \max_{H_k} \max_{H_{(1,k-1)}} F(M, (H_{(1,k-1)}, H_k)) = \max_{H_k} G(M_{(1,k)}, H_k)$$

40

であり、ここで、

【 0 1 2 3 】

【数 1 9】

$$\begin{aligned}
 G(M_{(1,n)}, H_n) &= \max_{H_{(1,n-1)}} F(M_{(1,n)}, (H_{(1,n-1)}, H_n) = \\
 &= \max_{H_{(1,n-1)}} F(M_{(1,n-1)}, H_{(1,n-1)}) * P(M_n | H_n) * PF(H_{n-1}, H_n) = \\
 &= P(M_n | H_n) * \max_{H_{(1,n-1)}} F(M_{(1,n-1)}, H_{(1,n-1)}) * PF(H_{n-1}, H_n) = \\
 &= P(M_n | H_n) * \max_{H_{n-1}} \max_{H_{(1,n-2)}} F(M_{(1,n-1)}, (H_{(1,n-2)}, H_{n-1})) * PF(H_{n-1}, H_n) = \\
 &= P(M_n | H_n) * \max_{H_{n-1}} PF(H_{n-1}, H_n) * G(M_{(1,n-1)}, H_{n-1})
 \end{aligned}$$

である。

10

【0 1 2 4】

これをまとめると：

【0 1 2 5】

【数 2 0】

$$\max_H F(M, H) = \max_{H_n} G(M_{(1,k)}, H_k)$$

であり、ここで、G が帰納的に見出すことができ：n = 2, . . . , k については、

【0 1 2 6】

【数 2 1】

$$G(M_{(1,n)}, H_n) = P(M_n | H_n) * \max_{H_{n-1}} [PF(H_{n-1}, H_n) * G(M_{(1,n-1)}, H_{n-1})]$$

20

および $G(M_{(1,1)}, H_1) = 0.25 * P(M_1 | H_1)$ である。

【0 1 2 7】

該アルゴリズムは以下の通りである：

n = 1 については：4 つの仮説 H_{1i} を作り出し、 $i = 1, . . . , 4$ について $G(M_1 | H_{1i})$ を計算する。n = 2 については： H_{2i} について4 つの仮説を作り出し、式：

【0 1 2 8】

【数 2 2】

$$G(M_{(1,2)}, H_{2i}) = P(M_2 | H_{2i}) * \max_{j=1, \dots, 4} [PF(H_{1j}, H_{2i}) * G(M_1, H_{1j})]$$

30

を用い、一定時間内に、 $G(M_{(1,2)} | H_{2i}), i = 1, . . . , 4$ を計算する。n = k については： H_{ki} について4 つの仮説を作り出し、

【0 1 2 9】

【数 2 3】

$$G(M_{(1,k)}, H_{ki}) = P(M_k | H_{ki}) * \max_{j=1, \dots, 4} [PF(H_{k-1j}, H_{ki}) * G(M_{(1,k-1)}, H_{k-1j})]$$

によって、 $G(M_{(1,k)} | H_{ki}), i = 1, . . . , 4$ を作成する。

いずれの時点においても、覚えておくべき4 つのみの仮説、および一定数の操作がある。

40

従って、アルゴリズムは、指数関数とは反対に、SNP の数 k において線形である。

【0 1 3 0】

直線的時間内における P (M) の解

P (M) について解いて、最良の仮説を得る必要はない。というのは、それは全ての H について一定だからである。しかしながら、条件付確率 $P(H | M) = P(M | H) * P(H) / P(M)$ についての現実的な意味のある数を得るためには、P (M) を導く必要もある。前記したように、

【0 1 3 1】

【数 2 4】

$$\begin{aligned}
 P(M) &= P(M_{(1,k)}) = \sum_{H_{(1,k)}} P(M_{(1,k)} | H_{(1,k)}) * P(H_{(1,k)}) \\
 &= \sum_{H_k} P(M_K | H_k) \sum_{H_{(1,k-1)}} P(M_{(1,k-1)} | H_{(1,k-1)}) * P(H_{(1,k-1)}) * PF(H_{k-1}, H_k) \\
 &= \sum_{H_k} P(M_K | H_k) * W(M_{(1,k-1)} | H_k)
 \end{aligned}$$

と書くことができ、ここで、

【0 1 3 2】

【数 2 5】

$$W(M_{(1,k-1)} | H_k) = \sum_{H_{(1,k-1)}} P(M_{(1,k-1)} | H_{(1,k-1)}) * P(H_{(1,k-1)}) * PF(H_{k-1}, H_k)$$

である。

帰納：

【0 1 3 3】

【数 2 6】

$$\begin{aligned}
 W(M_{(1,k-1)} | H_k) &= \sum_{H_{(1,k-1)}} P(M_{(1,k-1)} | H_{(1,k-1)}) * P(H_{(1,k-1)}) * PF(H_{k-1}, H_k) \\
 &= \sum_{H_{k-1}} P(M_{k-1} | H_{k-1}) \sum_{H_{(1,k-2)}} P(M_{(1,k-2)} | H_{(1,k-2)}) * P(H_{(1,k-2)}) * PF(H_{k-2}, H_{k-1}) * PF(H_{k-1}, H_k) \\
 &= \sum_{H_{k-1}} P(M_{k-1} | H_{k-1}) * PF(H_{k-1}, H_k) * W(M_{(1,k-2)} | H_{k-1})
 \end{aligned}$$

によって $W(M, H)$ について解くことができ、従って、簡単に述べると、サイズ k の問題は、

【0 1 3 4】

【数 2 7】

$$W(M_{(1,k-1)} | H_k) = \sum_{H_{k-1}} P(M_{k-1} | H_{k-1}) * PF(H_{k-1}, H_k) * W(M_{(1,k-2)} | H_{k-1})$$

$$\text{および } W(M_{(1,1)} | H_2) = \sum_{H_1} P(M_1 | H_1) * 0.25 * PF(H_1, H_2)$$

によってサイズ $(k - 1)$ の問題に変えられる。前記したように、 $n = 2 : k$ については、最後に、

【0 1 3 5】

【数 2 8】

$$P(M) = \sum_{H_k} P(M_K | H_k) * W(M_{(1,k-1)} | H_k)$$

を誘導することが可能となるまで、帰納的に $W(2), \dots, W(K) = W(M_{(1, k-1)} | H_k)$ を得る。

【0 1 3 6】

各レベルにおいて、4つの異なる仮説 H_k があるに過ぎず、従って、アルゴリズムは、再度、SNP k の数において線形である。

【0 1 3 7】

直線的時間内における個々の SNP 信頼性

一旦、最良の仮説 $H^* = (H_1^*, \dots, H_k^*)$ が計算されたならば、今度は、各 SNP についての最終的な解答における信頼性、すなわち、 $i = 1, \dots, k$ についての $P(H_i^* | M)$ を誘導することが望まれるであろう。前記したように、 $P(H_i^* | M$

10

20

30

40

50

) = P (M | H_i^{*}) P (H_i^{*}) / P (M) = W (H_i^{*} , M) / P (M) であり、こ
 ことで、 P (M) は既に知られている。

【 0 1 3 8 】

【 数 2 9 】

$$W(M, H_i^*) = \sum_{H, H_i=H_i^*} P(M|H) * P(H) = \sum_{H=(H_{(1,j-1)}, H_i^*, H_{(i+1,k)})} P(M|H) * P(H)$$

であり、すなわち、仮説 H は最初の i - 1 の SNP、 i 番目の SNP についての仮説、お
 よび i + 1 ないし k 番目の SNP についての仮説まで終えている。前記したように：

【 0 1 3 9 】

【 数 3 0 】

$$P(M_{(1,k)} | H_{(1,k)}) = \prod_{j=1}^k P(M_j | H_j) = \prod_{j=1}^{i-1} P(M_j | H_j) * P(M_i | H_i^*) * \prod_{j=i+1}^k P(M_j | H_j)$$

$$= P(M_{(1,i-1)} | H_{(1,i-1)}) * P(M_i | H_i^*) * P(M_{(i+1,k)} | H_{(i+1,k)})$$

および

【 0 1 4 0 】

【 数 3 1 】

$$P(H_{(1,k)}) = 1/4 * \prod_{j=2}^k PF(H_{j-1}, H_j)$$

$$= 1/4 * \prod_{j=2}^{i-1} PF(H_{j-1}, H_j) * PF(H_{i-1}, H_i^*) * PF(H_{i-1}, H_i^*) * \prod_{j=i+2}^k PF(H_{j-1}, H_j)$$

$$= 1/4 * T(H_{(1,i-1)}) * PF(H_{i-1}, H_i^*) * PF(H_{i-1}, H_i^*) * T(H_{(i+1,k)})$$

であり、したがって、

【 0 1 4 1 】

【 数 3 2 】

$$P(H_{(1,k)}) = 1/4 * T(H_{(1,k)}) = 1/4 * T(H_{(1,i-1)}) * PF(H_{i-1}, H_i^*) * PF(H_{i-1}, H_i^*) * T(H_{(i+1,k)})$$

であり、ここで、

【 0 1 4 2 】

【 数 3 3 】

$$T(H_{(1,k)}) = \prod_{j=2}^k PF(H_{j-1}, H_j)$$

である。これから、

【 0 1 4 3 】

10

20

30

【数 3 4】

$$\begin{aligned}
W(M_{(1,k)}, H_i^*) &= \sum_{H, H_i=H_i^*} P(M|H) * P(H) = \sum_{H, H_i=H_i^*} P(M|H) * 1/4 * T(H) \\
&= \sum_{H=(H_{(1,i-1)}, H_i^*, H_{(i+1,k)})} P(M_{(1,i-1)} | H_{(1,i-1)}) * P(M_i | H_i^*) * P(M_{(i+1,k)} | H_{(i+1,k)}) * \\
&\quad 1/4 * T(H_{(1,i-1)}) * PF(H_{i-1}, H_i^*) * PF(H_{i-1}, H_i^*) * T(H_{(i+1,k)}) \\
&= 4 * P(M_i | H_i^*) * \left(\sum_{H_{i-1}} P(M_{(1,i-1)} | H_{(1,i-1)}) * 1/4 * T(H_{(1,i-1)}) * PF(H_{i-1}, H_i^*) \right) \\
&\quad * \left(\sum_{H_{i+1}} P(M_{(i+1,k)} | H_{(i+1,k)}) * 1/4 * T(H_{(i+1,k)}) * PF(H_i^*, H_{i+1}) \right) \\
&= 4 * P(M_i | H_i^*) * \left(\sum_{H_{i-1}} W(M_{(1,i-1)}, H_{i-1}) * PF(H_{i-1}, H_i^*) \right) * \left(\sum_{H_{i+1}} W(M_{(i+1,k)}, H_{i+1}) * PF(H_i^*, H_{i+1}) \right)
\end{aligned} \tag{10}$$

を示すのは可能である。

【0 1 4 4】

再度、サイズ k の場合は、前記したよりも複雑なビットであるにもかかわらず、より小さなサイズの 2 つのピースに変えられている。ピースの各々は

【0 1 4 5】

【数 3 5】

$$\begin{aligned}
W(M_{(1,n)}, H_n) &= P(M_n | H_n) * \left(\sum_{H_{n-1}} W(M_{(1,n-1)}, H_{n-1}) * PF(H_{n-1}, H_n) \right) \\
W(M_{(m,k)}, H_m) &= P(M_m | H_m) * \left(\sum_{H_{m+1}} W(M_{(m+1,k)}, H_{m+1}) * PF(H_m, H_{m+1}) \right)
\end{aligned}$$

として計算することができる。従って、アルゴリズムは、4 つの異なる H_n 、 H_m の各々について $n = 1, \dots, k$ 、 $m = k, \dots, 1$ について、 $W(M_{(1,n)}, H_n)$ 、 $W(M_{(m,k)}, H_m)$ を計算し、次いで、必要に応じてそれらを組合せて $i = 1, \dots, k$ について $W(M_{(1,k)}, H_i^*)$ を計算する。操作の数は依然として k について直線的である。 30

【0 1 4 6】

データの小さなまたは異なるセットが利用可能である場合の、胚データへの開示された方法の適用

システムの 1 つの実施形態において、親のいずれかまたは双方からのハプロイドデータの有りまたは無しにて、かつそのデータがより高いまたはより低い程度の確実性まで知られている場合、1 人の親（恐らくは母親）からのジプロイドデータを利用する必要があるに過ぎない。例えば、卵の供与の厳しい性質を仮定すれば、母性ハプロイドデータが容易に入手できない場合があると予測される。この明細書を読んだ後に、どのようにして、特定の SNP の尤度を計算するための統計学的方法を限定されたデータセットを仮定して修飾できるかは当業者に明らかとなろう。 40

【0 1 4 7】

別のアプローチが、一方または双方の親の失われたジプロイドまたはハプロイドデータを補うためにより距離がある親族からのデータを用いる。例えば、個人の染色体の 1 つのセットは彼または彼女の親の各々に由来することが知られているので、母方祖父母からのジプロイドデータを用いて、失われたまたは貧弱にしか測定されていない母性ハプロイドデータを部分的に再構築できよう。

【0 1 4 8】

この方法の帰納的性質に注意し：適当な祖父母のジプロイドおよび/またはハプロイドデータと共に、単一細胞親ハプロイドデータの天然ではノイジーな測定を仮定し、開示さ 50

れた方法を用いて、親ハプロイドデータを清浄化することができ、これは、今度は、胚のより正確なゲノタイピングを供するであろう。これらの場合に用いる方法をどのようにして修飾するかは当業者に明らかなはずである。

【0149】

より少ないよりはむしろより多くの情報を用いるのが好ましい。というのはこれは所与のSNPにおいて正しい要求を行うチャンスを増大させることができ、かつそれらの要求において信頼性を増加させることができるからである。これは、システムの増大する複雑性とバランスしなければならない。というのは、データのさらなる技術および源を用いるからである。データを増大させるのに情報を用いるために利用できるさらなる情報、ならびに技術の多くの源がある。例えば、Hapmapデータ、またはゲノムデータの他のレパートリーで見出すことができる相関関係を利用するインフォマティクスを基礎としたアプローチがある。加えて、そうでなければイン・シリコにて再度作り出す必要がある遺伝子データの直接的測定を可能とできる生物学的アプローチがある。例えば、そうでなければ利用できないハプロイドデータは、フローサイトメトリー技術を用いてジプロイド細胞から個々の染色体を抽出して、蛍光タグド染色体を単離することによって測定可能であろう。別法として、細胞融合を用いて、一対立遺伝子ハイブリッド細胞を作り出して、ジプロイドからハプロイドへの変換を行うことができる。

10

【0150】

いずれの胚が着床するようであるかを選択することへの開示された方法の適用

1つの実施形態において、システムを用いて、母親に着床し、ベビーまで発生する胚の尤度を決定することができる。胚着床の尤度が胚のSNP、および/または母親のSNPに対するそれらの関係によって決定される程度まで、開示された方法は、いずれがクリーンなSNPデータに基づいて成功して着床するかの信頼性ある予測をなすことをベースとして、胚の選択を助けるにおいて重要であろう。尤度を最良に予測するためには、胚における遺伝子発現のレベル、母親における遺伝子発現のレベル、および/または母親の決定された遺伝子型と恐らくは組合された胚の決定された遺伝子型を考慮する必要がある。

20

【0151】

加えて、異数性胚はあまり着床しないようであり、成功した妊娠をもたらさないようであり、健康な子供をもたらさないようであることはよく知られている。結果として、異数体についてのスクリーニングは、成功した結果を最ももたらすようである胚の選択に対して重要な面である。このアプローチについてのより詳細は以下に掲げる。

30

【0152】

親ハプロイドデータの推定

該方法の1つの実施形態において、親のジプロイドデータの詳細な知識を仮定し、親はプロイドを推定する必要がある。これを行うことができる多数の方法がある。最も単純な場合において、ハプロタイプは、直接的関係(母親、父親、息子または娘)の単一ハプロイド細胞の分子アレイによって既に推定されている。この場合、分子からアッセイによって測定されたジプロイド遺伝子型からの公知のハプロイドを差し引くことによって姉妹ハプロイドを推定するのは当業者にとってたやすいことである。例えば、もし特定の遺伝子座がヘテロ接合性であれば、未知の親ハプロイドは公知の親ハプロタイプからの反対の対立遺伝子である。

40

【0153】

もう1つの場合において、親のノイジーなハプロイドデータは、精子のような個々の親ハプロイド細胞の分子生物学的ハプロタイピングから、または磁性ビーズおよびフローサイトメトリーを含めた種々の方法によって単離することができる個々の染色体から知ることができる。この場合、決定されたハプロタイプが測定されたハプロタイプと同程度にノイジーであることを除いて、同一手法を前記したように用いることができる。

【0154】

また、(公のHapmapプロジェクトで作られされたもののような)一般的集団にお

50

ける公知のハプロタイプブロックを利用する統計学的方法を用い、ジブroidデータから直接的にハプロイドデータセットを推定する方法もある。ハプロタイプブロックは、本質的には、種々の集団において反復して起こる一連の関連する対立遺伝子である。これらのハプロタイプブロックはしばしば古くかつ共通するので、それらを用いて、ジブroid遺伝子型からハプロイドを予測することができる。次いで、親の推定されたハプロイドブロックを本明細書中に記載された方法のために入力として用いて、胚からのノイジーなデータを清浄化することができる。この仕事を達成する公に入手可能なアルゴリズムは、不完全な系統発生アプローチ、共役事前分布、および集団遺伝学からの事前分布に基づくベイズアプローチを含む。これらのアルゴリズムのいくつかは隠れたMarkovモデルを用いる。1つの研究は、公のトリオおよび無関係な個々のデータを用いて、これらのアルゴリズムが1MBの配列にわたって0.05%と低い誤差率にて実行されることを示した。しかしながら、予測されるように、精度は稀なハプロタイプブロックを持つ個人についてより低い。1つの見積もりにおいて、計算方法は、20%のわずかな対立遺伝子頻度にて5.1%の遺伝子座と多くを同調できなかった。

10

【0155】

本発明の1つの実施形態において、IVFサイクルの間に異なる胚から取られた多数の胚盤胞からの遺伝子データを用いて、より大きな信頼性をもって親のハプロタイプブロックを推定する。

【0156】

高および中程度スループットのゲノタイピングを用いて異数性をスクリーニングするための技術

20

システムの1つの実施形態において、測定された遺伝子データを用いて、個体において異数体および/またはモザイク現象の存在について検出することができる。本明細書中に開示するのは、これらの試料からの増幅されたまたは増幅されていないDNAからの染色体の数またはDNAセグメントコピー数を検出するための中程度または高-スループットゲノタイピングを用いるいくつかの方法である。目標は、Illumina、AgilentおよびAffymetrixからのABI Taqman、MIPS、またはマイクロアレイのような異なる定量的および/または定性的ゲノタイピングプラットフォームを用いて異数性のあるタイプおよびモザイク現象のレベルの検出において達成することができる信頼性を見積もることである。これらの場合の多くにおいて、遺伝物質はゲノタイピングアレイ上のプローブへのPCRによって増幅して、特定の対立遺伝子の存在を検出する。これらのアッセイをゲノタイピングでどのようにして用いるかは本開示において他の箇所に記載されている。

30

【0157】

以下に記載するのは、欠失、異数体および/またはモザイク現象から生起するかに拘らず、異常な数のDNAセグメントについてスクリーニングするいくつかの方法である。該方法は以下のようにグループ分けされる：(i)対立遺伝子要求を行うことのない定量的技術；(ii)対立遺伝子要求を活用する定性的技術；(iii)対立遺伝子要求を活用する定量的技術；(iv)各遺伝子座における遺伝子データの増幅についての確率分布関数を用いる技術。全ての方法は、標的個体のゲノムにおける所与のセグメントの存在の数を決定するための、所与の染色体の所与のセグメント上の多数の遺伝子座の測定を含む。加えて、該方法は、所与のセグメントの存在の数についての1以上の仮説のセットを作り出し；所与のセグメント上の多数の遺伝子座における遺伝子データの量を測定し；標的個体遺伝子データの測定を仮定して、仮説の各々の相対的確率を決定し；次いで、所与のセグメントの存在の数を決定するために、各仮説に関連する相対的確率を用いることを含む。さらに、該方法は、全て、多数の遺伝子座における遺伝子データの量の測定の計算された関数である組合せ測定Mを作り出すことを含む。全ての方法において、閾値は、測定Mに基づいて各仮説 H_i の選択について決定され、測定すべき遺伝子座の数を見積もって、仮説の各々の偽検出の特定のレベルを有するようにする。

40

【0158】

50

測定Mを仮定して各仮説の確率は $P(H_i | M) = P(M | H_i) P(H_i) / P(M)$ である。 $P(M)$ は H_i から独立しているので、 $P(M | H_i) P(H_i)$ のみを考慮することによってMを仮定した仮説の相対的確率を決定することができる。以下において、技術の分析および比較を単純化するために、我々は、我々が $P(M | H_i)$ のみを考慮することによって全ての $P(H_i | M)$ の相対的確率を計算できるように、 $P(H_i)$ は全ての $\{H_i\}$ について同一である。その結果、閾値、および測定すべき遺伝子座の数の我々の決定は、 $P(H_i)$ が全ての $\{H_i\}$ について同一であるという仮定の下で偽仮説を選択する特定の確率を有することに基づく。この開示を読んだ後に、どのようにしてアプローチを修飾して、 $B(H_i)$ がセット $\{H_i\}$ において異なる仮説で変化するという事実を受け入れるであろうかは当業者に明瞭である。いくつかの実施形態において、全ての i にわたって $P(H_i | M)$ を最大化する仮説 H_{i^*} が選択されるように閾値を設定する。しかしながら、閾値は $P(H_i | M)$ を最大化するように必ずしも設定される必要はないが、むしろ、セット $\{H_i\}$ における異なる仮説の間の偽検出の確率の特定の比率を達成するように設定される必要がある。

10

【0159】

異数体を検出するための本明細書中で言及する技術は、片親二染色体、バランスしないトランスロケーションについて、および染色体の性別分け（男性または女性；XYまたはXX）について検出するのに等しく良く用いることができることに注意するのは重要である。概念の全ては、所与の試料に存在する染色体（または染色体のセグメント）の同一性および数を検出することに関連し、かくして、全ては、本処理に記載された方法によって

20

【0160】

マッチドフィルタリングの概念

ここに適用される方法は、デジタルシグナルの最適検出において適用されるのに同様である。正常に分布したノイズの存在下においてシグナル - ノイズ比率(SNR)を最大化する最適アプローチは、可能なノイズ - フリーシグナルの各々に対応する、理想化されたマッチングシグナル、またはマッチドフィルタを形成すること、およびこのマッチドシグナルを受け取られたノイズなシグナルと相関させることは、Schwartz不均衡を用いて示すことができる。このアプローチは、可能なシグナルのセット、ノイズの統計学的分布 - 平均および標準偏差(SD)が公知であることを必要とする。ここで、染色体、またはDNAのセグメントが試料中に存在するかまたは存在しないことを検出する一般的なアプローチを記載する。全染色体を調べることに、または挿入されたまたは欠失された染色体セグメントを調べることに間に差を設けない。この記載を読んだ後に、どのようにして、該技術を異数性および性別決定の多くのシナリオ、または胚、胎児または産まれた子供の染色体における挿入および欠失の検出まで拡大できるかは明らかである。このアプローチは、Taqman、qPCR、Illuminaアレイ、Affymetrixアレイ、Agilentアレイ、MIPSキット等を含めた広い範囲の定量的および定性的ゲノタイピングプラットフォームに適用することができる。

30

【0161】

一般的問題の公式化

2つの対立遺伝子変異が起こる(xおよびy)SNPにおいてプローブがあると仮定する。各遺伝子座 i 、 $i = 1 \dots N$ において、2つの対立遺伝子からの遺伝物質の量に対応するデータを収集する。Taqmanアッセイにおいて、これらの尺度は、例えば、各対立遺伝子 - 特異的色素のレベルが閾値を交差するサイクル時間 C_t であろう。どのようにして、このアプローチを、各遺伝子座における、または遺伝子座における各対立遺伝子に対応する遺伝物質の量の異なる測定まで拡大できるかは明らかであろう。遺伝物質の量の定量的測定は非線形であり、その場合、注目するセグメントの存在によって引き起こされた特定の遺伝子座の測定の変化は、どのようにして、その遺伝子座の多くの他のコピーが他のDNAセグメントからの試料に存在するかに依存するであろう。いくつかの場合に

40

50

において、技術が、注目するセグメントの存在によって引き起こされた特定遺伝子座の測定の変化が、どのようにしてその遺伝子座の多くの他のコピーが他のDNAセグメントからの試料に存在するか依存しないように、線形測定を必要とするであろう。アプローチを、どのようにしてTaqManまたはqPCRアッセイからの測定を線形化することができるかについて記載するが、異なるアッセイについて適応できる非線形測定を線形化するための多くの他の技術がある。

【0162】

遺伝子座1...Nにおける対立遺伝子xの遺伝物質の量の測定は、データ $d_x = [d_{x_1} \dots d_{x_N}]$ によって与えられる。同様に、対立遺伝子yについては、データ $d_y = [d_{y_1} \dots d_{y_N}]$ によって与えられる。各セグメントjは、各要素 a_{j_i} がxまたはyいずれかである対立遺伝子 $a_j = [a_{j_1} \dots a_{j_N}]$ を有すると仮定する。対立遺伝子xの遺伝物質の量の測定データを、 s_x がシグナルであって、 n_x が擾乱である $d_x = s_x + n_x$ として記載する。該シグナルは $s_x = [f_x(a_{1_1}, \dots, a_{j_1}) \dots f_x(a_{j_N}, \dots, a_{j_N})]$ であり、ここで、 f_x は測定に対する対立遺伝子からのセットのマッピングであり、jがDNAセグメントコピーの数である。擾乱ベクトル n_x は測定誤差によって引き起こされ、非線形測定の場合においては、注目するDNAセグメント以外の他の遺伝子物質の存在によって引き起こされる。測定誤差は通常正規分布し、それらは、非線形によって引き起こされた擾乱に対して大きく（線形化測定についてのセクション参照）、従って、 $n_{x_i} \sim N(0, R)$ であり、ここで、非 n_{x_i} が偏差 $n_{x_i}^2$ を有し、ベクトル n_x は正規分布する $\sim N(0, R)$ 、 $R = E(n_x n_x^T)$ と仮定する。さて、いくつかのフィルターhをこのデータに適応して、測定 $m_x = h^T d_x = h^T s_x + h^T n_x$ を行うと仮定する。ノイズに対するシグナルの比率($h^T s_x / h^T n_x$)を最大化するためには、hはマッチドフィルター $h = \mu R^{-1} s_x$ によって与えられ、ここで、 μ はスケール定数であることを示すことができる。対立遺伝子xについての議論は対立遺伝子yについて反復することができる。

【0163】

方法1a：各遺伝子座についての平均および標準偏差が知られている場合に、対立遺伝子要求を行わない定量的技術による異数性または性別の測定

このセクションでは、データは、（例えば、qPCRを用いる）対立遺伝子値に拘わらず遺伝子座における遺伝物質の量に関係し、またはデータは、集団において100%浸透度を有する対立遺伝子についてのみであると仮定し、あるいはデータは、各遺伝子座における多数の対立遺伝子において組合せて（線形化測定についてのセクション参照）、その遺伝子座における遺伝物質の量を測定すると仮定する。その結果、このセクションにおいては、データ d_x に言及でき、 d_y を無視することができる。また、2つの仮説：DNAセグメントの2つのコピーがある h_0 （これらは、典型的には同一のコピーではない）およびただ1つのコピーがある h_1 、があると仮定する。各仮説については、データは、各々、 $d_{x_i}(h_0) = s_{x_i}(h_0) + n_{x_i}$ および $d_{x_i}(h_1) = s_{x_i}(h_1) + n_{x_i}$ として記載でき、ここで、 $s_{x_i}(h_0)$ は、2つのDNAセグメントが存在する場合に、遺伝子座iにおける遺伝物質の予測される測定（予測されるシグナル）であり、 $s_{x_i}(h_1)$ は1つのセグメントについて予測されるデータである。仮説 $h_0 : m_{x_i} = d_{x_i} - s_{x_i}(h_0)$ についての予測されるシグナルを差分することによって各遺伝子座についての測定を構築する。もし h_1 が真であれば、測定の予測される値は $E(m_{x_i}) = s_{x_i}(h_1) - s_{x_i}(h_0)$ である。先に議論したマッチドフィルターを用い、 $h = (1/N) R^{-1} (s_{x_i}(h_1) - s_{x_i}(h_0))$ を設定する。測定は $m = h^T d_x = (1/N) \sum_{i=1}^N ((s_{x_i}(h_1) - s_{x_i}(h_0)) / \sqrt{n_{x_i}^2}) m_{x_i}$ と記載される。

【0164】

もし h_1 が真であれば、 $E(m | h_1) = m_1 = (1/N) \sum_{i=1}^N (s_{x_i}(h_1) - s_{x_i}(h_0)) / \sqrt{n_{x_i}^2}$ の予測される値、およびmの標準偏差は $\sigma_{m | h_1}^2 = (1/N^2) \sum_{i=1}^N ((s_{x_i}(h_1) - s_{x_i}(h_0))^2 / n_{x_i}^2)$

伝子要求を行わない定量的技術による異数性または性別の測定

各遺伝子座の特徴がよく知られていない場合、各遺伝子座における全てのアッセイが同様に挙動し、すなわち、その代わりに、 $E(m_x)$ および s_x のみに言及するのが可能であるように、 $E(m_{x_i})$ および s_{x_i} は全ての遺伝子座 e にわたって一定であるという単純化仮定をすることができる。この場合、マッチドフィルタリングアプローチ $m = h^T d_x$ は d_x の分布の平均を見出すことに変えられる。このアプローチは平均の比較といい、それは、真実のデータを用いる異なる種類の検出で必要とされる遺伝子座の数を見積もるのに用いられるであろう。

【0169】

前記したように、試料に存在する2つの染色体(仮説 h_0) または存在する1つの染色体(h_1)がある場合のシナリオを考える。 h_0 では、分布は $N(\mu_0, \sigma_0^2)$ であり、 h_1 については、分布は $N(\mu_1, \sigma_1^2)$ である。各々、測定された試料平均および $SD: m_1, m_0, s_1$ および s_0 を持つ N_0 および N_1 試料を用いて分布の各々を測定する。平均は、 $M_0 \sim N(\mu_0, \sigma_0^2 / N_0)$ および $M_1 \sim N(\mu_1, \sigma_1^2 / N_1)$ として正規分布するランダム変数 M_0, M_1 としてモデル化することができる。 $M_1 \sim N(m_1, s_1^2 / N_1)$ および $M_0 \sim N(m_0, s_0^2 / N_0)$ と仮定することができるように、 N_1 および N_0 は十分に大きい (> 30) と仮定する。分布が異なるか否かを検定するために、平均検定の差を用いることができ、ここで、 $D = m_1 - m_0$ である。ランダム変数 D の偏差は $\sigma_d^2 = \sigma_1^2 / N_1 + \sigma_0^2 / N_0$ であり、これは $\sigma_d^2 = s_1^2 / N_1 + s_0^2 / N_0$ と近似することができる。 h_0 を与えると、 $E(d) = 0$ となり； h_1 を与えると、 $E(d) = \mu_1 - \mu_0$ となる。 h_1 および h_0 の間の要求を行うための異なる技術をここに議論する。

【0170】

X染色体上の48 SNPを用いるTaqmanアッセイの異なる実行で測定されたデータを用いて、性能をキャリブレートした。試料1は、1つのX染色体を含有する混合男性起源のウェル当たりおよそ0.3 ngのDNAよりなり；試料0は、2つのX染色体を含有する混合女性起源のウェル当たりおよそ0.3 ngのDNAよりなるものであった。 $N_1 = 42$ および $N_0 = 45$ 。図7および図8は、試料1および0についてのヒストグラムを示す。これらの試料についての分布は $m_1 = 32.259$ 、 $s_1 = 1.460$ 、 $m_1 = s_1 / \text{sqrt}(N_1) = 0.225$ ； $m_0 = 30.75$ ； $s_0 = 1.202$ 、 $m_0 = s_0 / \text{sqrt}(N_0) = 0.179$ によって特徴付けられる。これらの試料では、 $d = 1.509$ および $\sigma_d = 0.2879$ である。

【0171】

このデータは混合男性および女性試料に由来するので、標準偏差の多くは、混合試料中の各SNPにおける異なる対立遺伝子頻度によるものである。SDは、多数の実行にわたり、一定時刻における1つのSNPについての C_t における変動を考慮することによって見積もられる。このデータを図9に示す。ヒストグラムは0の周りに対称である。というのは、各SNPについての C_t は2つの実行または実験で測定され、各SNPについての C_t の平均値は差し引かれるからである。2つの実行を用いる混合男性試料中の20のSNPにわたる平均標準偏差は $s = 0.597$ である。このSDは男性および女性双方の試料で保存的に用いられる。というのは、女性試料についてのSDは男性試料についてよりも小さいだろうからである。加えて、混合試料は全てのSNPについてヘテロ接合性であると推定されるので、ただ1つの色素からの測定が用いられていることを注記する。双方の色素の使用は、遺伝子座における各対立遺伝子の測定が組み合わせられることを必要とし、これはより複雑である(線形化測定についてのセクション参照)。双方の色素についての測定の組合せはシグナルの振幅を2倍とし、およそ $\text{sqrt}(2)$ によってノイズ振幅を増大させ、その結果、およそ $\text{sqrt}(2)$ または3 dBのSNR改良がもたらされる。

【0172】

モザイク現象なしおよび参照試料なしを仮定する検出

10

20

30

40

50

m_0 が多くの実験から完全に知られており、かつかく実験の実行は、 m_1 を計算して m_0 と比較するのにただ 1 つの試料を実行すると仮定する。 n_1 はアッセイの数であり、各アッセイは異なる SNP 遺伝子座であると仮定する。閾値 t は m_0 および m_1 の間に設定して、偽陰性の尤度を偽陰性の数と等しくすることができ、もしそれが閾値を超えれば、試料は異常であると記される。 $s_1 = s_2 = s = 0.597$ であると仮定し、偽陰性または陽性の確率が $1 - \text{normcdf}(5, 0, 1) = 2.87e-7$ となるように 5 - シグマアプローチを用いる。目標は $5s_1 / \text{sqr t}(N_1) < (m_1 - m_0) / 2$ 、従って、 $N_1 = 100s_1^2 / (m_1 - m_0)^2 = 16$ についてのものである。さて、有害なシナリオである、偽陽性の偽陰性の確率よりも高くされるアプローチを用いることもできる。もし陽性を測定すれば、実験は再度行うことができる。その結果、偽陰性の確率は偽陽性の確率の平方と等しいはずであるということが可能である。図 3 を考え、 $t =$ 閾値とし、 $\text{シグマ}0 = \text{シグマ}1 = s$ と仮定する。かくして、 $1 - \text{normcdf}((t - m_0) / s, 0, 1))^2 = 1 - \text{normcdf}((m_1 - t) / s, 0, 1)$ である。これを解き、 $t = m_0 + 0.32(m_1 - m_0)$ であることを示すことができる。よって、目標は $5s / \text{sqr t}(N_1) < m_1 - m_0 - 0.32(m_1 - m_0) = (m_1 - m_0) / 1.47$ 、よって、 $N_1 = (5^2)(1.47^2)s^2 / (m_1 - m_0)^2 = 9$ についてのものである。

【0173】

参照試料を実行することのないモザイク現象での検出

目標は 97.7% の確率でモザイク現象を検出することである (すなわち、2 - シグマアプローチ) 以外は前記したのと同様状況を仮定する。これは、およそ 20 の細胞を抽出し、それらの写真を撮る羊水穿刺に対する標準アプローチよりも良好である。もし 20 細胞のうち 1 が異数体であって、これは 100% の信頼性でもって検出されると仮定するならば、標準アプローチを用いる異数体である群の少なくとも 1 つを有する確率は $1 - 0.95^{20} = 64\%$ である。もし細胞の 0.05% が異数体であれば (この試料 3 を要求する、 $m_3 = 0.95m_0 + 0.05m_1$ および $\text{var}(m_3) = (0.95s_0^2 + 0.05s_1^2) / N_1$ である。かくして、 $\text{std}(m_3)^2 < (m_3 - m_0) / 2 = \text{sqr t}(0.95s_0^2 + 0.05s_1^2) / \text{sqr t}(N_1) < 0.05(m_1 - m_2) / 4 = N_1 = 16(0.95s_2^2 + 0.05s_1^2) / (0.05^2(m_1 - m_2)^2) = 1001$ である。慣用的アプローチを用いて達成することができるよりも依然として良好な (すなわち、84.1% 確率での検出) 1 - シグマ統計学の目標を用い、同様にして $N_1 = 250$ であると示すことができる。

【0174】

モザイク現象がなく、参照試料を用いる検出

このアプローチは必要でないかもしれないが、各実験は 2 つの試料を実行して、 m_1 を真実の試料 m_2 と比較すると仮定する。 $N = N_1 = N_0$ と仮定する。 $d = m_1 - m_0$ を計算し、 $s_1 = s_2 = s$ と仮定し、閾値 $t = (m_0 + m_1) / 2$ を設定し、従って偽陽性および偽陰性の確率は等しい。偽陰性の確率を $2.87e-7$ とし、それは、 $(m_1 - m_2) / 2 > 5 \text{sqr t}(s_1^2 / N + s_2^2 / N) = N = 100(s_1^2 + s_2^2) / (m_1 - m_2)^2 = 32$ があてはまらなければならない。

【0175】

モザイク現象での検出および参照試料の実行

前記したように、偽陰性の確率は 2.3% であると仮定する (すなわち、2 - シグマアプローチ) もし細胞の 0.05% が異数体であれば (これを試料 3 と呼ぶ)、 $m_3 = 0.95m_0 + 0.05m_1$ および $\text{var}(m_3) = (0.95s_0^2 + 0.05s_1^2) / N_1$ である。 $d = m_3 - m_2$ および $d^2 = (1.95s_0^2 + 0.05s_1^2) / N$ である。 $\text{std}(m_3)^2 < (m_0 - m_2) / 2 = \text{sqr t}(1.95s_2^2 + 0.05s_1^2) / \text{sqr t}(N) < 0.05(m_1 - m_2) / 4 = N = 16(1.95s_2^2 + 0.05s_1^2) / (0.05^2(m_1 - m_2)^2) = 2002$ でなければならない。再度 1 - シグマアプローチを用い、 $N = 500$ であることが同様に示すことができる

。

【0176】

目標が、現在の技術水準におけるように、64%の確率でもって5%モザイク現象を検出するにすぎない場合を考える。従って、偽陰性の確率は36%となろう。換言すれば、 $1 - \text{normcdf}(x, 0, 1) = 36\%$ となるような x を見出す必要があるだろう。かくして、2-シグマアプローチについては $N = 4(0.36^2)(1.95s_2^2 + 0.05s_1^2) / (0.05^2(m_1 - m_2)^2) = 65$ であり、または1-シグマアプローチについては $N = 33$ である。この結果、取り組むことが必要な、非常に高いレベルの偽陽性をもたらされることに注意されたし。というのは、偽陽性のそのようなレベルは現在実行可能な代替ではないからである。

10

【0177】

また、もし N が384に限定され、(すなわち、染色体当たり384ウェルTaqmanプレート)、かつ目標が97.72%の確率でモザイク現象を検出することになれば、1-シグマアプローチを用いて8.1%のモザイク現象を検出することが可能であろうことを注記する。84.1%の確率でもって(または15.9%偽陰性率でもって)モザイク現象を検出するには、1-シグマアプローチを用いて5.8%のモザイク現象を検出するのが可能である。97.72%の信頼性でもって19%のモザイク現象を検出するには、およそ70の遺伝子座を必要とするであろう。かくして、単一プレート上で5つの染色体についてスクリーニングできよう。

【0178】

これらの異なるシナリオの各々のまとめを表2に供する。また、この表2は、qPCRから得られた結果、およびSYBRアッセイを含める。前記した方法を用い、各遺伝子座についてのqPCRアッセイの性能は同一であるという単純化した仮定を行った。図10および図11は、前記したような、試料1および0についてのヒストグラムを示す。 $N_0 = N_1 = 47$ 。これらの試料についての測定分布は、 $m_1 = 27.65$ 、 $s_1 = 1.40$ 、 $m_0 = 26.64$ 、 $s_0 = 1.146$ 、 $m_1 = s_1 / \text{sqrt}(N_1) = 0.204$ 、 $m_0 = s_0 / \text{sqrt}(N_0) = 0.167$ によって特徴付けられる。これらの試料について、 $d = 1.01$ および $d = 0.2636$ である。図12は、0.75の全ての遺伝子座にわたる差の標準偏差での各遺伝子座についての男性および女性試料に対する C_t の間の差を示す。SDは、男性または女性試料での各遺伝子座の各測定について $0.75 / \text{sqrt}(2) = 0.53$ と近似した。

20

30

【0179】

方法2：対立遺伝子要求を用いる定性的技術

このセクションにおいては、アッセイは定量的であるという仮定をしない。その代わり、仮定は、対立遺伝子要求は定性的であって、アッセイに由来する意味のある定量的データはないというものである。このアプローチは、対立遺伝子要求を行ういずれのアッセイについても適当である。図13は、どのようにして、異なるハプロイド配偶子が減数分裂の間に形成されるか、およびそれを用いて、このセクションに関連する異なる種類の異数性を記載するのに用いる。最良のアルゴリズムは、検出されるべき異数性のタイプに依存する。

40

【0180】

異数性が、他の2つのセグメントのいずれかのコピーであるセクションを有しない第3のセグメントによって引き起こされる状況を考える。図13より、例えば、もし p_1 および p_4 、または p_2 および p_3 の双方が、他の親からの1つのセグメントに加えて、子供の細胞中で生起するならば、該状況は起きるであろう。これは、異数性を引き起こすメカニズムを仮定すれば、非常に普通である。1つのアプローチは、細胞中に2つのセグメントがある仮説 h_0 、およびこれらの2つのセグメントは何であるかでもって開始することである。説明の目的で、 h_0 は図13からの p_3 および m_4 についてのものであると仮定する。好ましい実施形態において、この仮説は本書類中の他の箇所に記載されたアルゴリズムに由来する。仮説 h_1 は、他のセグメントのコピーであるセクションを有しないさら

50

なるセグメントがあるというものである。これは、例えば、もしこの p_2 または m_1 もまた存在するならば正直であろう。 p_3 および m_4 においてホモ接合性である全ての遺伝子座を同定することが可能である。異数性は、ホモ接合性であると予測される遺伝子座におけるヘテロ接合性遺伝子型要求をサーチすることによって検出することができる。

【0181】

各遺伝子座は2つの可能な対立遺伝子 x および y を有すると仮定する。各々、対立遺伝子 x および y の確率は一般に p_x および p_y であり、および $p_x + p_y = 1$ であるとする。もし h_1 が真であれば、それについて p_3 および m_4 がホモ接合性である各遺伝子座 i について、非ホモ接合性要求の確率は、遺伝子座が、各々、 x または y においてホモ接合であるかに依存して p_y または p_x である。注意：親データ、すなわち、 p_1 、 p_2 、 p_4 および m_1 、 m_2 、 m_3 の知識に基づいて、各遺伝子座において非ホモ接合性対立遺伝子 x または y を有する確率をさらに改良することが可能である。これは、同一数の SNP での各仮説についてより信頼性のある測定を可能とするが、標記方法を複雑化し、従って、この延長は明示的には取り扱わない。どのようにしてこの情報を用いて、仮説の信頼性を増加させるかは当業者に明らかなはずである。

10

【0182】

対立遺伝子ドロップアウトの確率は p_d である。遺伝子座 i においてヘテロ接合性遺伝子型を見出す確率は仮説 h_0 を仮定すれば、 p_{0i} であり、仮説 h_1 を仮定すれば p_{1i} である。

【0183】

h_0 : $p_{0i} = 0$ とする。

20

【0184】

遺伝子座が x または y に対してホモ接合性であるかに依存して、 h_1 : $p_{1i} = p_x (1 - p_d)$ 、または $p_{1i} = p_y (1 - p_d)$ とする。

【0185】

測定 $m = 1 / N_h \quad i = 1 \dots N_h \quad I_i$ を作り出し、ここで、 I_i はインジケータ変数であり、もしヘテロ接合性要求がなされたならば、1 であって、その他の場合は 0 である。 N_h はホモ接合性遺伝子座の数である。 $p_x = p_y$ であって、全ての遺伝子座について p_{0i} 、 p_{1i} が同一の2つの値 p_0 および p_1 であると仮定することによって、説明を簡略化することができる。 h_0 を与えて、 $E(m) = p_0 = 0$ 、および ${}^2 m | h_0 = p_0 (1 - p_0) / N_h$ となる。 h_1 を与えて、 $E(m) = p_1$ および ${}^2 m | h_1 = p_1 (1 - p_1) / N_h$ となる。5シグマ - 統計学を用い、偽陽性の確率を偽陰性の確率と等しくし、 $(p_1 - p_0) / 2 > 5 \quad m | h_1$ 、よって、 $N_h = 100 (p_0 (1 - p_0) + p_1 (1 - p_1)) / (p_1 - p_0)^2$ と示すことができる。5 - シグマ信頼性の代わりに2 - シグマ信頼性では、 $N_h = 4 \cdot 2^2 (p_0 (1 - p_0) + p_1 (1 - p_1)) / (p_1 - p_0)^2$ と示すことができる。

30

【0186】

信頼性が少なくとも97.7%であるように(2 - シグマ)十分な入手可能なホモ接合性遺伝子座 $N_{h - available}$ があることは、十分な遺伝子座 N をサンプリングするのに必要である。 $N_{h - available} = i = 1 \dots N \quad J_i$ を特徴付け、ここで、 J_i は、もし遺伝子座がホモ接合性であれば値1のインジケータ変数であり、そうでなければ、0である。ホモ接合性である遺伝子座の確率は $p_x^2 + p_y^2$ である。その結果、 $E(N_{h - available}) = N (p_x^2 + p_y^2)$ 、および $N_{h - available}^2 = N (p_x^2 + p_y^2) (1 - p_x^2 + p_y^2)$ となる。 N が97.7%信頼性でもって十分に大きいことを補償するためには、 $E(N_{h - available}) - 2 \quad N_{h - available} = N_h$ でなければならない。ここで、 N_h は前記から見出される。

40

【0187】

例えば、もし $p_d = 0.3$ 、 $p_x = p_y = 0.5$ を仮定するならば、5 - シグマ信頼性について、 $N_h = 186$ および $N = 391$ を見出すことができる。同様に、2 - シグマ信頼性、すなわち、偽陰性および偽陽性における97.7%信頼性について、 $N_h = 30$ で

50

あって、 $N = 68$ であることを示すのは可能である。

【0188】

同様なアプローチを、 h_0 が2つの公知の染色体セグメントが存在する仮説であって、 h_1 が染色体セグメントの一方が失われている仮説である場合、セグメントの欠失を探すことに適応することができることを注記する。例えば、前記でなされたように、対立遺伝子ドロップアウトの効果をコードし、ヘテロ接合性であるが、ホモ接合性である遺伝子座を探すことが可能である。

【0189】

また、アッセイが定性的であったとしても対立遺伝子ドロップアウト率を用いて、存在するDNAセグメントの数についての定量的尺度のタイプを供することができることを注記する。

10

【0190】

方法3：参照配列の公知の対立遺伝子、および定量的対立遺伝子測定の使用

ここで、セグメントの清浄なまたは予測されるセットは知られていると仮定する。これらの染色体についてチェックするためには、各染色体の2つを仮定して、第一の工程はデータを正常化することである。本発明の好ましい実施形態において、第一の工程におけるデータ正常化は、本書類の他の箇所に記載された方法を用いてなされる。次いで、予測される2つのセグメントに関連するシグナルは、測定されたデータから差し引かれる。次いで、残りのシグナル中のさらなるセグメントを探すことができる。マッチドフィルタリングアプローチを用いて、さらなるセグメントを特徴付けるシグナルは、存在すると信じられるセグメントの各々、ならびにそれらの相補的染色体に基づく。例えば、図13をコードし、もしPSの結果が、セグメントp2およびm1が存在することを示すならば、本明細書中に記載された技術を用いて、さらなる染色体上でのp2、p3、m1、およびm4の存在をチェックすることができる。もし存在するさらなるセグメントがあれば、それは、これらのテストシグナルの少なくとも1つと共通する、50%を超える対立遺伝子を有することが保証される。ここに詳細に記載されていないもう1つのアプローチは、染色体の異常な番号、すなわち、1、3、4、および5染色体を仮定し、書類の他の箇所に記載されたアルゴリズムを用いて、データを正常化し、次いで、本明細書中で議論した方法を適用することができる。このアプローチの詳細は、本書類を読んだ後に当業者に明瞭なはずである。

20

30

【0191】

仮説 h_0 は、対立遺伝子ベクトル a_1 、 a_2 をもつ2つの染色体があるというものである。仮説 a_1 は、対立遺伝子ベクトル a_3 を持つ第三の染色体があるというものである。遺伝子データを正常化するために本書類に記載した方法、またはもう1つの技術を用い、各要素 a_{j_i} が x または y いずれかである $h_0 : a_1 = [a_{11} \dots a_{1N}]$ および $a_2 = [a_{21} \dots a_{2N}]$ によって予測される2つのセグメントの対立遺伝子を決定することが可能である。予測されるシグナルは、 f_x 、 f_y が各対立遺伝子の測定に対する対立遺伝子のセットからのマッピングを記載する仮説 $h_0 : s_{0x} = [f_{0x}(a_{11}, a_{21}) \dots f_{0x}(a_{1N}, a_{2N})]$ 、 $s_{0y} = [f_{0y}(a_{11}, a_{21}) \dots f_{0y}(a_{1N}, a_{2N})]$ について作り出される。 h_0 を仮定すれば、データは $d_{xi} = s_{0xi} + n_{xi}$ 、 $n_{xi} \sim N(0, \sigma_{xi}^2)$ ； $d_{yi} = s_{0yi} + n_{yi}$ 、 $n_{yi} \sim N(0, \sigma_{yi}^2)$ と記載することができる。データおよび参照シグナルを差分することによって測定を作り出す： $m_{xi} = d_{xi} - s_{xi}$ ； $m_{yi} = d_{yi} - s_{yi}$ 。十分な測定ベクトルは $m = [m_x^T m_y^T]^T$ である。

40

【0192】

さて、注目するセグメント、その存在が疑われるセグメントについてのシグナルを作り出し、それを、このセグメントの推定される対立遺伝子に基づいて求める： $a_3 = [a_{31} \dots a_{3N}]$ 。残りについてのシグナルを： $s_r = [s_{rx}^T s_{ry}^T]^T$ と記載し、ここで、 $s_{rx} = [f_{rx}(a_{31}) \dots f_{rx}(a_{3N})]$ 、 $s_{ry} = [f_{ry}(a_{31}) \dots f_{ry}(a_{3N})]$ であり、ここで、もし $a_{3i} = x$ であれば、 $f_{rx}($

50

a_{3i}) = x_i であって、その他の場合は 0 であり、もし $a_{3i} = y$ であれば f_{ry} (a_{3i}) = y_i であり、そうでなければ 0 である。この分析は、遺伝子座 i における対立遺伝子 x の 1 つのコピーの存在がデータ $x_i + m_{x_i}$ を作り出し、遺伝子座 i における対立遺伝子 x の x コピーの存在はデータ $x_i + n_{x_i}$ を作り出すように、測定は線形化されている (後記セクション参照) と仮定する。しかしながら、この仮定は本明細書中に記載された一般的なアプローチでは必要ないことに注意されたし。 h_1 を仮定すれば、もし対立遺伝子 $a_{3i} = x$ であれば、 $m_{x_i} = x_i + n_{x_i}$ 、 $m_{y_i} = n_{y_i}$ であり、もし $a_{3i} = y$ であれば、 $m_{x_i} = n_{x_i}$ 、 $m_{y_i} = y_i + n_{y_i}$ である。その結果、マッチドフィルタ $h = (1/N) R^{-1} s_r$ を作り出すことができ、ここで、 $R = \text{diag}([x_1^2 \dots x_N^2 \quad y_1^2 \dots y_N^2])$ である。測定は $m = h^T d$ である。

10

$$h_0 : m = (1/N) \sum_{i=1}^N s_r x_i n_{x_i} / (x_i^2 + s_{ry_i} n_{y_i} / y_i^2)$$

$$h_1 : m = (1/N) \sum_{i=1}^N s_r x_i (x_i + n_{x_i}) / (x_i^2 + s_{ry_i} (y_i + n_{y_i}) / y_i^2)$$

必要な SNP の数を見積もるためには、全ての対立遺伝子および全ての遺伝子座についての全てのアッセイが同様な特徴を有し、すなわち、 $i = 1 \dots N$ について $x_i = y_i$ および $x_i = y_i =$ であるという単純化仮定を行う。次いで、平均および標準編纂は以下のように見出すことができる

$$h_0 : E(m) = m_0 = 0 ; \quad m | h_0^2 = (1/N^2) (N/2) (2 + 2) = 2 / (N^2) \quad 20$$

$$h_1 : E(m) = m_1 = (1/N) (N/2) (2 + 2) = 2 / N ; \quad m | h_1^2 = (1/N^2) (N) (2 + 2) = 2 / (N^2)$$

さて、 h_1 - 対 - h_0 のこのテストについてシグナル - 対 - ノイズ比率 (SNR) を計算する。シグナルは $m_1 - m_0 = 2 / N$ であって、この測定のノイズの偏差は $m | h_0^2 + m | h_1^2 = 2^2 / (N^2)$ である。その結果、このテストについての SNR は $(4 / N^2) / (2^2 / (N^2)) = N^2 / (2^2)$ である。

【0193】

この SNR を、対立遺伝子要求に基づいてマッチドフィルタリングを行うことなく、遺伝子情報を各遺伝子座において単純に合計するシナリオと比較する。

30

【0194】

【数36】

$$h = (1/N) \bar{1}$$

と仮定し、ここで、

【0195】

【数37】

$$\bar{1}$$

は N のそのベクトルであり、 $i = 1 \dots N$ について $x_i = y_i =$ および $x_i = y_i =$ であると前記したように単純化仮定をする。このシナリオについては、もし $m = h^T d$ であれば：

40

$$h_0 : E(m) = m_0 = 0 ; \quad m | h_0^2 = N^2 / N^2 + N^2 / N^2 = 2^2 / N$$

$$h_1 : E(m) = m_1 = (1/N) (N / 2 + N / 2) = ; \quad m | h_1^2 = (1/N^2) (N^2 + N^2) = 2^2 / N$$

であることを直接的に示すことができる。その結果、このテストについての SNR は $N^2 / (4^2)$ である。言い換えれば、セグメント a_3 について予測される対立遺伝子測定を単に合計するマッチドフィルタを用いることによって必要な SNP の数は 2 倍だけ低下する。これは、各遺伝子座におけるアッセイの異なる効率を説明するためにマッチドフィルタリングを用いることによって達成された SNR 利得を無視する。

【0196】

50

もし参照シグナル s_{x_i} および s_{y_i} を正しく特徴付けなければ、得られた測定シグナル m_{x_i} および m_{y_i} についてのノイズまたは擾乱の SD は増加するであろう。これはもし \ll であれば有意でなく、そうでなければそれは偽検出の確率を増加させるであろう。その結果、この技術は、3つのセグメントが存在し、2つのセグメントは相互の正確なコピーであると推定される仮説をテストするのによく適合している。この場合、 s_{x_i} および s_{y_i} は、他の箇所に記載された定性的対立遺伝子要求に基づくデータ正常化の技術を用いて信頼性よく知られるであろう。1つの実施形態において、方法3は、定性的ゲノタイピングを用い、対立遺伝子ドロップアウトからの定量的測定とは別に、セグメントの第二の正確なコピーの存在を検出することができない方法2と組合せて用いられる。

【0197】

さて、対立遺伝子要求を用いるもう1つの定量的技術に記載する。該方法は、所与の対立遺伝子についての4つの登録の各々におけるシグナルの相対的量を比較することを含む。ホモ接合性増幅が起こる、(または増幅の相対的量が正規化される)、単一の正常な細胞を含む理想化された場合において、4つの可能な状況が起こり得ると想像することができる： (i) ヘテロ接合性対立遺伝子の場合には4つの登録の相対的強度はほぼ1:1:0:0であり、シグナルの絶対的強度は1つの塩基対に対応し； (ii) ホモ接合性対立遺伝子の場合には、相対的強度はほぼ1:0:0:0であり、シグナルの絶対的強度は2つの塩基対に対応する； (iii) ADOが対立遺伝子のうち1つについて起こる対立遺伝子の場合において、相対的強度はほぼ1:0:0:0であって、シグナルの絶対強度は1つの塩基対に対応し；および (ix) ADOが対立遺伝子の双方について起こる対立遺伝子の場合において、相対強度はほぼ0:0:0:0であって、シグナルの絶対的強度は塩基対に対応しないであろう。

【0198】

しかしながら、異数体の場合には、異なる状況が観察されるであろう。例えば、トリソミーの場合には、ADOはなく、3つの状況の1つが起こり： (i) 三重にヘテロ接合性である対立遺伝子の場合には、4つの登録の相対的強度はほぼ1:1:1:0であり、シグナルの絶対的強度は1つの塩基対に対応し； (ii) 対立遺伝子の2つがホモ接合性である場合には相対的強度はほぼ2:1:0:0であり、シグナルの絶対的強度は、各々、2つおよび1つの塩基対に対応し； (iii) 対立遺伝子がホモ接合性である場合には、相対的強度はほぼ1:0:0:0であって、シグナルの絶対的強度は3つの塩基対に対応するであろう。もし対立遺伝子ドロップアウトがトリソミーを持つ細胞における対立遺伝子の場合で起こるならば、正常な細胞で期待される状況のうちの1つが観察されるであろう。モノソミーの場合には、4つの登録の相対的強度はほぼ1:0:0:0であって、シグナルの絶対的強度は1つの塩基対に対応するであろう。この状況は、ADOにおける対立遺伝子の1つが起こった正常な細胞の場合に対応するが、正常な細胞の場合には、これは対立遺伝子のいくらかのパーセンテージで観察されるのに過ぎないであろう。2つの同一の染色体が存在する片親二染色体の場合には、4つの登録の相対的強度はほぼ1:0:0:0であって、シグナルの絶対的強度は2つの塩基対に対応するであろう。1つの親からの2つの異なる染色体が存在するUPDの場合には、この方法は、本特許に記載された他の方法を用いるデータのさらなる分析はこれを明らかにするであろうが、細胞は正常であることを示す。

【0199】

これらの場合の全てにおいて、正常であり、異数体またはUPDを有する細胞いずれかにおいて、1つのSNPからのデータは、細胞の状態について決定するのに適切ではないであろう。しかしながら、もし前記仮説の各々の確率を計算し、それらの確率を所与の染色体上の十分な数のSNPと組み合わせるならば、1つの仮説が支配的であり、高い信頼性をもって染色体の状態を決定することが可能であろう。

【0200】

定量的測定を線形化するための方法

多くのアプローチを採用して、異なる対立遺伝子からのデータを容易に合計し、または

10

20

30

40

50

差分できるように、特定の遺伝子座における遺伝物質の量の測定を線形化することができる。まず、上位概念的なアプローチを議論し、次いで、特定のタイプのアッセイについて設計されるアプローチを議論する。

【0201】

データ d_{x_i} は遺伝子座 i における対立遺伝子 x の遺伝物質の量の非線形測定をいうと仮定する。Nの測定を用いてデータの訓練セットを作り出し、ここに各測定については、データ d_{x_i} に対応する遺伝物質の量は x_i であると見積もられ、またはそのように知られている。この訓練セット $x_i, i = 1 \dots N$ は、現実遭遇するであろう全ての異なる量の遺伝物質にわたるよう選択される。標準回帰技術を用いて、線形測定 $E(x_i)$ を期待して、非線形測定 d_{x_i} からマップされる関数を訓練することができる。例えば、線形回帰を用いて、 c が係数 $c = [c_0 \ c_1 \ \dots \ c_p]^T$ のベクトルである $E(x_i) = [1 \ d_{x_i} \ d_{x_i}^2 \ \dots \ d_{x_i}^p]^T c$ であるように、次元 P の多項関数を訓練することができる。この線形化関数を訓練するために、Nの測定 $x = [x_1 \ \dots \ x_N]^T$ についての遺伝物質の量のベクトル、およびパワーに生じられた測定されたデータのマトリックス $0 \dots P : D = [[1 \ d_{x_1} \ d_{x_1}^2 \ \dots \ d_{x_1}^p]^T [1 \ d_{x_2} \ d_{x_2}^2 \ \dots \ d_{x_2}^p]^T \dots [1 \ d_{x_N} \ d_{x_N}^2 \ \dots \ d_{x_N}^p]^T]^T$ を作り出す。次いで、最小二乗フィット $c = (D^T D)^{-1} D^T x$ を用いて係数を見出すことができる。

10

【0202】

フィットした多項式のような上位概念的関数に依存するよりはむしろ、特定のアッセイの特徴について特殊化された関数を作り出すことができる。例えば、TaqmanアッセイまたはqPCRアッセイを考える。いくつかの閾値と交差する点までの時間の関数としての、対立遺伝子 x およびいくつかの遺伝子座 i についてのダイの量を、 x_i がバイアスオフセットであり、 x_i が指数関数的成長速度であって、 x_i が遺伝物質の量に対応するバイアスオフセット： $g_{x_i}(t) = x_i + x_i \exp(x_i t)$ を持つ指数関数曲線として記載することができる。 x_i の項における測定をキャストするためには、曲線の漸近限界 $g_{x_i}(\infty)$ を探すことによってパラメーター x_i を計算し、次いで、曲線のLOGを取って、 $\log(g_{x_i}(t) - x_i) = \log(x_i) + x_i t$ が得られ、標準的な線形回帰を行うことによって x_i および x_i を見出すことができる。一旦 x_i および x_i についての値を有すれば、もう1つのアプローチは、閾値 g_x がその時点で超過する時間 t_x から x_i を計算することである。 $x_i = (g_x - x_i) \exp(-x_i t_x)$ 。これは、特定の対立遺伝子の遺伝子データの真実の量のノイジーな測定であろう。

20

30

【0203】

どのような技術を用いても、線形化測定を $x_i = x \ x_i + n_{x_i} g_{x_i}(\infty)$ としてモデル化することができ、ここに x は対立遺伝子 x のコピーの数であり、 x_i は対立遺伝子 x および遺伝子座 i についての定数であり、 $n_{x_i} \sim N(0, x^2)$ であり、ここで、 x^2 は経験的に測定することができる。

【0204】

方法4：各遺伝子座における遺伝子データの増幅のための確率分布の関数の使用
特定のSNPについての物質の量は、その上にそのSNPが存在する細胞中の初期セグメントの数に依存するであろう。しかしながら、増幅およびハイブリダイゼーションプロセスのランダムな性質のため、特定のSNPからの遺伝物質の量は、セグメントの出発数に直接的に比例しないであろう。 $q_{s,A}, q_{s,G}, q_{s,T}, q_{s,C}$ が、対立遺伝子を構成する4つの核酸(A, C, T, G)の各々についての特定のSNP s に対する遺伝物質の増幅された量を表すものとする。これらの量は、増幅で用いる技術に依存して、正確にゼロであり得ることを注記する。また、これらの量は、典型的には、特定のハイブリダイゼーションプローブからのシグナルの強度から測定されることも注記する。この強度測定を量の測定の代わりに用いることができるか、あるいは発明の性質を変化させることなく標準的な技術を用いて量の見積もりに変換することができる。 q_s を特定のSN

40

50

Pの全ての対立遺伝子から生じた全ての遺伝物質の合計とする： $q_s = q_{s,A} + q_{s,G} + q_{s,T} + q_{s,C}$ 。NをSNP s を含有する細胞中のセグメントの数とする。Nは典型的には2であるが、0、1または3以上であってよい。議論したいずれの高または中程度スループットのゲノタイピング方法についても、遺伝物質の得られた量は $q_s = (A + A_{\theta,s})N + \theta_s$ として表すことができ、ここで、Aは事前に見積もられたか、または経験的に容易に測定される合計増幅であり、 $A_{\theta,s}$ はSNP s についてのAの見積もりにおける誤差であって、 θ_s はそのSNPについての増幅、ハイブリダイゼーションおよび他のプロセスで導入される相加的ノイズである。ノイズの項 $A_{\theta,s}$ および θ_s は、典型的には、 q_s がNの信頼性がある測定ではないのに十分に大きい。しかしながら、これらのノイズの項の効果は、染色体上の多数のSNPを測定することによって緩和することができる。Sを、染色体21のような特定の染色体上で測定されるSNPの数とする。以下のように、特定の染色体上の全てのSNPにわたる遺伝物質の平均量を得ることが可能である：

10

【0205】

【数38】

$$q = \frac{1}{S} \sum_{s=1}^S q_s = NA + \frac{1}{S} \sum_{s=1}^S A_{\theta,s} N + \theta_s \tag{16}$$

$A_{\theta,s}$ および θ_s は正規分布したランダム変数であり、平均0、および偏差

20

【0206】

【数39】

$$\sigma^2_{A_{\theta,s}} \text{ および } \sigma^2_{\theta_s}$$

であると仮定し、 $q = NA + \dots$ をモデル化することができ、ここで、 \dots は正規分布したランダム変数であり、平均0および偏差

【0207】

【数40】

$$\frac{1}{S} (N^2 \sigma^2_{A_{\theta,s}} + \sigma^2_{\theta})$$

である。その結果、もし十分な数のSNPが、

30

【0208】

【数41】

$$S \gg (N^2 \sigma^2_{A_{\theta,s}} + \sigma^2_{\theta})$$

となるように染色体上で測定されるならば、 $N = q / A$ は正確に見積もることができる。

【0209】

もう1つの実施形態において、増幅は、1つのSNPからのシグナルレベルが $s = a + \dots$ であり、ここで、 $(a + \dots)$ が図14左側の図に似た分布を有するモデルに従うと仮定する。0における関数はおよそ30%の対立遺伝子ドロップアウトの速度をモデル化し、平均はaであり、もし対立遺伝子ドロップアウトがなければ、増幅は0ないし a_0 の均一な分布を有する。この分布の平均の項において、 a_0 は $a_0 = 2.86a$ であることが判明する。さて、図14右側の図面を用いての確率密度関数をモデル化する。 s_c をc遺伝子座から生起するシグナルとし；nをセグメントの数とし； \dots_i を、遺伝子座iからのシグナルに寄与する図14に従って分布したランダム変数とし；および \dots を全ての $\{ \dots_i \}$ についての標準偏差とする。 $s_c = a n c + \dots_{i=1} \dots_{i=n} c \dots_i$ ；平均 $(s_c) = a n c$ ； $std(s_c) = sqrt(nc)$ 。もし \dots を図14左側における分布に従って計算すれば、それは $\dots = 0.907 a^2$ であることが判明する。 $n = s_c / (ac)$ からのセグメントの数を見出すことができ、<5-シグマ統計学>については、 $std(n) < 0.1$ 、従って、 $std(s_c) / (ac) = 0.1 = > 0.95 a \cdot sqrt(nc) / (ac) = 0.1$ 、従って、 $c = 0.95^2 n / 0.1^2 = 181$ を必要とする。

40

50

【 0 2 1 0 】

要求における信頼性を見積もるためのもう1つのモデル、およびどのようにして多くの遺伝子座またはSNPを測定して、所与の程度の信頼性を確保としなければならないかは、相加的ノイズ源、すなわち、 $s = a (1 + \alpha)$ の代わりに増幅のマルチプライア-としてのランダム変数を取り込む。logを取り、 $\log (s) = \log (a) + \log (1 + \alpha)$ となる。さて、新しいランダムな変数 $\alpha = \log (1 + \alpha)$ を作り出し、この変数は、正規分布していると仮定することができる $\sim N (0 , \sigma^2)$ 。このモデルにおいて、増幅は、 α に依存して非常に小さいないし非常に大きいを範囲とすることができるが、決して負ではない。従って、 $s = e^{\alpha} a$ であり；および $s_c = \prod_{i=1}^{cn} a (1 + \alpha_i)$ である。表記方法については、平均 (s_c) および予測値 $E (s_c)$ を相互交換的に用いる。

10

【 0 2 1 1 】

【 数 4 2 】

$$E(s_c) = acn + aE\left(\sum_{i=1...cn} \alpha_i\right) = acn + aE\left(\sum_{i=1...cn} \alpha_i\right) = acn(1 + E(\alpha))$$

$E (\alpha)$ を見出すためには、確率密度関数 (p d f) が、可能である α について見出されなければならない。というのは、 α は公知のガウス p d f を有する α の関数だからである。 $p (\alpha) = p (\gamma) (d \gamma / d \alpha)$ である。従って、

20

【 0 2 1 2 】

【 数 4 3 】

$$p_\gamma(\gamma) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{\gamma^2}{2\sigma^2}} \quad \text{および} \quad \frac{d\gamma}{d\alpha} = \frac{d}{d\alpha}(\log(1+\alpha)) = \frac{1}{1+\alpha} = e^{-\gamma}$$

および：

$$p_\alpha(\alpha) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{\gamma^2}{2\sigma^2}} e^{-\gamma} = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{-(\log(1+\alpha))^2}{2\sigma^2}} \frac{1}{1+\alpha}$$

30

である。これは、 $\alpha = 1$ についての図 1 5 で示された形態を有する。さて、 $E (\alpha)$ は、この p d f にわたって、多数の異なる α について数値的に行うことができる。

【 0 2 1 3 】

【 数 4 4 】

$$E(\alpha) = \int_{-\infty}^{\infty} \alpha p_\alpha(\alpha) d\alpha$$

を積分することによって見出すことができる。これは、 $E (s_c)$ または平均 (s_c) を α の関数として与える。さて、この p d f を用いて、 $\text{var} (s_c)$ を見出すこともできる：

40

【 0 2 1 4 】

【数 4 5】

$$\begin{aligned}
\text{var}(s_c) &= E(s_c - E(s_c))^2 = E\left(\sum_{i=1..cn} a(1 + \alpha_i) - acn - aE\left(\sum_{i=1..cn} \alpha_i\right)\right)^2 \\
&= E\left(\sum_{i=1..cn} a\alpha_i - aE\left(\sum_{i=1..cn} \alpha_i\right)\right)^2 \\
&= a^2 E\left(\sum_{i=1..cn} \alpha_i - cnE(\alpha)\right)^2 \\
&= a^2 E\left(\left(\sum_{i=1..cn} \alpha_i\right)^2 - 2cnE(\alpha)\left(\sum_{i=1..cn} \alpha_i\right) + c^2n^2E(\alpha)^2\right) \\
&= a^2 E\left(cn\alpha^2 + cn(cn-1)E(\alpha_i\alpha_j) - 2cnE(\alpha)\left(\sum_{i=1..cn} \alpha_i\right) + c^2n^2E(\alpha)^2\right) \\
&= a^2c^2n^2\left(E(\alpha^2) + (cn-1)E(\alpha_i\alpha_j) - 2cnE(\alpha)^2 + cnE(\alpha)^2\right) \\
&= a^2c^2n^2\left(E(\alpha^2) + (cn-1)E(\alpha_i\alpha_j) - cnE(\alpha)^2\right)
\end{aligned}$$

ここで、これは多数の異なる α について $p(\alpha)$ を用いて数值的に解いて、 $p(\alpha)$ の関数として $\text{var}(s_c)$ を得ることもできる。次いで、公知の数の遺伝子座 c および公知の数のセグメント n を持つ試料から一連の測定を取ることができ、このデータから $\text{std}(s_c) / E(s_c)$ を見出すことができる。それにより、 α についての値を計算することが可能となる。 n を見積もるためには、 $E(s_c) = nac(1 + E(\alpha))$ 、従って、

【0 2 1 5】

【数 4 6】

$$\hat{n} = \frac{s_c}{ac(1 + E(\alpha))}$$

を、

【0 2 1 6】

【数 4 7】

$$\text{std}(\hat{n}) = \frac{\text{std}(s_c)}{ac(1 + E(\alpha))} \text{std}(n)$$

となるように測定することができる。0 ないし平均の十分に多数の独立したランダムな変数を合計すると、分布はガウス形態に近づき、かくして、 s_c (および

【0 2 1 7】

【数 4 8】

 \hat{n}

) は正規分布しているとして処理することができ、前記したように、5 - 統計学：

【0 2 1 8】

【数 4 9】

$$\text{std}(\hat{n}) = \frac{\text{std}(s_c)}{ac(1 + E(\alpha))} < 0.1$$

を用いて、 $2 \text{normcdf}(5, 0, 1) = 2.7e-7$ の誤差確率を有するようになることができる。これより、遺伝子座 c の数について解くことができる。

【0 2 1 9】

雌雄鑑別

システムの 1 つの実施形態において、遺伝子データを用いて、標的個体の性別を決定することができる。本明細書中に開示した該方法を用いて、親からのいずれの染色体のいずれのセグメントが標的の遺伝物質に貢献したかを決定した後、性染色体のいずれが父親から遺伝したかを見るためにチェックすることによって標的の性別を決定することができる：X は女性を示し、および Y は男性を示す。この方法をどのようにして用いて、標的の性

別を決定するかは当業者に明らかなはずである。

【 0 2 2 0 】

仮説の確証

システムのいくつかの実施形態において、1つの決定は、最高の可能な信頼性をもって正しい遺伝子状態の予測を行うためには、各可能な状態について仮説を立てる必要があることである。しかしながら、遺伝子状態の可能な数が指数関数的に大きくなり、計算時間が制限されるにつれ、各仮説を検定するのは合理的でないであろう。これらの場合において、別のアプローチは、仮説確証の概念を用いることである。これは、ある値、値のセット、もしある仮説、または仮説のクラスが真実であるならば測定されたデータにおいて観察されることが期待される特性またはパターンに対する制限を見積もることを含む。次いで、測定された値を検定して、それらが予測された制限に入るか、および/またはある予測された特性またはパターンを検定することができるか、および予測が適合しないかを見ることができ、次いで、アルゴリズムはさらなる調査のための測定に警告を与えることができる。

10

【 0 2 2 1 】

例えば、染色体の1つのアームの端部が標的DNAにおいて破壊されている場合、最もありそうな仮説は(例えば、「異数体」とは反対に)「正常である」と計算することができる。これは、遺伝物質の真の状態に対応し、すなわち、染色体の1つの端部が破壊された特定の仮説は、その状態の尤度が非常に低いので検定されていないからである。もし確証の概念を用いれば、アルゴリズムは、多数の値、染色体の破壊されたセクションに存在する対立遺伝子に対応するものは、測定の期待された限界の外にあることを注記するであろう。フラグが生起され、この場合についてのさらなる調査を促し、遺伝物質の真の状態が発見される尤度を増大させる。

20

【 0 2 2 2 】

どのようにして、開示された方法を修飾して、確証技術を含ませるかは当業者に明らかなはずである。開示された方法を用いて検出するのは非常に困難であると予測される1つの異常は、バランスしたトランスロケーションであることを注記する。

【 0 2 2 3 】

汚染されたDNAでの方法の適用

システムの1つの実施形態において、外来性DNAで明確にまたは可能性として汚染された標的DNAからの遺伝子データもまた、開示された方法を用いて正常化することができる。先に概説した概念、仮説確証のそれを用いて、予測される限界の外になる遺伝子試料を同定することができ；汚染された試料の場合には、この確証は警告を生起させ、試料を汚染したものとして同定できると予測される。

30

【 0 2 2 4 】

標的DNAの大きなセグメントは親遺伝子データから知られておるので、かつ汚染の程度は十分に低く、十分なSNPが測定されるものとする、外来性遺伝物質による誤ったデータを同定しかねない。本明細書中に開示された方法は、依然として、より低い信頼性のレベルに拘わらず、標的ゲノムの再構築を可能とするはずである。汚染のレベルが十分に低いものとするれば、最もありそうであると計算される仮説は、依然として、標的DNA試料中の遺伝物質の真の状態に対応すると予測される。

40

【 0 2 2 5 】

どのようにして、外来性DNAにより誤ったシグナルで汚染された遺伝子データを清浄化する目的でこれらの方法を最適化するかは当業者に明らかなはずである。

【 0 2 2 6 】

実施例

システムの1つの実施形態において、前記した方法は、関連SNPのリストにおける各SNPの最もありそうな同一性、ならびに各SNP要求についての信頼性レベルを計算するアルゴリズムのセットを用いて実行することができる。本明細書中に記載するのは、この特許に開示した方法を実行するための1つの可能な方法である。図16および図17は

50

、開示された方法のこの実施を頓挫、入力要件、および出力のフォーマットを実質的に表す。

【0227】

図16は入力データ(1601)およびそのフォーマットおよび要件、ならびに出力データ(1605)およびそのフォーマットに焦点を当てる。アルゴリズムへの入力は、ユーザーによる入力を含めた測定されたデータ(1602)、および結果的には新しく収集されたデータによって更新されるデータベースに保存された現存データ(1603)よりなる。測定されたデータ(MD, 1602)は胚、および父性および母性対立遺伝子についての所望のSNPについて測定された遺伝子データ、ならびに対立遺伝子の各々が知られている精度または信頼性よりなる。現存データ(1603)は集団頻度データ(FD)、測定バイアスデータ(BD)および交差データ(CD)よりなる。

10

【0228】

集団頻度データ(FD)は利用可能なSNPの各々について(値A、C、T、Gの各々についての)対立遺伝子頻度を含有する。これらのデータは従前に知られているか、または測定することができ、本書類中の他の箇所に記載されたように新しく収集されたデータで更新することができる。

【0229】

測定バイアスデータ(BD)は、ある種の値に向けての測定プロセスのバイアスを捕獲する。例えば、対立遺伝子の真の値が $X = A$ であって、正しい測定の確率は p_x であると仮定し、測定された値 x の分布は：

20

【0230】

【数50】

	A	C	T	G
確率	p_x	p_c	p_t	p_g
バイアスの無い確率	p_x	$(1-p_x)/3$	$(1-p_x)/3$	$(1-p_x)/3$

であり、ここで、 $p_x + p_c + p_t + p_g = 1$ である。もし値のいずれかに向けての測定のバイアスがなければ、 $p_c = p_t = p_g = (1 - p_x) / 3$ である。この情報は、測定プロセスのメカニズムおよび関連機器についての経験的および理論的知識から区別することができる。

30

【0231】

交差データ(CD)は、HAPMAPデータから収集された、スニップの対の間の遺伝子距離および交差確率のデータベースよりなる。

【0232】

一緒にすると、(MD)、(FD)、(BD)、(CD)は、開示された方法(「親サポート」, 1604という)アルゴリズムに対する必要な入力をなす。次いで、このアルゴリズム(1604)を入力データとして操作して、出力データ(1605)を生じさせ、これは測定値を仮定した標的の遺伝子データのもっともありそうな「真の」値、ならびに親対立遺伝子に関する各SNPの最もありそうな起源を記載する。

40

【0233】

図17は(「親サポート」という)アルゴリズムそれ自体の構造、およびどのようにしてこれらの入力データの各々がアルゴリズムによって利用されるかに焦点を当てる。逆に作業し、最もありそうな仮説を見出すためには、全ての可能な仮説Hについての、測定を仮定した仮説の確率 $P(H | M)$ 1707を計算する必要がある。

先に記載したように：

【0234】

【数51】

$$P(H|M) = \frac{P(M|H)}{P(M)} P(H), \quad P(M) = \sum_{h \in S_H} P(M|h)P(h)$$

である。

$P(H|M)$ (1710)を見出すためには、全ての仮説Hについて、 $P(M|H)$ (1707)および $P(H)$ (1708)を見出すことがまず必要である。これは、先に示した方程式による $P(M)$, 1709の計算を可能とする。仮説の確率 $P(H)$ (1708)は、先に説明したようにどれくらい多くの交差が推定されるか、およびこれらの交差の各々の尤度(CD, 1704)に依存する。

10

【0235】

$P(M|H)$ は、先に説明したように、以下の方程式

【0236】

【数52】

$$P(M|H) = \sum_t P(M|H \& t)P(t)$$

を用いて計算することができる。

【0237】

$P(T)$, 1706は父性および母性対立遺伝子についての特定の値tの頻度であり、集団頻度データ(FD, 1703)に由来する。 $P(M|H \& t)$, 1705は、特定の「真の」値tを仮定し、胚、父親および母親の対立遺伝子を正しく測定する確率である。ユーザーによってエンターされた測定データおよび精度(MD, 1701)、および測定バイアスデータベース(BD, 1702)は、 $P(M|H \& t)$, 1705を計算するのに必要な出力である。

20

【0238】

該方法のより詳細な記載を以下に掲げる。id S_1, \dots, S_k で同定される、kのSNPについての、 $SNP \ R = \{r_1, \dots, r_k\}$ 、(kのSNPのセット)、および親および胚の対応する測定された同一性 $M = (e_1, e_2, p_1, p_2, m_1, m_2)$ で開始し、ここに：

30

$e_1 = (e_{11}, e_{12}, \dots, e_{1k})$ は、全てのSNPについての、胚の染色体の1つでの測定であり(それらは、全てが、同一親染色体に由来する必要はない)、

$e_2 = (e_{21}, e_{22}, \dots, e_{2k})$ は胚の他の染色体での測定であり、

$p_1 = (p_{11}, p_{12}, \dots, p_{1k})$ は、(全て同一染色体に由来する)父親の第一の染色体での測定であり、

$p_2 = (p_{21}, p_{22}, \dots, p_{2k})$ は、(全て同一染色体に由来する)父親の第二の染色体での測定であり、

$m_1 = (m_{11}, m_{12}, \dots, m_{1k})$ は、(全て同一染色体に由来する)母親の第一の染色体での測定であり、

$m_2 = (m_{21}, m_{22}, \dots, m_{2k})$ は、(全て同一染色体に由来する)母親の第二の染色体での測定である。

40

【0239】

また、 $M = (M_1, \dots, M_k)$ を書くことができ、ここで、 $M_i = (e_{1i}, e_{2i}, p_{1i}, p_{2i})$ である。

【0240】

該方法の目標は、「真の」胚値 $T = (E_1, E_2)$ 、すなわち、測定Mを仮定した最もありそうな場合を決定することであり、ここに：

$E_1 = (E_{11}, E_{12}, \dots, E_{1k})$ は、父性染色体に対応する胚の第一の染色体での測定、 $E_{1i} = \{p_{1i}, p_{2i}\}$ であり、

$E_2 = (E_{21}, E_{22}, \dots, E_{2k})$ は、母性値に対応する胚の第二の染色体で

50

の測定、 $E_{2i} = \{m_{1i}, m_{2i}\}$ である。

【0241】

また、 $T = \{T_1, \dots, T_k\}$ を書くことができ、ここで、 $T_i = (E_{1i}, E_{2i})$ である。

【0242】

効果的には、親染色体値 (p_1, p_2, m_1, m_2) は、 (E_1, E_2) の測定された値をチェックし、確認し、および修正するためのサポートとして用いる、よって、用語「親サポートアルゴリズム」。

【0243】

この目標を達成するためには、胚値の起源についての全ての可能な仮説を開発し、測定Mを仮定して最もありそうなものを選択する。仮説空間は $S_H = \{H^1, \dots, H^q\} = \{\text{全ての仮説のセット}\}$ であり、ここで、各仮説はフォーマット $H^j = (H^j_1, \dots, H^j_k)$ のものであり、ここで、 H^j_1 は、 $p_i^* = \{p_{1i}, p_{2i}\}$ および $m_i^* = \{m_{1i}, m_{2i}\}$ であるフォーマット $H^j_1 = (p_i^*, m_i^*)$ の、SNP i についての「ミニ」仮説である。4つの異なる「ミニ」仮説 H^j_1 、特に：

$H^j_{i1} : (e_{1i}, e_{2i}) = \{(p_{1i}, m_{1i}) \text{ または } (m_{1i}, p_{1i})\}$

$H^j_{i2} : (e_{1i}, e_{2i}) = \{(p_{1i}, m_{2i}) \text{ または } (m_{2i}, p_{1i})\}$

$H^j_{i3} : (e_{1i}, e_{2i}) = \{(p_{2i}, m_{1i}) \text{ または } (m_{1i}, p_{2i})\}$

$H^j_{i4} : (e_{1i}, e_{2i}) = \{(p_{2i}, m_{2i}) \text{ または } (m_{2i}, p_{2i})\}$

がある。

【0244】

理論において、 S^H は $q = 4^k$ の異なるメンバーを有して、ピックアップすることができるが、後に、この空間は父性および母性染色体の最大数の交差で限定されるであろう。

【0245】

最もありそうな仮説 H^* は：

【0246】

【数53】

$$H^* = \arg \max_{H \in S_H} P(H | M)$$

であると選択される。特定のHについては：

【0247】

【数54】

$$P(H | M) = \frac{P(M | H)}{\sum_{h \in S_H} P(M | h) P(h)} P(H)$$

である。

各仮説についてのそのような由来：

(1) $P(M | H)$ は、特定の仮説Hを仮定した測定Mの確率である。

(2) $P(H)$ は特定の仮説Hの確率である。

(3) $P(M)$ は測定Mの確率である。

全てのHについて $P(H | M)$ を導いた後、最大の確率を持つものを選択する。

【0248】

$P(M | H)$ の誘導

各SNPについての測定は、全てのkのSNPでの、 $M = (M_1, \dots, M_k)$ および特定の仮説 $H = (H_1, \dots, H_k)$ について独立しているので： $P(M | H) = P(M_1 | H_1) \dots P(M_k | H_k)$ である。特定のSNP r では、 $P(M_r | H_r)$ を誘導する。 $X = \{A, C, T, G\} \times \{A, C, T, G\} \times \{A, C, T, G\} \times \{A, C, T, G\}$ については、ベイズ式による「真の親値 $(P_{1r}, P_{2r}, M_{1r}, M_{2r})$ 」についての全ての可能な空間は：

10

20

30

40

50

【 0 2 4 9 】

【 数 5 5 】

$$P(M_r / H_r) = \sum_{t \in \Omega} P(M_r / H_r \& (P_{1r}, P_{2r}, M_{1r}, M_{2r}) = t) * P((P_{1r}, P_{2r}, M_{1r}, M_{2r}) = t)$$

である。

【 0 2 5 0 】

P (M_r | H_r & (P_{1r} , P_{2r} , M_{1r} , M_{2r}) = t) の誘導

M_r = (e_{1r} , e_{2r} , p_{1r} , p_{2r} , m_{1e} , m_{2r}) はこの SNP での所与の測定である。

10

【 0 2 5 1 】

T = (E_{1r} , E_{2r} , P_{1r} , P_{2r} , M_{1r} , M_{2r}) は、仮説による T から固定された t = (P_{1r} , P_{2r} , M_{1r} , M_{2r}) および (E_{1r} , E_{2r}) での推定された「真の」値である。(E_{1r} は P_{1r}、P_{2r} の一方であり、E_{2r} は M_{1r}、M_{2r} の一方である)。

$$P (M_r = (e_{1r} , e_{2r} , p_{1r} , p_{2r} , m_{1r} , m_{2r}) / T = (E_{1r} , E_{2r} , P_{1r} , P_{2r} , M_{1r} , M_{2r})) =$$

$$P (e_{1r} / E_{1r}) * P (e_{2r} / E_{2r}) * P (p_{1r} / P_{1r}) * P (p_{2r} / P_{2r}) * P (m_{1r} / M_{1r}) * P (m_{2r} / M_{2r})$$

p_{e_{1r}} = P (SNP_r についての胚値 i を正確に測定)

20

p_{p_{1r}} = P (SNP_r についての父親値 i を正確に測定)

p_{m_{1r}} = P (SNP_r についての母親値 i を正確に測定) とすれば、

【 0 2 5 2 】

【 数 5 6 】

$$P(e_{1r} / E_{1r}) = \begin{cases} p_{er1} & e_{1r} = E_{1r} \\ (1-p_{er1}) * p(e_{1r}, E_{1r}, r) & e_{1r} \neq E_{1r} \end{cases}$$

$$= I_{e_{1r}=E_{1r}} * p_{er1} + I_{e_{1r} \neq E_{1r}} * (1-p_{er1}) * p(e_{1r}, E_{1r}, r) = F(e_{1r}, E_{1r}, p_{er1}, r)$$

であり、ここで、測定バイアスがなければ、p (e_{1r} , E_{1r} , r) = 1 / 3 であり、そうでなければ、それは H a p m a p プロジェクトからのデータのような実験データから決定することができる。

30

【 0 2 5 3 】

P ((P_{1r} , P_{2r} , M_{1r} , M_{2r}) = t) の誘導

t = (t₁ , t₂ , t₃ , t₄) については：

P ((P_{1r} , P_{2r} , M_{1r} , M_{2r}) = (t₁ , t₂ , t₃ , t₄)) = P (P_{1r} = t₁) * P (P_{2r} = t₂) * P (M_{1r} = t₃) * P (M_{2r} = t₄) である。(P₁ , P₂ , M₁ , M₂) の n の試料があると仮定し、全ての父性および母性値は独立しており、{ A , C , T , G } における t_i については t = (t₁ , t₂ , t₃ , t₄) であると推定される。

40

【 0 2 5 4 】

t₁ = A について特定の p_{1A} = P (P₁ = t₁) を得るためには、いずれものデータの不存在下において、この確率は 0 および 1 の間の何かであり得ると推定され、従って、それは U (0 , 1) の値が割り当てられる。データの獲得に関しては、これは新しい値で更新され、このパラメーターの分布はベータ分布となる。P₁ の n の観察のうち、h の値 P₁ = A、および w = (事象 P₁ = A) および D = (所与のデータ) がある。先のセクションにおいて、p (w | データ) について = h + 1、 = n - h + 1 での 分布 B (,) の形式が記載されている (方程式 (8) 参照)。予測された値および X ~ B (,) 分布の偏差は：

【 0 2 5 5 】

50

【数57】

$$EX = \frac{\alpha}{\alpha + \beta}$$

$$VX = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$

である。従って、パラメーターの事後平均値は $p_{1rA} = P(P_{1r} = A | Data) = (h + 1) / (n + 2)$ である。同様に、 $p_{1rB} = (\#(p_{1r} = B) + 1) / (n + 2)$ 、 \dots 、 $m_{2rG} = (\#(m_{2r} = G) + 1) / (n + 2)$ などである。かくして、全ての値 p_{1rA}, \dots, m_{2rG} が導かれ：

10

【0256】

【数58】

$$P((P_{1r}, P_{2r}, M_{1r}, M_{2r}) = (t_1, t_2, t_3, t_4)) = p_{1rt_1} * p_{2rt_2} * m_{1rt_3} * m_{2rt_4}$$

である。

【0257】

P(H)の誘導

$H_i = (p_i^*, m_i^*)$ での仮説 $H = (H_1, \dots, H_k)$ の確率は、染色体交差の量に依存する。例えば、

$P(\text{交差}) = 0$ であれば、もし $\{(p_{11}, p_{21}, \dots, p_{s1}), (p_{12}, p_{22}, \dots, p_{s2})\}$ における p^* 、 $\{(m_{11}, m_{21}, \dots, m_{s1}), (m_{12}, m_{22}, \dots, m_{s2})\}$ における m^* であれば、 $P(H = 1/4)$ であって、 $H = (p^*, m^*)$ であり、そうでなければ0であり

20

$P(\text{交差}) > 0$ であれば、各SNPの間の交差の確率を一体化させるのは重要である。

【0258】

仮説Hは、各SNPについての父性および母性染色体についての仮説、独立している、 $p_i^* \{p_{1i}, p_{2i}\}$ および $m_i^* \{m_{1i}, m_{2i}\}$ 、すなわち、 $H = (H_p, H_m)$ よりなり、ここで、 $H_p = (p_1^*, \dots, p_k^*)$ および $H_m = (m_1^*, \dots, m_k^*)$ である。

$P(H = P(H_p) * P(H_m))$ 。SNPはロケーションを増大させることによって秩序化され、

30

【0259】

【数59】

$$P(H_p) = \frac{1}{4} \prod_{i=2}^k ((PC_i * (1 - I_i) + (1 - PC_i) * I_i)$$

であると仮定し、ここで、 $PC_i = P(\text{交差}(r_{i-1}, r_i))$ 、すなわち、SNP r_{i-1}, r_i の間のどれかの交差の確率であり、もし p_i^*, p_{i-1}^* が共に p_1 または p_2 に由来するならば、 $I_i = 1$ であり、そうでなければそれは0である。

【0260】

40

P(交差(a, b))の誘導

(塩基で与えた)塩基ロケーション $1_a, 1_b$ におけるSNP a, b を仮定すれば、交差の確率は：

$$P(1_a, 1_b) = 0.5 (1 - \exp(-2G(1_a, 1_b)))$$

として近似され、ここで、 $G(1_a, 1_b) = \text{ロケーション } 1_a, 1_b \text{ の間のモルガンで表した遺伝子距離}$ 。Gについての正確な閉じた形態の関数はないが、それは $G(1_a, 1_b) = |1_a - 1_b| * 1e^{-8}$ として緩く見積もられる。良好な近似は、全てのロケーションにわたってのiスパンニングについての、塩基ロケーション s_i および距離 $G(s_i, s_{i+1})$ のHapMapデータベースを利用することによって用いることができる。特に、

50

【 0 2 6 1 】

【 数 6 0 】

$$G(l_a, l_b) = \sum_{l_a < s_i < l_b} G(s_i, s_{i+1})$$

であり、従って、それは交差確率で用いることができる。

【 0 2 6 2 】

P (M) の誘導

一旦 P (M | H) が知られていれば、 P (H) は S_H における全ての異なる H について見出すことができる。

10

【 0 2 6 3 】

【 数 6 1 】

$$P(M) = \sum_{H \in S_H} P(M | H) P(H)$$

。

【 0 2 6 4 】

最大確率の仮説を導くためのより便宜な方法

コンピュータ時間の制限、および前記した方法の複雑性の指数関数スケーリングを仮定すれば、SNP の数が増加するにつれ、ある場合には、より便宜な方法を用いて、最大確率の仮説を決定し、かくして、関連する SNP 要求をなすのが必要であろう。これを達成するためのより迅速な方法は以下の通りであり：

20

以前より：

$P(H | M) = P(M | H) * P(H) / P(M)$ 、 $\arg \max_H P(H | M) = \arg \max_H$ および $P(M | H) * P(H) = \arg \max_H F(M, H)$ であり、目的は $F(M, H)$ を最大化する H を見出すことである。

【 0 2 6 5 】

$M_{(s, k)}$ = スニップ s ないし k についての測定、 $H_{(s, k)}$ = スニップ s ないし k についての仮説、および短いものについて、 $M_{(k, k)} = M_k$ を仮定すれば、 $H_{(k, k)} = H_k$ = スニップ k についての測定および仮説である。先に示したように：

30

【 0 2 6 6 】

【 数 6 2 】

$$P(M_{(1,k)} | H_{(1,k)}) = \prod_{i=1}^k P(M_i | H_i) = P(M_k | H_k) * \prod_{i=1}^{k-1} P(M_i | H_i) = P(M_k | H_k) * P(M_{(1,k-1)} | H_{(1,k-1)})$$

であり、また、

【 0 2 6 7 】

【 数 6 3 】

$$P(H_{(1,k)}) = 1/4 * \prod_{i=2}^k PF(H_{i-1}, H_i) = PF(H_{k-1}, H_k) * 1/4 * \prod_{i=2}^{k-1} PF(H_{i-1}, H_i) = PF(H_{k-1}, H_k) * P(H_{(1,k-1)})$$

40

であり、ここで、

【 0 2 6 8 】

【 数 6 4 】

$$PF(H_{i-1}, H_i) = \begin{cases} 1 - PC(H_{i-1}, H_i) & H_{i-1} = H_i \\ PC(H_{i-1}, H_i) & H_{i-1} \neq H_i \end{cases}$$

および $PC(H_{i-1}, H_i) = H_{i-1}, H_i$ の間の交差の確率。

【 0 2 6 9 】

従って、最後には、n のスニップについては：

$$F(M, H) = P(M | H) * P(H) = P(M_{(1, n)}, H_{(1, n)}) * P(H_{(1, n)})$$

50

$$= P(M_{(1,n-1)}, H_{(1,n-1)}) * P(H_{(1,n-1)}) * P(M_n | H_n) * PF(H_{n-1}, H_n)$$

であり、従って： $F(M, H) = F(M_{(1,n)}, H_{(1,n)}) = F(M_{(1,n-1)}, H_{(1,n-1)}) * P(M_n | H_n) * PF(H_{n-1}, H_n)$ である。かくして、 n のスニップについての計算を $n-1$ スニップについての計算に代えることが可能である。

n についてのスニップの n についての $H = (H_1, \dots, H_n)$ 仮説では：

【 0 2 7 0 】

【 数 6 5 】

$$\max_H F(M, H) = \max_{(H_{(1,n-1)}, H_n)} F(M, (H_{(1,n-1)}, H_n)) = \max_{H_n} \max_{H_{(1,n-1)}} F(M, (H_{(1,n-1)}, H_n)) = \max_{H_n} G(M_{(1,n)}, H_n)$$

10

であり、ここに

【 0 2 7 1 】

【 数 6 6 】

$$G(M_{(1,n)}, H_n) = \max_{H_{(1,n-1)}} F(M_{(1,n)}, (H_{(1,n-1)}, H_n)) =$$

$$\max_{H_{(1,n-1)}} F(M_{(1,n-1)}, H_{(1,n-1)}) * P(M_n | H_n) * PF(H_{n-1}, H_n) =$$

$$P(M_n | H_n) * \max_{H_{(1,n-1)}} F(M_{(1,n-1)}, H_{(1,n-1)}) * PF(H_{n-1}, H_n) =$$

$$P(M_n | H_n) * \max_{H_{n-1}} \max_{H_{(1,n-2)}} F(M_{(1,n-1)}, (H_{(1,n-2)}, H_{n-1})) * PF(H_{n-1}, H_n) =$$

$$P(M_n | H_n) * \max_{H_{n-1}} PF(H_{n-1}, H_n) * G(M_{(1,n-1)}, H_{n-1})$$

20

である。まとめると：

【 0 2 7 2 】

【 数 6 7 】

$$\max_H F(M, H) = \max_{H_n} G(M_{(1,n)}, H_n)$$

30

であり、ここで、 G は帰納的に見出すことができ： $i = 2, \dots, n$ については、

【 0 2 7 3 】

【 数 6 8 】

$$G(M_{(1,n)}, H_n) = P(M_n | H_n) * \max_{H_{n-1}} [PF(H_{n-1}, H_n) * G(M_{(1,n-1)}, H_{n-1})]$$

$$\text{および } G(M_{(1,1)}, H_1) = 0.25 * P(M_1 | H_1)$$

である。

【 0 2 7 4 】

最良の仮説は以下のアルゴリズムに従って見出すことができる：

40

工程 1： $I = 1$ では、 H_1 についての 4 つの仮説を作り出し、これらの各々についての $G(M_1 | H_1)$ を作り、 G_1, G_2, G_3, G_4 を覚える。

工程 2： $I = 2$ では： H_2 についての 4 つの仮説を作り出し、前記式を用いて $G(M_{(1,2)} | H_2)$ を作成し：

【 0 2 7 5 】

【 数 6 9 】

$$G(M_{(1,2)}, H_2) = P(M_2 | H_2) * \max_{H_1} [PF(H_1, H_2) * G(M_1, H_1)]$$

これらの新しい 4 つの G_n を覚える。

$k = n$ まで、 $k_i = k_{i-1} + 1$ にて $I = k$ につき工程 2 を反復し： H_k について 4 つの

50

仮説を作り出し、 $G(M_{(1,k)} | H_k)$

【0276】

【数70】

$$G(M_{(1,k)}, H_k) = P(M_k | H_k) * \max_{H_{k-1}} [PF(H_{k-1}, H_k) * G(M_{(1,k-1)}, H_{k-1})]$$

を作成し、これらの4つの G_n を覚える。

【0277】

いずれかの時点において覚えるべきただ4つの仮説、および一定数の操作があるので、アルゴリズムは線形である。

【0278】

$P(M) : P(H | M) = P(M | H) * P(H) / P(M) = F(M, H) / P(M)$)) を見出すために、前記したように：

【0279】

【数71】

$$\begin{aligned} P(M) &= P(M_{(1,k)}) = \sum_{H_{(1,k)}} P(M_{(1,k)} | H_{(1,k)}) * P(H_{(1,k)}) \\ &= \sum_{H_k} P(M_K | H_k) \sum_{H_{(1,k-1)}} P(M_{(1,k-1)} | H_{(1,k-1)}) * P(H_{(1,k-1)}) * PF(H_{k-1}, H_k) \\ &= \sum_{H_k} P(M_K | H_k) * W(M_{(1,k-1)} | H_k) \end{aligned}$$

10

20

であり、ここで、

【0280】

【数72】

$$W(M_{(1,k-1)} | H_k) = \sum_{H_{(1,k-1)}} P(M_{(1,k-1)} | H_{(1,k-1)}) * P(H_{(1,k-1)}) * PF(H_{k-1}, H_k)$$

である。

$W(M, H)$ は帰納を用いることによって解くことができる：

【0281】

【数73】

$$\begin{aligned} W(M_{(1,k-1)} | H_k) &= \sum_{H_{(1,k-1)}} P(M_{(1,k-1)} | H_{(1,k-1)}) * P(H_{(1,k-1)}) * PF(H_{k-1}, H_k) \\ &= \sum_{H_{k-1}} P(M_{k-1} | H_{k-1}) \sum_{H_{(1,k-2)}} P(M_{(1,k-2)} | H_{(1,k-2)}) * P(H_{(1,k-2)}) * PF(H_{k-2}, H_{k-1}) * PF(H_{k-1}, H_k) \\ &= \sum_{H_{k-1}} P(M_{k-1} | H_{k-1}) * PF(H_{k-1}, H_k) * W(M_{(1,k-2)} | H_{k-1}) \end{aligned}$$

30

従って：

【0282】

【数74】

$$W(M_{(1,k-1)} | H_k) = \sum_{H_{k-1}} P(M_{k-1} | H_{k-1}) * PF(H_{k-1}, H_k) * W(M_{(1,k-2)} | H_{k-1})$$

40

$$\text{および } W(M_{(1,1)} | H_2) = \sum_{H_1} P(M_1 | H_1) * 0.25 * PF(H_1, H_2)$$

である。

【0283】

アルゴリズムは前記の場合に同様であり、ここで、 $i = 2 : n$ であって、各工程において、 $W(i)$ の新しいセットを、最終工程が最適化された W を生じるまで作り出される。

【0284】

50

d_1 、 d_2 、 h 、 p_{d_1} 、 p_{d_2} 、 p_h からの p_1 、 p_2 、 p_{p_1} 、 p_{p_2} 値の誘導説明の目的で、このセクションは父親のジブroidおよびハブroidデータに焦点を合わせるが、同一アルゴリズムを母親に適用することができることに注意するのは重要である。

- d_1 、 d_2 - ジブroid測定での対立遺伝子要求
 - h - ハブroid測定についての対立遺伝子要求
 - p_{d_1} 、 p_{d_2} - ジブroid測定の各々についての正しい対立遺伝子要求の確率
 - p_h - ハブroid測定についての正しい対立遺伝子要求の確率
- これらのデータは開示されたアルゴリズムについての以下の入力パラメーターにマッピングすべきである：

- p_1 - ハブroid細胞および1つのジブroid細胞に対応する対立遺伝子
 - p_2 - 残りのジブroid細胞に対応する対立遺伝子
 - p_{p_1} 、 p_{p_2} - 正しい対立遺伝子要求の確率
- h は d_1 に対応するので、 p_1 の値を見出すためには、 h および d_1 を用いる必要がある。次いで、 p_2 は自動的に d_2 に対応する。同様に、もし h が d_2 に対応すれば、 p_1 の値を見出すためには、 h および d_2 を用いる必要があり、次いで、 p_2 は d_1 に対応するであろう。

【0285】

用語「対応する」を用いる。というのは、それは、異なる測定結果および集団頻度に依存して、「等しい」または「より高い確率で由来する」を意味することができるからである。

【0286】

アルゴリズムの目標は、生の測定 h 、 d_1 、 d_2 、 p_h 、 p_{d_1} 、 p_{d_2} および集団頻度の結果を超えて隠された「真の」対立遺伝子値の確率を計算することである。

【0287】

基本的なアルゴリズム工程は以下の通りである：

- (i) h 、 d_1 、 d_2 、 p_h 、 p_{d_1} 、 p_{d_2} 値、および集団頻度データに基づいて h が d_1 または d_2 に対応するかを決定する、
- (ii) 対立遺伝子要求を p_1 および p_2 に帰属させ；工程(1)に基づいて確率 p_{p_1} および p_{p_2} を計算する。

【0288】

h の d_1 または d_2 への帰属

2つの仮説：

H_1 ： h は d_1 に対応する（ h は d_1 に由来する）

H_2 ： h は d_2 に対応する（ h は d_2 に由来する）

を確率する。仕事は、測定 M ：を仮定してこれらの2つの仮説の確率を計算することである：

$P(H_1 / M(h, d_1, d_2, p_h, p_{d_1}, p_{d_2}))$ および $P(H_2 / M(h, d_1, d_2, p_h, p_{d_1}, p_{d_2}))$ 。

以下、(テキストを単純化するために、これらを $P(H_1 / M)$ および $P(H_2 / M)$)という)。

【0289】

これらの確率を計算するために、ベイズ則：

【0290】

【数75】

$$P(H_1 | M) = \frac{P(M | H_1) * P(H_1)}{P(M)}; P(H_2 | M) = \frac{P(M | H_2) * P(H_2)}{P(M)}$$

を適用し、ここで、 $P(M) = P(M / H_1) * P(H_1) + P(M / H_2) * P(H_2)$ である。仮説 H_1 および H_2 は同等にありそうなので、 $P(H_1) = P(H_2) = 0$ 。

10

20

30

40

50

5 であり、従って：

【 0 2 9 1 】

【 数 7 6 】

$$P(H_1 | M) = \frac{P(M | H_1)}{P(M | H_1) + P(M | H_2)} \text{ および } P(H_2 | M) = \frac{P(M | H_2)}{P(M | H_1) + P(M | H_2)}$$

である。

【 0 2 9 2 】

P (M / H ₁) および P (M / H ₂) を計算するためには、ジブroid結果 d ₁ および d ₂ の全ての可能な値のセット、 Ω = { A A , A C , . . . , G G }、すなわち、A、C、T、Gのいずれかの組合せ、いわゆる基礎となる状態を考慮しなければならない。仮説を基礎となる状態に適用する場合（すなわち、仮説 H ₁ または H ₂ に基づいて h の推定値を値 d ₁ および d ₂ に伴わせる）、h、b ₁ および d ₂ についての「真の値」H、D ₁ および D ₂ の全ての可能な組合せ（状態 S = { s ₁ , s ₂ , . . . , s ₁₆ }）の以下の表を、各々、作成することができる。

【 0 2 9 3 】

【 数 7 7 】

仮説	H ₁ : h=d ₁	Ω={AA,AC,...,G G}	
状態	H	D ₁	D ₂
s ₁	A	A	A
s ₂	A	A	C
s ₃	A	A	T
s ₄	A	A	G
s ₅	C	C	A
s ₆	C	C	C
s ₇	C	C	T
s ₈	C	C	G
s ₉	T	T	A
s ₁₀	T	T	C
s ₁₁	T	T	T
s ₁₂	T	T	G
s ₁₃	G	G	A
s ₁₄	G	G	C
s ₁₅	G	G	T
s ₁₆	G	G	G

仮説	H ₂ : h=d ₂	Ω = {AA,AC,...,GG }	
状態	H	D ₁	D ₂
s ₁	A	A	A
s ₂	C	A	C
s ₃	T	A	T
s ₄	G	A	G
s ₅	A	C	A
s ₆	C	C	C
s ₇	T	C	T
s ₈	G	C	G
s ₉	A	T	A
s ₁₀	C	T	C
s ₁₁	T	T	T
s ₁₂	G	T	G
s ₁₃	A	G	A
s ₁₄	C	G	C
s ₁₅	T	G	T
s ₁₆	G	G	G

【 0 2 9 4 】

10

20

30

40

50

「真の値」 H 、 D_1 および D_2 は知られておらず、かつ生の測定結果 h 、 d_1 、 d_2 、 p_h 、 p_{d_1} 、 p_{d_2} のみが知られているので、全セット にわたる $P(M/H_1)$ および $P(M/H_2)$ の計算は以下のように行わなければならない：

【0295】

【数78】

$$P(M|H_1) = \sum_{(D_1, D_2) \in \Omega} P(M(h, d_1, d_2) | H_1 \& D_1, D_2) * P(D_1, D_2)$$

$$P(M|H_2) = \sum_{(D_1, D_2) \in \Omega} P(M(h, d_1, d_2) | H_2 \& D_1, D_2) * P(D_1, D_2)$$

10

もし、計算の目的で、 d_1 および d_2 、ならびに p_{d_1} および p_{d_2} が独立した変数であると仮定すれば：

【0296】

【数79】

$$P(M|H_1) = \sum_{\Omega} P(M(h, d_1, d_2) | H_1 \& D_1, D_2) * P(D_1, D_2) =$$

$$\sum_s P(M(h) | H) * P(M(d_1) | D_1) * P(M(d_2) | D_2) * P(D_1) * P(D_2)$$

を示すことができる。 $\{h, d_1, d_2\}$ における h についての、前記した最後の合計：
 $P(M(x)/X)$ 下で3項を計算する。

20

【0297】

(「真の対立遺伝子値」をヒットさせる) 正しい対立遺伝子要求の確率の計算は、対立遺伝子 X の真の値を仮定した結果 x の測定に基づく。もし測定された値 x および真の値 X が等しいならば、確率は p_x である(正しい測定の確率)。もし x および X が異なるならば、その確率は $(1 - p_x) / 3$ である。例えば、 $X = C$ 、および測定された値が $x = A$ である条件下で「真の値」 C が見出される確率を計算する。 A を得る確率は p_x である。 C 、 T また G を得る確率は $(1 - p_x) / 3$ である。従って、 C がヒットする確率は $(1 - p_x) / 3$ である。というのは、 C 、 T および G は等しくありそうと仮定することができるからである。

30

【0298】

もしインジケータ変数 I_x が計算に含まれ、ここで、もし $x = X$ であれば $I_x = 1$ であり、もし $x \neq X$ ならば $I_x = 0$ であれば、確率は以下の通りである：

$$P(M(x)/X) = I_{\{x=X\}} * p_x + (1 - I_{\{x=X\}}) * (1/3) * (1 - p_x)、\{h, d_1, d_2\} \text{ における } x。$$

さて、 $P(M|H_1)$ における最後の2つの項を考える。 $P(D_1)$ および $P(d_2)$ は、事前の知識から知ることができる、対立遺伝子 A 、 C 、 T および G の集団頻度である。

【0299】

特定の測定 $M(h = A, d_1 = G, d_2 = C)$ を仮定して、特定の状態 s_2 について先に示した表現を考慮する：

40

【0300】

【数80】

$$\begin{aligned} P(M(h)|H) * P(M(d_1)|D_1) * P(M(d_2)|D_2) * P(D_1) * P(D_2) &= \\ = P(M(h)=A|H=A) * P(M(d_1)=G|D_1=A) * P(M(d_2)=C|D_2=C) * P(D_1=A) * P(D_2=C) &= \\ = p_h * ((1-p_{d1})/3) * p_{d2} * f(D_1=A) * f(D_2=C) \end{aligned}$$

同様に、残りの15の状態、およびセット にわたる合計について、特定の測定(この場合、 $M(h = A, d_1 = G, d_2 = C)$)を仮定して(1)を計算する。

【0301】

さて、 $P(M/H_1)$ および $P(M/H_2)$ は計算された。最後に、前記したように $P(H_1/M)$ および $P(H_2/M)$ を計算する：

50

【 0 3 0 2 】

【 数 8 1 】

$$P(H_1 | M) = \frac{P(M | H_1)}{P(M | H_1) + P(M | H_2)}$$

$$P(H_2 | M) = \frac{P(M | H_2)}{P(M | H_1) + P(M | H_2)}$$

。

【 0 3 0 3 】

10

対立遺伝子要求の帰属および対応する確率

さて、4つの異なる仮説：

H_{p2A} : p₂ の「真の値」はAである、H_{p2C} : p₂ の「真の値」はCである、H_{p2T} : p₂ の「真の値」はTである、H_{p2G} : p₂ の「真の値」はGである、を確立し、P(H_{p2A} / M)、P(H_{p2C} / M)、P(H_{p2T} / M)、P(H_{p2G} / M)を計算する。最高の値は、特定の対立遺伝子の要求および対応確率を決定する。

【 0 3 0 4 】

p₂ の起源は未知である（それは、P(H₂ / M)の確率でもってd₁から、および確率P(H₁ / M)をもってd₂から由来する）ので、p₂対立遺伝子がd₁またはd₂に由来する双方の場合を考慮しなければならない。仮説H_Aについては、ベイズ則を適用し、

20

【 0 3 0 5 】

【 数 8 2 】

$$P(H_{p2A} | M) = P(H_{p2A} | M, H_1) * P(H_1 | M) + P(H_{p2A} | M, H_2) * P(H_2 | M)$$

が得られる。

【 0 3 0 6 】

P(H₁ / M)およびP(H₂ / M)は工程1においてすでに決定されている。ベイズ則によると、

30

【 0 3 0 7 】

【 数 8 3 】

$$P(H_{p2A} | H_1, M) = \frac{P(M | H_1, H_{p2A}) * P(H_1, H_{p2A})}{P(H_1, M)}$$

である。H₁は、p₂がd₂に由来することを示唆するので、前記したように、

【 0 3 0 8 】

【 数 8 4 】

$$P(M | H_1, H_{p2A}) = P(M(d_2) | D_2 = A) = I_{\{d_2=D_2\}} * p_{d2} + (1 - I_{\{d_2=D_2\}}) * (1 - p_{d2}) / 3$$

40

$$P(H_1, M / H_{p2A}) = P(M(d_2) / D_2 = A) = I_{\{d_2=D_2\}} * p_{d2} + (1 - I_{\{d_2=D_2\}}) * (1/3) * (1 - p_{d2})$$

である。

P(H_{p2A}) = P(D₂ = A) = f_{d2}(A)であり、ここで、f_{d2}(A)は集団頻度データから得られる。

$$P(H_1, M) = P(H_1, M / H_{p2A}) * P(H_{p2A}) + P(H_1, M / H_{p2C}) * P(H_{p2C}) + P(H_1, M / H_{p2T}) * P(H_{p2T}) + P(H_1, M / H_{p2G}) * P(H_{p2G})$$

同様に、P(H_{p2A} & H₂ / M)を計算する。

50

$P(H_{p_2 A} / M) = P(H_{p_2 A} \& H_1 / M) + P(H_{p_2 A} \& H_2 / M)$ であり、したがって、 p_2 が A に等しい確率は計算された。C、T および G についての計算を反復する。最高の値は、 p_2 対立遺伝子要求および対応する確率の回答を与えるであろう。

【0309】

対立遺伝子要求の p_1 への帰属（ハプロイド細胞、および1つのジプロイド細胞に対応する対立遺伝子）

前記したように、4つの異なる仮説：

- $H_{p_1 A}$: p_1 の「真の値」は A である、
- $H_{p_1 C}$: p_1 の「真の値」は C である、
- $H_{p_1 T}$: p_1 の「真の値」は T である、
- $H_{p_1 G}$: p_1 の「真の値」は G である、

10

を確立し、 $P(H_{p_1 A} / M)$ 、 $P(H_{p_1 C} / M)$ 、 $P(H_{p_1 T} / M)$ 、 $P(H_{p_1 G} / M)$ を計算する。

【0310】

これは、 $H_{p_1 A}$ の仕上げである。「真の場合」の場合においては、ハプロイドおよび対応するジプロイド細胞が A と等しい場合にのみ p_1 は A と等しい。したがって、 p_1 および p_{p_1} を計算するためには、ハプロイドおよび対応するジプロイド細胞が等しい状況を考慮しなければならない。したがって、仮説 $H_{p_1 A}$: p_1 の「真の値」は A であって、 H_{hdA} となる：ハプロイド細胞および対応するジプロイド細胞の「真の値」は A である。

20

【0311】

h の起源は未知である（それは、 $P(H_1 / M)$ の確率でもって d_1 から、および確率 $P(H_2 / M)$ でもって d_2 から由来する）ので、 h 対立遺伝子が d_1 または d_2 に由来する双方の場合を考慮し、 p_1 の決定におけるそれを実行しなければならない。それは、ベイズ則を用いると：

$$P(H_{hdA} | M) = P(H_{hdA} | M, H_1) * P(H_1 | M) + P(H_{hdA} | M, H_2) * P(H_2 | M)$$

を意味する。

【0312】

前記したように、 $P(H_1 / M)$ および $P(H_2 / M)$ は先の計算から知られている。

30

【0313】

【数85】

$$P(H_{hdA} | H_1, M) = \frac{P(H_1, M | H_{hdA}) * P(H_{hdA})}{P(H_1, M)}$$

$$P(H_1, M / H_{hdA}) = P(M(h) / H = A) * P(M(d_1) / D_1 = A) = [I_{\{h=H\}} * p_h + (1 - I_{\{h=H\}}) * (1/3) * (1 - p_h)] * [I_{\{d_1=D_1\}} * p_{d_1} + (1 - I_{\{d_1=D_1\}}) * (1/3) * (1 - p_{d_1})]$$

である。というのは、 H_1 は、 p_1 が d_1 に由来することを示唆するからである。 $P(H_{hdA}) = P(h = A) * P(D_1 = A) f_h(A) * f_{d_1}(A)$ であり、ここで、 $f_h(A)$ および $f_{d_2}(A)$ は集団頻度データから得られる。 $P(H_1, M) = P(H_1, M / H_{hdA}) * P(H_{hdA}) + P(H_1, M / H_{hdC}) * P(H_{hdC}) + P(H_1, M / H_{hdT}) * P(H_{hdT}) + P(H_1, M / H_{hdG}) * P(H_{hdG})$ 。

40

【0314】

同様に、 $P(H_{hdA} \& H_2 / M)$ を計算する。

$P(H_{hdA} / M) = P(H_{hdA} \& H_1 / M) + P(H_{hdA} \& H_2 / M)$ であり、今や、我々は p_1 が A と等しい確率を計算した。C、T および G についての計算を反復する。最高の値は p_1 対立遺伝子要求および対応する確率の回答を与えるであろう。

【0315】

例としての入力

50

2つの入力の例を示す。最初の例は、共分離する低い傾向があるSNPのセットのものであり、すなわち、SNPは染色体を通過して拡大し、入力データを表3に示す。第二の例は、共分離する高い傾向があるSNPのセットのものであり、すなわち、SNPは染色体上にクラスター形成し、インプットデータを表4に示す。双方のデータのセットは個体の測定されたSNPデータ、個体の親のSNPデータおよび対応する信頼性値を含む。このデータは現実の人々から測定された現実のデータであることを注記する。各列は、1つの特定のSNPロケーションについての測定を表す。行は、行の見出しによって示されるデータを含む。行の見出し中の略語に対する鍵は以下の通りである：

- f a m i l y _ i d = 各人についてのユニークなid (事務的理由について含む)
- s n p _ i d = SNP同定番号
- e 1 , e 2 = 胚についてのSNPヌクレオチド値
- p 1 , p 2 = 父親についてのSNPヌクレオチド値
- m 1 , m 2 = 母親についてのSNPヌクレオチド値
- p e 1 , p e 2 = e 1 , e 2 についての測定精度 p p 1 , p p 2 = p 1 , p 2 についての測定精度
- p m 1 , p m 2 = m 1 , m 2 についての測定精度

10

例としての出力

出力データの2つの例を表5および表6に示し、これは各々、表3および表4に掲げたデータからの出力データに対応する。双方の表は、個体の測定されたSNPデータ、個体の親のSNPデータ、個体のSNPデータの最もありそうな真の値、および対応する信頼性を示す。各列は、1つの特定のSNPに対応するデータを表す。行は行の見出しによって示されるデータを含まれる。行の見出し中の略語に対する鍵は以下の通りである：

20

- s n p _ i d = SNP同定番号
- t r u e _ v a l u e = e 1 , e 2 についての提案されたヌクレオチド値
- t r u e _ h y p = e 1 , e 2 の起源についての仮説
- e e = e 1 , e 2 についての測定されたSNPヌクレオチド値
- p p = p 1 , p 2 についての測定されたSNPヌクレオチド値
- m m = m 1 , m 2 についての測定されたSNPヌクレオチド値
- H y p P r o b = 最終仮説の確率。出力についてはただ1つの数があるが、優れた行構造のため、この数字は全ての列中に複製される。

30

【0316】

このアルゴリズムは手動で、またはコンピュータによって実施することができることを注記する。表3および表4は、該方法のコンピュータで実施されたバージョンについての入力データの例を示す。表5は表3に示された入力データに対する出力データを示す。表6は、表4に示された入力データに対する出力データを示す。

【0317】

シミュレーションアルゴリズム

以下に、システムの一体性を確実にし、およびより広く種々の状況におけるアルゴリズムの現実の効率を評価するためになされた第二のシミュレーションを示す。これを行うために、1,000のフルシステムシミュレーションを実行した。これは、親遺伝子データをランダムに作り出し、イン・シリコにて減数分裂を模倣して、胚データが得られ、胚データの不完全な測定をシミュレートし、次いで、本明細書中に開示された方法を実行して、シミュレートされた測定胚データを清浄化し、次いで、その「清浄化された」データを「現実の」データと比較することを含む。シミュレーションのより詳細な説明を以下に掲げ、事象のフローの目に見える表示を図18に掲げる。理論の2つの異なる実施を検定した。より十分な説明を以下に掲げる。

40

【0318】

DHおよびPSについてのシミュレーションアルゴリズムおよび結果

双方のアルゴリズムについて、初期入力変数は：

- (i) 検定すべきSNPのリスト、

50

(i i) 母性 (p o p f r e q l i s t M M) および父性 (p o p f r e q l i s t P P) 染色体の集団頻度、

(i i i) ハプロイド測定 (p h , p e) についての、および秩序立っていないジプロイド測定 (p d) についての正しい対立遺伝子要求の確率、
である。

【 0 3 1 9 】

これらの値は、関連する S N P についての経験的なデータ (集団頻度) からの、および測定機器性能 (p h , p d , p e) からの結果に基づいて固定すべきである。シミュレーションは、最もありそうな (通知された) 、均一な (通知されていない) および非常にありそうにない (極端な場合) のようないくつかのシナリオについて実行した。

10

【 0 3 2 0 】

一旦、前記した静的なパラメーターが固定されれば、特定の S N P を仮定した交差確率はすべてのシミュレーションについて同一であり、スニップロケーション (S N I P L O C _ N A M E _ M A T) および遺伝子距離 (H A P L O C _ N A M E _ M A T) についてのデータベースを仮定して該時点に先立って誘導されるであろう。

[c r o s s p r o b , s n i p s] =

GetCrossProb (スニップ , S N I P L O C _ N A M E _ M A T , パラメーター , H A P L O C _ N A M E _ M A T)

予備的シミュレーションループ

予備的シミュレーションループは、十分なシミュレーションで用いられるであろう遺伝子データが現実的であることを示すものである。工程 1 ないし 5 を 1 0 , 0 0 0 回反復した。このシミュレーションが、いずれかのまたは双方の親について実行することができ；該工程は同一であることを注記する。この場合、シミュレーションは説明目的のために父性ケースで実行され、図 1 8 への言及はカッコに入れた図 1 8 中の対応する母性エントリーも含む。

20

【 0 3 2 1 】

工程 1 : オリジナルの親ジプロイド細胞 (P 1 , P 2) の創製

[P 1 , P 2] = オリジナルの染色体の創製 (s n i p s , p o p f r e q l i s t P P) ; 1 8 0 1 (1 8 0 2)

父親細胞についての各 S N P に対する集団頻度に依存して、オリジナルの父性細胞を創製する。

30

【 0 3 2 2 】

工程 2 : D H A l g o についてのハプロイドおよび秩序立っていないジプロイドデータの創製

親染色体 1 8 0 3 の交差をシミュレートして、染色体、交差の 2 つのセット : P 1 C 1 、 P 2 C 1 および P 1 C 2 、 P 2 C 2 ; 1 8 0 4 (1 8 0 5) を得る。ハプロイド対立遺伝子 H P 1 8 0 7 (1 8 0 8) 、この場合は、P 1 (というのは、いずれについても差はないからである) についての、(第一のセットからの) 交差 1 8 0 6 後に父親対立遺伝子のうちの 1 つをピックアップし、ジプロイド対立遺伝子中の順序を混合して、(D 1 P , D 2 P) 1 8 0 7 (1 8 0 8) を得る。

40

H P = P i c k O n e (P 1 C 1 , P 2 C 1) ;

[D 1 P , D 2 P] = J u m b l e (P 1 , P 2) 。

【 0 3 2 3 】

工程 3 : オリジナルなデータセットへエラーを導入して、測定をシミュレートする。

【 0 3 2 4 】

正しい測定 (p h - ハプロイド、p d - ジプロイド測定) の所与の確率に基づき、エラーを測定を導入して、シミュレートされた測定親データ 1 8 1 1 (1 8 1 2) を得る。

h p = M a k e E r r o r (H P , p h) ;

d 1 p = M a k e E r r o r (D 1 P , p d) ;

d 2 p = M a k e E r r o r (D 2 P , p d) 。

50

【0325】

工程4：DHALgoを適用して、(p1, p2)、(pp1, pp2)を得る。

【0326】

DHALgoは、ハプロイド細胞からの対立遺伝子、およびジプロイド細胞からの秩序立っていない対立遺伝子を取り、これらを生起したもつともありそうな秩序立ったジプロイド対立遺伝子を戻す。DHALgoは(P1, P2)を再形成するよう試み、また、父親についての見積もり誤差(pp1, pp2)を戻す。比較のために、単純な対立遺伝子マッチングを行う経験的アルゴリズムを用いる。目標は、単純な経験的アルゴリズムと比較して、どれくらい開示されたアルゴリズムが良好であるかを比較することである。

[p1, p2, pp1, pp2] = DHALgo(hp, d1p, d2p, ph, pd, snips, popfreqlistPP, 'DH');

[p1s, p2s, pp1s, pp2s] = DHALgo(hp, d1p, d2p, ph, pd, snips, popfreqlistPP, 'ST');

。

【0327】

工程5：実行のための統計学の収集

(P1, P2)を誘導された(p1, p2)と比較する。

[P1cmp(:, i), P2cmp(:, i), P1prob(:, i), P2prob(:, i), P1mn(i), P2mn(i)] = DHSimValidate(P1, P2, p1, p2, pp1, pp2);

注意：(P1Si, P2Si, P1Pi, P2Pi, P1Ai, P2Ai) = (I_{P1=p1}, I_{P2=p2}, Pp1, Pp2, P1acc, P2acc)であり、ここで、I_{P1=p1}は全てのSNPについての、同様に、I_{P2=p2}についての、DHアルゴリズム精度の見積もり用のバイナリインジケータアレイである。Pp1, Pp2は該アルゴリズムに由来する正しい対立遺伝子要求およびp1acc = 平均(I_{P1=p1})、すなわち、p2accについてと同様な、p1についてのこの実行に対する平均精度の確率である。

【0328】

予備的シミュレーションの結果

10,000のシミュレーションを用いて、P1, P2からのDHアルゴリズムの総じての精度を示す、アルゴリズム精度DHAaccuracy.P1 = 平均(P1Ai)、DHAaccuracy.P2 = 平均(P2Ai)を見積もった。個々のSNPに基づき、各SNP SNPACC.P1 = 平均(P1Si)についての平均精度は、SNP, SNPProb.P1 = 平均(P2Pi)であると正しく測定する見積もられた確率の平均に合致すべきであり、すなわち、もしアルゴリズムが正しく作動すれば、SNPACC.P1に対する値はSNPPro.P1に密接に対応すべきである。これらの2つの間の関係はそれらの相関によって反映される。

【0329】

シミュレーションの10000ループは異なる設定シナリオで実行した：

(1) 基礎となる集団頻度は、より現実的である現存のゲノタイピングデータ、およびA、C、T、Gが各SNPについて同一の確率を有する均一な集団頻度によって与えられた。

(2) ハプロイドおよび秩序立っていないジプロイド測定(PH, PD)についての測定精度に対するいくつかの組合せ。種々の仮定を行った；測定は共に非常に精度があり(0.95, 0.95)、精度が低く(0.75, 0.75)、および精度なしまたはランダムであり(0.25, 0.25)、ならびに(0.9, 0.5)、(0.5, 0.9)のバランスが取れていない組合せである。現実にもっと近いであろうものは、ほぼ0.6ないし0.8の精度であろう。

(3) シミュレーションを、DHALgorithmおよび単純なマッチングSTAlgorithm双方についてのすべてのこれらの場合に実行して、開示されたアルゴリズム

10

20

30

40

50

の性能を評価した。

これらの全ての実行の結果を表 7 にまとめる。

【 0 3 3 0 】

開示されたアルゴリズムは、これらのシミュレーションにおいて、特に、不均一な集団頻度、および正しい測定のアンバランスな、または低下した確率の現実の場合について、現存の経験的アルゴリズムよりも良好に実行される。また、個々の SNP についてのアルゴリズム精度の 1 つの見積もりはこれらの場合において非常に良好であることが確認された。というのは、正しい対立遺伝子要求の見積もられた精度、およびシミュレーション平均精度の間の相関は 99% 程度であり、平均比率は 1 だからである。

【 0 3 3 1 】

最も現実の場合において、データ集団頻度および $(p_h, p_d) = (0.6, 0.8)$ については、 (P_1, P_2) についての正しく検索された SNP の平均パーセントは実行 1 において $(0.852, 0.816)$ であって、実行 2 において $(0.601, 0.673)$ である。

【 0 3 3 2 】

表 7 および表 8 については、「データ」使用集団頻度データで始まる列は経験的結果から取られたものであり、他方、「均一」で始まる列は均一な集団を仮定することに注意されたし。

【 0 3 3 3 】

表 7 および表 8 においては、精度は、正しい SNP 要求がなされ、正しい元の染色体が同定された SNP の平均パーセントとして定義されることに注意するのは重要である。また、これらのシミュレーションはアルゴリズムの 2 つの可能な実行を反映するのに注意するのも重要である。良好な結果を与えることができるアルゴリズムを実行する他の方法があり得る。このシミュレーションは、該方法が実施できることを示すつもりだけである。

【 0 3 3 4 】

十分なシミュレーションループ

工程 1 ないし 8 を 10000 回反復した。これは、関連する個体、この場合は、親から測定された遺伝子データを用いて標的固体についての測定された遺伝子データを清浄化する十分に開示された方法を検定するためのシミュレーションである。

【 0 3 3 5 】

工程 1：オリジナルの親ジブroid細胞 (P_1, P_2) , (M_1, M_2) の創製

$[P_1, P_2] =$ オリジナルな染色体の創製 (snips, popfreqlist PP)
); (1801)

$[M_1, M_2] =$ オリジナルな染色体の創製 (snips, popfreqlist MM)
); (1802)

母親および父親細胞についての各 SNP に対する集団頻度に依存して、オリジナルな親細胞を創製する。

【 0 3 3 6 】

工程 2：交差親細胞 (P_1C, P_2C) , (M_1C, M_2C) (1803)

交差を持つ父性細胞の 2 つのセットを創製して：第一に、DHA1go で用いる (P_1C_1, P_2C_1) が得られ、第二に、PSA1go で用いる (P_1C_2, P_2C_2) を得る。(1804)

交差を持つ母性細胞の 2 つのセットを創製して：第一に、DHA1go で用いる (M_1C_1, M_2C_1) 、および PSA1go で用いる (M_1C_2, M_2C_2) を得る。(1805)

$[P_1C_1, P_2C_1] = (P_1, P_2, \text{snips}, \text{fullprob})$ を交差させる
;

$[P_1C_2, P_2C_2] = (P_1, P_2, \text{snips}, \text{fullprob})$ を交差させる
;

$[M_1C_1, M_2C_1] = (M_1, M_2, \text{snips}, \text{fullprob})$ を交差させる

10

20

30

40

50

;

[M 1 C 2 , M 2 C 2] = (M 1 , M 2 , s n i p s , f u l l p r o b) を交差させる

;

。

【 0 3 3 7 】

工程 3 D H A l g o についてのハプロイド細胞および無秩序ジプロイド細胞を作成する (1 8 0 6) 。

【 0 3 3 8 】

ハプロイド細胞 H P についての父性細胞のセットのうち 1 つ (1 8 0 4 , 第一のセット) をピックアップし、ジプロイド細胞中の順序を混合して、(D 1 P , D 2 P) (1 8 0 7) を得る。母性細胞 (1 8 0 5 , 第一のセット) についても同様にして、M H , (D 1 M , D 2 M) を得る。 (1 8 0 8)

H P = 1 つの (P 1 C 1 , P 2 C 1) をピックアップする ;

H M = 1 つの (M 1 C 1 , M 2 C 1) をピックアップする ;

[D 1 P , D 2 P] = (P 1 , P 2) を乱雑とする ;

[D 1 M , D 2 M] = (M 1 , M 2) を乱雑とする ;

。

【 0 3 3 9 】

工程 4 : ジプロイド胚細胞の作成 (1 8 0 9)

胚細胞について父性細胞の 1 つ (1 8 0 4 , 第二のセット) および母性細胞の 1 つ (1 8 0 5 , 第二のセット) をピックアップする。測定目的で順序を混合する。

E 1 = 1 つの (P 1 C 2 , P 2 C 2) をピックアップする ;

E 2 = 1 つの (M 1 C 2 , M 2 C 2) をピックアップする ;

[E 1 J , E 2 J] = (E 1 , E 2) を乱雑とする ; (1 8 1 0) 。

【 0 3 4 0 】

工程 5 : 測定 (1 8 1 1 , 1 8 1 2 , 1 8 1 3) に誤差を導入する

所与の測定誤差 (H P - ハプロイド細胞 , P D - 無秩序ジプロイド細胞 , p e - 胚細胞) に基づいて、測定に誤差を導入する。

h p = 誤差 (H P , p h) を作りだす ; (1 8 1 1)

d 1 p = 誤差 (D 1 P , p d) を作りだす ; (1 8 1 1)

d 2 p = 誤差 (D 2 P , p d) を作りだす ; (1 8 1 1)

h m = 誤差 (H M , p h) を作りだす ; (1 8 1 2)

d 1 m = 誤差 (D 1 M , p d) を作りだす ; (1 8 1 2)

d 2 m = 誤差 (D 2 M , p d) を作りだす ; (1 8 1 2)

e 1 = 誤差 (E 1 J , p e 1) を作りだす ; (1 8 1 3)

e 2 = 誤差 (E 2 J , p e 2) を作りだす ; (1 8 1 3) 。

【 0 3 4 1 】

工程 6 : D H A l g o を適用して、(p 1 , p 2)、(m 1 , m 2)、(p p 1 , p p 2)、(p m 1 , p m 2) を得る。

【 0 3 4 2 】

D H A l g o はハプロイド細胞および無秩序ジプロイド細胞を取り、これらを生起させた最もありそうな秩序立ったジプロイド細胞を戻す。D H A l g o は父親染色体について (P 1 C 1 , P 2 C 1)、および母親染色体について (M 1 C 1 , M 2 C 1) を再形成するよう試み、また父親 (p p 1 , p p 2) および母親 (p m 1 , p m 2) 細胞についての見積もり誤差を戻す。

[p 1 , p 2 , p p 1 , p p 2] = D H A l g o (h p , d 1 p , d 2 p , s n i p s , p o p f r e q l i s t P P) ; (1 8 1 4)

[m 1 , m 2 , p m 1 , p m 2] = D H A l g o (h m , d 1 m , d 2 m , s n i p s , p o p f r e q l i s t M M) ; (1 8 1 5) 。

【 0 3 4 3 】

10

20

30

40

50

工程7: PSALgoを適用して、(DE1, DE2)(1816)を得る。

【0344】

PSALgoは再形成された親細胞(p1, p2, m1, m2)および無秩序な測定胚細胞(e1, e2)を取って、最もありそうな秩序立った真の胚細胞(DE1, DE2)を戻す。PSALgoは、(E1, E2)を再形成するよう試みる。

```
[DE1, DE2, alldata] = PSALgo(snips, e1, e2, p1,
p2, m1, m2, pe, pp1, pp2, pm1, pm2, parameters, c
rossprob, popfreqlistPP, popfreqlistMM);
```

。

【0345】

工程8: このシミュレーション実行からの望まれる統計学の収集
実行についての統計学を得る:

```
simdata = SimValidate(alldata, DE1, DE2, P1, P
2, M1, M2, E1, E2, p1, p2, m1, m2, e1, e2, pe, pe, pp
1, pp2, pm1, pm2);
```

。

【0346】

シミュレーションの結果

10000のシミュレーションを実行し、E1, E2からのPSアルゴリズムの全精度を我々に告げる、アルゴリズム精度についての最終見積もりPSAccuracy.E1 = 平均(E1Ai)、PSAccuracy.E2 = 平均(E2Ai)を計算した。個々のSNPに基づき、各SNP SNPAcc.E1 = 平均(E1Si)についての平均精度は、SNP, SNPprob.E1 = 平均(E2Pi)であると正しく測定する見積もられた確率の平均に合致するはずであり、すなわち、もしアルゴリズムが正しく書かれれば、SNPAcc.E1は、SNPprob.E1に相関するように観察されるはずである。これらの2つの間の関係はそれらの相関によって反映される。

【0347】

シミュレーションの10000ループを異なる設定シナリオについて実行した:

(1) より現実的である現存のゲノタイピングデータ、およびA、C、T、Gが各SNPにおいて同一の確率を有する均一な集団頻度によって与えられる基礎となる集団頻度

(2) ハプロイド、無秩序ジブロイドおよび胚測定(ph, pd, pe)についての測定精度のいくつかの組合せ。種々の精度をシミュレートした: 非常に精度がある(0.95, 0.95, 0.95)、精度が低い(0.75, 0.75, 0.75)、および精度なしまたはランダム(0.25, 0.25, 0.25)、ならびに(0.9, 0.5, 0.5)、(0.5, 0.9, 0.9)のアンバランスな組合せ。現実に最も近いであろうものは、ほぼ(0.6, 0.8, 0.8)である。

(3) すべてのこれらの場合において、我々のPSAlgorithmおよび単純なマッチングSTPSAlgorithm双方についてシミュレーションを行って、開示されたアルゴリズムの性能を評価した。

これらの実行の結果を表8にまとめる。

【0348】

開示されたアルゴリズムは、これらのシミュレーションにおいて、特に、不均一な集団頻度および正しい測定のアンバランスな、または低下した確率の現実の場合について、現存の経験的アルゴリズムよりも良好に実行される。また、このSNPについてのアルゴリズム精度の見積もりはこれらの場合において非常に良好であることが示された。というのは、正しい対立遺伝子要求の見積もられた精度、およびシミュレーション平均精度の間の相関は99%程度であり、平均率は1だからである。

【0349】

最も現実的な場合において、データ集団頻度および(ph, pd, pe) = (0.6, 0.8, 0.8)については、(E1, E2)については正しく検索されたSNPの平均バ

10

20

30

40

50

ーセントは実施1において(0.777, 0.788)および実施2において(0.835, 0.828)である。前記したように、アルゴリズムの平均精度を示す数は正しいSNPの要求のみならず、SNPの正しい親起源の同定もいう。効果的であるためには、アルゴリズムは、それが測定されるにつれデータを単純に許容するアルゴリズムよりも良好な結果を戻さなければならない。ある場合には、アルゴリズムの精度には測定のリストされた精度よりも低いを見て驚くであろう。このシミュレーションの目的では、もしそれが共に正しく要求され、また、その親および元の染色体が正しく同定された場合のみ、SNPの要求は正確であると考えられる。偶然にこれを正しくするチャンスは測定精度よりもかなり低い。

【0350】

出生前および胚遺伝物質を得るのに必要な実験室的技術

ゲノタイピングのための細胞およびDNA断片の単離を可能とする多くの利用できる技術がある。本明細書中に記載されたシステムおよび方法をこれらの技術、特に、母性血液からの胎児細胞またはDNAの単離、またはIVFの関係で胚からの胚盤胞の単離を含むものの中でいずれにも適応することができる。それはイン・シリコにてゲノムデータに同等に適応することができ、すなわち、遺伝物質から直接的に測定できない。

【0351】

システムの1つの実施形態において、このデータは以下に記載するように獲得することができる。

【0352】

細胞の単離

成人ジプロイド細胞はバルク組織または血液試料から得ることができる。成人ジプロイド単一細胞は、FACS、または蛍光活性化細胞ソーティングを用い、全血液試料から得ることができる。成人はプロイド単一精子細胞もまた、FACSを用いて精子試料から単離することができる。成人ハプロイド単一卵細胞は、IVF手法の間に卵収穫に関して単離することができる。

【0353】

ヒト胚からの標的単一胚盤胞の単離は、体外受精クリニックにおいて普通の技術に従って行うことができる。母性血液中の標的胎児細胞の単離は、モノクローナル抗体、あるいはFACSまたは密度勾配遠心のような他の技術を用いて達成することができる。

【0354】

DNA抽出は本出願についての標準的でない方法も含むであろう。DNA抽出についての種々の方法を比較する文献の報告は、いくつかの場合において、N-ラウロイルサルコシンの添加の使用のような新規なプロトコルは、より効果的であることが判明し、最も少ない偽陽性を生じることを見出している。

【0355】

ゲノムDNAの増幅

ゲノムの増幅は、連結-媒介PCR(LM-PCR)、縮重オリゴヌクレオチドプライマーPCR(DOP-PCR)、および多重置換増幅(MDA)を含めた多数の方法によって達成することができる。これらの方法のうち、DOP-PCRは、染色体の単一コピーを含めた、少量のDNAから多量のDNAを信頼性よく生じさせ；この方法は、データ忠実度が臨界的である親ジプロイドデータをゲノタイピングするために最も適しているであろう。MDAは最速な方法であり、数時間以内にDNAの100倍増幅を生じる；この方法は、胚細胞をゲノタイピングするのに、あるいは時間が必須である他の状況において最も適切であろう。

【0356】

バックグラウンド増幅はこれらの方法の各々で問題である。というのは、各方法は、潜在的に、汚染DNAを増幅するだろうからである。非常に少量の汚染はアッセイを不可逆的に毒し、偽データを与えかねない。従って、増幅前および後ワークフローが完全に物理的に分離されたクリーンな実験室条件を用いるのが非常に重要である。DNA増幅のため

10

20

30

40

50

のクリーンな汚染なしのワークフローは、今日、産業的分子生物学においてルーチン的であって、単に詳細に対して思慮深い注意を必要とする。

【0357】

ゲノタイピングアッセイおよびハイブリダイゼーション

増幅されたDNAのゲノタイピングは、Affymetrix's Genflex Tag Arrayのような分子逆転プローブ(MIP)、Affymetrix's 500KアレイまたはIllumina Bead Arraysのようなマイクロアレイ、またはApplied Bioscience's TaqmanアッセイのようなSNPゲノタイピングアッセイを含めた多くの方法によって行うことができる。Affymetrix's 500Kアレイ、MIPs/Genflex、TaqmanおよびIlluminaアッセイは、全て、マイクログラム量のDNAを必要とし、従って、いずれかのワークフローでの単一細胞のゲノタイピングはいくつかの種類の増幅を必要とするであろう。これらの技術の各々は、とりわけ、コスト、データの質、定量的vs定性的データ、慣用化性、アッセイを完了するための時間、および測定可能なSNPの数の点で種々の釣り合いを有する。500KおよびIlluminaアッセイの利点は、10,000のSNPのオーダーで検出できるMIP、およびより少数さえを検出できるTaqmanアッセイとは反対に、それがデータを集めることができるSNPの大きな数、およそ250,000である。500Kアレイよりも優れたMIP、TaqmanおよびIlluminaアッセイの利点は、それらが固有に慣用化可能であり、ユーザーがSNPを選択するのを可能とすることであり、他方、500Kアレイはそのような慣用化を可能としない。

10

20

【0358】

IVFの間における着床前の診断の関係では、固有の時間の制限は重要であり；この場合、応答時間に換えてデータの質を犠牲にするのは有利であろう。それは他の明瞭な利点を有するが、標準MIPアッセイプロトコルは、典型的には、完了するのに2.5ないし3日かかる比較的時間を消費するプロセスである。MIPにおいて、DNAを標的とするためのプローブのアニーリング、および増幅後ハイブリダイゼーションは特に時間を消費し、これらの時間からのいずれの偏差もデータ質の劣化をもたらす。プローブはDNA試料に一晚アニーリングさせる(12ないし16時間)。増幅後ハイブリダイゼーションはアレイに一晚アニーリングさせる(12ないし16時間)。アニーリングおよび増幅双方の前および後の多数の他の工程は、プロトコルの合計標準タイムラインを2.5日とする。スピードについてのMIPアッセイの最適化は、潜在的に、プロセスを36時間未満に低下させることができよう。500KアレイおよびIlluminaアッセイは共により速い応答時間：ほぼ1.5ないし2日を有して、標準的プロトコルにおいて高度に信頼性があるデータを生じる。これらの方法の双方は最適化可能であり、500Kアレイについてのゲノタイピングアッセイおよび/またはIlluminaアッセイのための応答時間は24時間未満まで低下させることができようと思積もられる。なおより速いのはTaqmanアッセイであり、これは3時間以内に行うことができる。これらの方法の全てについて、アッセイ時間の低下の結果、データの質の低下をもたらすが、それは、正確には、開示された発明が何を取り組むように設計されているかである。より速いいくつかの利用可能な技術は、特に高-スループットではなく、従って、この時点において高度に平衡な出生前遺伝子診断で使用できない。

30

40

【0359】

当然に、IVFの間における胚盤胞のゲノタイピングのような、時間が臨界的である状況においては、より速いアッセイはより遅いアッセイよりも明瞭な利点を有し、他方、IVFの前に出生前DNAをゲノタイピングすることが開始されている場合のような、そのような時間圧力を有しない場合には、他の因子が適当な方法を選択するのに支配的であろう。例えば、もう1つの技術に対する1つの技術から出てくるもう1つの釣り合いは、価格vsデータ質のものである。より重要な測定のための高い質のデータを与えるより効果的な技術および忠実度が臨界的でない測定用のより低い質のデータを与える安価な技術を用

50

いるのは理にかなっているであろう。十分に迅速な高 - スループットゲノタイピングの点まで開発されたいずれの技術を用いて、この方法で用いる遺伝物質をゲノタイピングすることもできよう。

【 0 3 6 0 】

該方法の関連数例

どのようにして、開示された方法を、I V F手法の時間拘束内に全ての生きた胚の十分なゲノタイピングを可能とするであろう。I V F実験室の関係で用いることができるかの例をここに記載する。卵受精から胚着床までの、I V F実験室に必要な応答時間は3日下である。これは、関連する実験室的作業、データの清浄化および表現型予測がその時間内に完了させなければならないことを意味する。このシステムの模式的ダイアグラムを図19に示し、本明細書中に記載する。このシステムは、ゲノタイピングシステムを用いてI V F l a b 1 9 0 4で分析されるI V Fユーザー(母親)1902およびI V Fユーザー(父親)1903からの親遺伝子試料1901よりなることができる。それは、母親1902から収穫され、父親1903からの精子で受精させて、多数の受精した胚1905を作り出す多数の卵を含むことができる。それは、各胚について胚盤胞を抽出し、各胚盤胞のDNAを増幅し、高スループットゲノタイピングシステム1906を用いてそれらを分析する実験室技術者を含むことができる。それは、親からの、および胚盤胞からの遺伝子データをデータ保護プロセッシングシステム1907に送ることを含む、該システムは胚遺伝子データを確証し、清浄化する。それは、フェノタイピングアルゴリズム1909によって操作されて、各胚の表現型感受性を予測する清浄化胚データ1908を含むことができる。それは、I V Fユーザー1902および1903が母親1901における着床について胚を選択するのを助ける医師1910に送られる関連信頼性レベルと共にこれらの予測を含むことができる。

【 0 3 6 1 】

遺伝子データの清浄化に関連する雑多な注意

本明細書中に記載される方法は遺伝子データの清浄化に関することに注意するのは有用であり、すべての生き物は遺伝子データを含むので、該方法は親から染色体を受け継ぐいずれのヒト、動物または植物にも等しく適用することができる。動物および植物のリストは、限定されるものではないが、ゴリラ、チンパンジー、ピグミーチンパンジー、ネコ、イヌ、パンダ、ウマ、ウシ、ヒツジ、ヤギ、ブタ、チーター、トラ、ライオン、サケ、サメ、クジラ、ラクダ、バイソン、マナティー、ウナギ、メカジキ、イルカ、アルマジロ、カリバチ、ゴキブリ、虫、コンドル、ワシ、スズメ、チョウ、セコイア、トウモロコシ、小麦、米、ペチュニア、カウズベッチ、ヒマワリ、ブタクサ、カシノキ、栗の木およびアタマジラミを含む。

【 0 3 6 2 】

遺伝子データの測定は、特に、遺伝物質の試料が少量である場合に完全なプロセスではない。測定は、しばしば、正しくない測定、不明瞭な測定、誤った測定、および失われた測定を含む。本明細書中に記載された方法の目的は、これらの誤差のいくつかまたはすべてを検出し、修正することにある。この方法を用い、遺伝子データがかなり知られる信頼性を改良することができる。例えば、現行の技術を用い、単一細胞から増幅されたDNAからの不明瞭な測定遺伝子データは、20%および50%の間の未測定領域、または対立遺伝子ドロップアウトを含み得る。いくつかの場合において、遺伝子データは20%および99%の間の未測定領域、または対立遺伝子ドロップアウトを含み得るであろう。加えて、所与の測定SNPの信頼性は同様に誤差に従う。

【 0 3 6 3 】

未清浄化データがほぼ50%の対立遺伝子ドロップアウト率を有する場合において、本明細書中に開示された方法を適応した後に、清浄化されたデータは少なくとも90%の場合において正しい対立遺伝子要求を有し、理想的な状況下では、これは99%またはそれを超えるまで上昇し得ると予測される。未清浄化データがほぼ80%の対立遺伝子ドロップアウト率を有する場合において、本明細書中に開示された方法を適応した後に、清浄化

10

20

30

40

50

データは少なくとも95%の場合において、正しい対立遺伝子要求を有し、理想的な状況下ではこれは99%またはそれを超えるまで上昇し得ると予測される。未清浄化データがほぼ90%の対立遺伝子ドロップアウト率を有する場合、本明細書中に開示された方法を適用した後に、清浄化データは少なくとも99%の場合において正しい対立遺伝子要求を有し、理想的な状況下ではこれは99%以上まで上昇し得ると予測される。特定のSNP測定が90%近くの信頼性率でもってなされる場合、清浄化データは95%を超える、および理想的な場合には、99%を超える、またはそれを超える信頼性率でもってSNP要求を有すると予測される。特定のSNP測定が99%近くの信頼性率でもってなされる場合において、清浄化データは、99.9%を超えるおよび理想的な場合には99.99%を超える、またはそれよりも高い信頼性率でもってSNP要求を有すると予測される。

10

【0364】

また、1つの胚盤胞からの増幅されたDNAを測定することによって創製することができる胚遺伝子データは、多数の目的で使用することができるのみ注意するのも重要である。例えば、それは、異数体、片親二染色体を検出し、個体の性別を鑑定し、ならびに複数の表現型予測を行うのに用いることができる。現在、IVF実験室においては、用いる技術のため、しばしば、それは、胚盤胞が異数性のような1つの障害、または特定の単一遺伝子病についてテストするのに十分な遺伝物質を供することができるに過ぎない場合である。本明細書中に開示された方法は、なされる予測のタイプに拘わらず、胚盤胞からSNPの大きなセットを測定する通常の最初の工程を有するので、医師または親は、スクリーニングすべき限定された数の障害を選択することを強制されない。その代わりに、医療的知識の状態が許容する程度に多くの遺伝子および/または表現型についてスクリーニングするオプションが存在する。開示された方法では、胚盤胞のゲノタイピングに先立ってスクリーニングするための特定の条件を同定する唯一の利点は、もしあるPSNPが特に関連すると決定されたならば、注目するPSNPとより共分離するようなNSNPのより適切なセットを選択することができ、かくして、注目する対立遺伝子の要求の信頼性を増大させることである。SNPが先立って個人化されない場合においてさえ、信頼性は、本明細書中に記載された種々の目的で適切なものを超えると予測されることを注記する。

20

【0365】

表現型および臨床的予測

遺伝子型および臨床的情報から表現型データを予測するのに利用できる多くの方法がある。異なるモデルは、利用できるデータの量およびタイプに基づいて、異なる状況においてより適切である。表現型予測のための最も適切な方法を選択するためには、テストデータのセットについて多数の方法をテストし、テストデータの測定された結果と比較する場合に、予測の最良の精度をいずれの方法が提供するかを決定するのがしばしば最良である。本明細書中に記載されたある実施形態は、組合せて採用され、かつテストデータでの性能に基づいて選択された場合に、正確な表現型予測を行う高い尤度を供する方法のセットを含む。まず、(ii)偶発事象表を用いるシナリオでの遺伝子型-表現型モデリングのための技術を記載する。次に、(iii)凸最適化によって形成された回帰モデルを用いるシナリオにおける遺伝子型-表現型モデリングのための技術を記載する。次いで、予測すべき特定の表現型、特定の患者のデータ、およびモデルを訓練し、テストするためのデータの特定のセットを仮定して最良のモデルを選択するための技術を記載する。

30

40

【0366】

今日のデータ：偶発事象表に基づく表現型結果のモデリング

公知の遺伝子的欠陥、および病気表現型の確率を増加させる対立遺伝子がある場合、およびプレディクターの数が十分に少数である場合、表現型確率は偶発事象表でモデル化することができる。もした1つの関連遺伝子対立遺伝子があれば、特定の対立遺伝子の存在/不存在はA+/A-として記載することができ、病気表現型の存在/不存在はD+/D-として記載することができる。(f₁, N₁, f₂, N₂)を含有する偶発事象表は：

【0367】

50

【数86】

	D+	D-	#
G+	f ₁	1-f ₁	N ₁
G-	f ₂	1-f ₂	N ₂

$$S^2 = \frac{N_1 N_2 (N_1 + N_2) (p_1 (1 - p_2) - (1 - p_1) p_2)^2}{(p_1 N_1 + p_2 N_2) ((1 - p_1) N_1 + (1 - p_2) N_2)}$$

である。ここで、f₁ および f₂ は測定された頻度または異なる結果の確率を表し、対象の合計数は N = N₁ + N₂ である。この表から、独立変数 (IV) G+ または G- を有する2つの場合において病気状態 D+ を有する確率についてのオッズ比は、95%信頼区間を持つ OR = f₁ (1 - f₂) / f₂ (1 - f₁) : S が標準偏差である OR^{1 ± 1.96 / S} として報告することができる。例えば、10,000 の個体における乳癌の実験を用い、ここで、N+ は BRCA1 または BRCA2 対立遺伝子の存在を表す。

10

【0368】

【数87】

	D+	D-	#
M+	.563	.437	1720
M-	.468	.532	8280

このデータの結果、信頼区間 [1.31 ; 1.62] でのオッズ比率 OR = 1.463 がもたらされ、これを用いて、所与の突然変異を持つ乳癌の出現の増大した確立を予測することができる。2 × 2 よりも大きな偶発事象表を用いて、より独立した変数または結果変数を収容することができることを注記する。例えば、乳癌の場合には、偶発事象 M+ および M- は4つの偶発事象 : BRCA1 および BRCA2、BRCA1 および BRCA2 ではない、および BRCA1 ではなくおよび BRCA2、および最後に BRCA1 でなく BRCA2 でない ; で置き換えることができよう。どのようにして 2 × 2 を超える偶発事象表についての信頼区間を決定するのは当業者によってよく理解されるこの技術は、独立変数の異なる偶発事象によって定義される異なる群における患者をカウントすることによって低い標準偏差を持つモデルを形成するのに十分に少数の IV および十分なデータがある場合に用いられる。このアプローチは、回帰モデルを構築する場合に必要なように異なる IV をモデル化すべき結果に関連させる数学モデルを設計する困難性を回避する。

20

30

【0369】

特定の SNP からの遺伝子データは、特に、HapMap プロジェクトで認識される SNP の異なるパターンのような独立した変数の他の空間へ投影することもできることを注記する。HapMap 投影は個体をピンにクラスター化し、各ピンは SNP の特定のパターンによって特徴づけられるであろう。例えば、1つのピン (B1) は BRCA1 および BRCA2 を含有する SNP パターンを有し、もう1つのピン (B2) は BRCA1 を含有するが、BRCA2 を含有しない SNP パターンを有し、および第三のピンは、突然変異のすべての他の組合せに関連する SNP パターン (B3) を含有すると考える。これらの SNP のすべての異なる組合せを表す偶発事象表を作成するよりはむしろ、偶発事象 B1、B2 および B3 を表す偶発事象表を作成することができる。

40

【0370】

HapMap 投影によって記載されるように、ある SNP が一緒におこる傾向を用いて、プレディクターとして多数の SNP を用いるモデルを作成することができ、次いで、データは患者の別々の群よりなり、ここで、各群はただ1つの測定された SNP を有することにさらに注意されたし。この問題は、OMIM から入手可能なもののような公に入手可能な研究論文からモデルを作成する場合に普通に遭遇し、多数の SNP は表現型を予測するものではあるが、各論文は唯一の測定された関連 SNP を有するコホートについてのデ

50

ータを含有する。今日利用可能なデータを用いて予測モデルを形成するのに有用なこの態様を説明するために、IV：アルツハイマー病の家族履歴、性別、人種、年齢、3つの遺伝子、すなわち、APOE、NOS3、およびACEの種々の対立遺伝子に基づいて予測モデルを形成することができるアルツハイマー病に特に言及する。この病気との関係では、アルツハイマー病を超える多くの病気に適用される普及した論点を議論し：多くの遺伝子は特定の表現型についての特性の決定に関与するが、履歴研究のほとんど大部分は特定の遺伝子の対立遺伝子をサンプリングしたのに過ぎなかった。アルツハイマー病の場合においては、ほとんど全ての研究コホートは唯一の遺伝子をサンプリングしたに過ぎなかった；すなわち、APOE、NOS3、またはACE。それにも拘わらず、利用可能なデータの大部分が、唯一の遺伝子を調べる研究から由来する場合でさえ、多数の遺伝子対立遺伝子を入力するモデルを形成するのが重要である。この問題は、2つの表現型状態の単純化された場合、および各々が丁度2つの状態を持つ、2つの関連遺伝子を表す唯2つの独立した変数を考慮することによって説明される1つの態様において取り込まれる。病気表現型を記載するランダム変数D [D+, D-]、および遺伝子を記載する2つのランダムな変数A [A+, A-]およびB [B+, B-]を仮定すれば、目標はP(D/A, B)の最良の可能な見積もりを見出すことである。これは、 $P(D/A, B) = P(A, B/D)P(D)/P(A, B)$ を用いてベイズ則を適用することによって見出すことができる。P(D)およびP(A, B)は公のデータから入手可能である。特に、P(D)とは、集団における病気的全罹患率をいい、これは公に入手可能な統計学から見出すことができる。加えて、P(A, B)とは、個体において一緒に起こる遺伝子AおよびBの特定の状態の罹患率をいい、これは、異なる人種群における多数の個体での測定された多くの異なるSNPを有するHapMap Projectのような公のデータベースから見出すことができる。好ましい実施形態においては、これらの確率の全ては、全ヒト集団についてよりはむしろ、確率バイアスがある、特定の人種群および特定の性別について計算することができることを注記する。一旦、これらの確率が決定されたならば、挑戦は正確にP(A, B/D)を見積もることから由来する。というのは、コホートデータの大部分はP(A/D)およびP(B/D)の見積もりを供するからである。関連情報は、異なる遺伝子対立遺伝子間の統計学的関連についての、すなわち、P(A/B)についての、HapMap Projectのような種々の公のデータベースで見出すことができる。しかしながら、P(A/B)、P(A/D)、P(B/D)のみを仮定し、依然として、P(A, B/D)については何も言うことができない。というのは、拘束されない自由度があるからである。それにも拘わらず、もしなんらかの情報が、(A-, B-)のような丁度単一偶発事象についてさえ、遺伝子AおよびB双方をそれにつきサンプリングしたコホートからのP(A, B/D)について知られていれば、P(A/D)、P(B/D)、P(A/B)についての情報の価値を利用して、P(A, B/D)の見積もりを改良することができる。この概念は、偶発事象の表を用いて説明されるであろう。

【0371】

遺伝子状態A+およびA-に従う結果D+およびD-の確率を表す以下の2つの偶発事象表を考える。この実験はAと言及される。Aについての測定された頻度はFと言及され、見積もりを求める現実の確率はpを伴って言及される。

【0372】

【数88】

A	D+	D-
A+	f ₁	f ₂
A-	f ₃	f ₄

A	D+	D-
A+	p ₁	p ₂
A-	p ₃	p ₄

10

20

30

40

50

ここで、 $f_3 = 1 - f_1$ 、 $f_4 = 1 - f_2$ および $p_3 = 1 - p_1$ 、 $p_4 = 1 - p_2$ である。 K_1 が、A についての場合の群における対象の数、すなわち、結果 D + を有する対象の数を表すものとする。 K_2 が、A についての対照群における数、すなわち、結果 D - を有する対象の数であるとする。

【0373】

同様に、遺伝子状態 B + および B - に従う結果 D + および D - の確率を表す以下の2つの偶発事象の表を考える。この実験は B と言及される。測定された頻度は f を伴って言及され、見積もりを求める現実の確率は p を伴って言及される。

【0374】

【数89】

B	D+	D-
B+	f_5	F_6
B-	f_7	F_8

B	D+	D-
B+	p_5	p_6
B-	p_7	p_8

ここで、 $f_7 = 1 - f_5$ 、 $f_8 = 1 - f_6$ および $p_7 = 1 - p_5$ 、 $p_8 = 1 - p_6$ である。 K_3 が B についての場合の群における数を表すものとし、 K_4 が B についての対照群における数であるとする。前記偶発事象の表は、遺伝子状態 A および B が別々に測定される試験を表す。しかしながら、理想的に求められる偶発事象表は、組み合わせられた A および B の異なる状態を含む。偶発事象表は、A B という仮定実験について以下に示し、ここで、 f は測定された確率を表し、および p は現実の確率を表す。

【0375】

【数90】

AB	D+	D-
A+B+	f_9	f_{10}
A+B-	f_{11}	f_{12}
A-B+	f_{13}	f_{14}
A-B-	f_{15}	f_{16}

AB	D+	D-
A+B+	p_9	p_{10}
A+B-	p_{11}	p_{12}
A-B+	p_{13}	p_{14}
A-B-	p_{15}	p_{16}

ここで、 $f_{15} = 1 - f_9 - f_{11} - f_{13}$ 、 $f_{16} = 1 - f_{10} - f_{12} - f_{14}$ および $p_{15} = 1 - p_9 - p_{11} - p_{13}$ 、 $p_{16} = 1 - p_{10} - p_{12} - p_{14}$ である。 K_5 が A B についての場合の群における数とし、 K_6 が A B についての対照群における数とする。

【0376】

表記方法目的では、 $K_7 = K_9 = K_5$ および $K_8 = K_{10} = K_6$ であることを注記する。従って、事実、群のサイズは：

【0377】

10

20

30

40

【数 9 1】

#	D+	D-
A	K_1	K_2
B	K_3	K_4
AB	K_5	K_6

である。

【0 3 7 8】

統計学の基本則を用いて、仮定偶発事象表 A B の細胞の間の依存性を強制することができる。この例においては、D + に対応する細胞について、以下の関係を強制することができる：

$$P(A + B - / D +) = P(A + / D +) - P(A + B + / D +),$$

$$P(A - B + / D +) = P(B + / D +) - P(A + B + / D +)$$

$$P(A - B - / D +) = 1 - P(A + / D +) - P(B + / D +) + P(A + B + / D +)$$

同様に、D - に対応する細胞については：

$$P(A + B - / D -) = P(A + / D -) - P(A + B + / D -)$$

$$P(A - B + / D -) = P(B + / D -) - P(A + B + / D -)$$

$$P(A - B - / D -) = 1 - P(A + / D -) - P(B + / D -) + P(A + B + / D -)$$

である。

前記偶発事象表中の表記方法を用い、および余分な最後の関係を残し、これらの関係は：

$$p_{11} = p_1 - p_9$$

$$p_{13} = p_5 - p_9$$

$$p_{12} = p_2 - p_{10}$$

$$p_{14} = p_6 - p_{10}$$

に移され、あるいは同等に、

$$p_1 = p_9 + p_{11}$$

$$p_2 = p_{10} + p_{12}$$

$$p_5 = p_9 + p_{13}$$

$$p_6 = p_{10} + p_{14}$$

に移される。

全ての関係をまとめると、 p_9, \dots, p_{16} に対する p_1, \dots, p_{16} の全ての依存性の表を以下に掲げる。値の間の依存性を得るために、列内の確率は、値 = 1 を有する行内の確率の合計であり、例えば、第一列は $p_1 = p_9 + p_{11}$ を与える。

【0 3 7 9】

10

20

30

【数 9 2】

	p ₉	p ₁₀	p ₁₁	p ₁₂	p ₁₃	p ₁₄	p ₁₅	p ₁₆
p ₁	1		1					
p ₂		1		1				
p ₃					1		1	
p ₄						1		1
p ₅	1				1			
p ₆		1				1		
p ₇			1				1	
p ₈				1				1
p ₉	1							
p ₁₀		1						
p ₁₁			1					
p ₁₂				1				
p ₁₃					1			
p ₁₄						1		
p ₁₅							1	
p ₁₆								1

10

20

頻度および確率の関係から、 $n = 9 \dots 16$ についての測定方程式 $f_i = p_i + n_i$ を作成でき、ここで、 n_i は、出現 f_i の頻度に基づいた確率 p_i の不完全な測定を表すノイズ項である。前記した関係にこれを適用し、かつ偶発事象 A B の細胞の全ては測定されていると仮定し（これは、丁度説明目的のためであり、以下に議論する）、これらの 10 の観察を表すことができる。

30

これらの測定方程式は：

$$F = X P + N$$

として行列表記方法で表すことができる。ここで、 $F = [F_1, \dots, F_{16}]^T$ 、 $P = [p_9, \dots, p_{16}]^T$ および $N = [n_9, \dots, n_{16}]^T$ であり、 X は前記表中に表した行列である。この行列方程式を用いて、8 つの未知の係数、 $p_9 \dots p_{16}$ 、を解くことができる。この特別な場合において、我々は、全てのパラメータ $p_9 \dots p_{16}$ について解く。もし我々が組み合わされた A, B 遺伝子について全ての測定を有しないならば、我々は、D+ についての少なくとも 1 つの測定、および D- についての 1 つの測定を必要とする。前記関係を仮定すれば、次いで、我々は表の残りを満たすことができる。言い換えれば、仮定実験 A B についての偶発事象表を埋めることができるためには、望ましくは、A および B の特定の状態が、D+ および D- の結果を有する対象について同時に測定される少なくとも 1 つの例がある。これは、なされた測定を表す行列 X について十分なランクを達成することを可能とし、従って、値 $p_9 \dots p_{16}$ を解き、偶発事象表 A B に満たす。もしより多くの実験データが存在すれば、さらなる列を、前記で示したのと同様な構造を持つ行列 X の底部に加えることができる。

40

【0 3 8 0】

正確な回帰を行うためには、群試料のサイズによって決定される各観察 f_i についての重みを持つ重み付け回帰が望ましく、従って、さらに多くの観察を持つ実験および細胞は

50

より重み付けされる。測定方程式 $f_i = p_i + n_i$ では、 n_i は全て同一の偏差を有さず、回帰は等分散性でない。具体的には、 $f_i = 1/K_i * \text{Binomial}(p_i, K_i) \sim N(p_i, p_i(1-p_i)/K_i)$ であり、ここで、 $\text{Binomial}(p_i, K_i)$ は、各テストがケース結果 p_i の確率を有し、および K_i テストを行う二項分布を表す。この二項分布は $N(p_i, p_i(1-p_i)/K_i)$ によって近似することができ、これは平均 p_i および偏差 $p_i(1-p_i)/K_i$ を持つ二項分布である。この結果、ノイズは、理論的偏差 $V_i = p_i * (1-p_i)/K_i$ を有する正規変数 $n_i \sim N(0, p_i(1-p_i)/K_i)$ としてモデル化することができる。この偏差は、試料頻度 $v_i = f_i * (1-f_i)/K_i$ で近似することができる。

【0381】

10

偏差 v_i に逆比例する各観察 i についての重みを持つ重み付け回帰を行った。V が直交要素 $[v_9, \dots, v_{16}]$ を持つ行列であって、全ての他の要素は0である $\sim N(0, V)$ としてのノイズ行列 N の分布は、今や、記載することができる。これは $V = \text{diag}([v_9, \dots, v_{16}])$ として示される。同様に、 $W = \text{diag}([1/v_9, \dots, 1/v_{16}])$ とする。さて、重み付け回帰：

$$P = (X'WX)^{-1}X'WY$$

を用いて P について解くことが可能である。

P の偏差は

$$\text{Var}(P) = (X'WX)^{-1}$$

であることは直接的に示され、これを用いて、P の決定における信頼性を示すことができる。

20

【0382】

まとめると、我々は、A および B の組合せからのデータ (A B : f_9, \dots, f_{16}) と共に、個々の遺伝子からのデータ (A : f_1, \dots, f_4 , B : f_5, \dots, f_8) を用いて、A および B の組合せについての確率 (p_9, \dots, p_{16}) およびそれらの偏差 (v_9, \dots, v_{16}) を見積もるのを助けた。最後に、我々の研究においては、我々は、確率ではなく log オッズ比をほとんど取り扱い、従って、我々は、これらの確率を LOR に移す必要がある。一般に、事象 H について確率および偏差を以下のように仮定する。

【0383】

30

【数93】

	D+	D-
H+	p1	p2
H-	1-p1	1-p2
V	v1	v2

LOR についての式は (デルタ方法によって) 偏差を伴って、 $LOR = [\log(p_1) - \log(1-p_1)] - [\log(p_2) - \log(1-p_2)]$ である。 $V = [(p_1)^{-1} + (1-p_1)^{-1}]^{-2} * V(p_1) + [(p_2)^{-1} + (1-p_2)^{-1}]^{-2} * V(p_2)$ 。以下の表は、A, B の組合せについての確率、対応する LOR および偏差を示す。

40

【0384】

【数 9 4】

	D+	D-	LOR	Var
A+B+	p_9	p_{10}	lor_1	$V_1 = [1/p_9 + 1/(1-p_9)]^2 v_9 + [1/p_{10} + 1/(1-p_{10})]^2 v_{10}$
A+B-	p_{11}	p_{12}	lor_2	$V_2 = [1/p_{11} + 1/(1-p_{11})]^2 v_{11} + [1/p_{12} + 1/(1-p_{12})]^2 v_{12}$
A-B+	p_{13}	p_{14}	lor_3	$V_3 = [1/p_{13} + 1/(1-p_{13})]^2 v_{13} + [1/p_{14} + 1/(1-p_{14})]^2 v_{14}$
A-B-	p_{15}	p_{16}	lor_4	$V_4 = [1/p_{15} + 1/(1-p_{15})]^2 v_{15} + [1/p_{16} + 1/(1-p_{16})]^2 v_{16}$

10

これは、log オッズ比および各偏差の見積もりを提供する。

【0385】

この方法の説明として、該技術を使用して、 $P(A, B/D)$ の改良された見積もりが得られ、ここで、D はアルツハイマー病を有する状態を表し、およびここで、A および B は、各々、APOE および ACE 遺伝子の 2 つの異なる状態を表す。表 9 は、唯一の遺伝子 A がサンプリングされた 1999 年に Alvarez によって；唯一の遺伝子 B がサンプリングされた 1998 年に Labet によって；および遺伝子 A および B がサンプリングされた 2005 年 Farrer によって行われた 3 つの異なる実験を表す。結果の 2 つのセットはこれらの実験から作成されたものであり、表 10 に示す。最初のセット（表 10、行 2、3、4 および 5 参照）は、全てのコホートを分析し、本明細書中に開示された方法を用いて $P(A/D)$ $P(B/D)$ を仮定して $P(A, B/D)$ の見積もりを改良する。第二のセット（表 10、行 6、7、8 および 9 参照）は、 $P(A, B/D)$ についての Farrer (2005) の近代コホートから生じた結果のみを用い、そこでは、双方の遺伝子がサンプリングされた。前者の場合における予測の信頼限界は低下したと考えられる。これらの予測は、公の源からの $P(A/B)$ を記載するデータを用いてさらに改良することができ、これらの測定は前記したように X 行列に加えることができることを注記する。また、本明細書中に記載された技術を用いて、前記した $p_1 = p_5 + p_7$ のような関係を用い、 $P(A+/D+)$ 、 $P(A+/D-)$ 、 $P(B+/D+)$ および $P(B-/D-)$ のような別々の A、B 確率についての見積もりを改良することができることも注記する。

20

30

【0386】

この方法はただ 2 つの変数 A および B について説明してきたが、偶発事象の表は、アルツハイマー予測の関係で前記したもの：アルツハイマー病の家族履歴、性別、人種、年齢、および 3 つの遺伝子、すなわち、APOE、NOS3、および ACE の種々の対立遺伝子のような多くの異なる IV を含むことができることに注意すべきである。年齢のような連続的変数は、値のピンにカテゴリー化することによってカテゴリーを作成して、偶発事象表の処方に適当とすることができる。好ましい実施形態において、最大数を用いて、結果の確率をモデル化し、確率の標準偏差は、典型的には、いくつかの特定の閾値未満である。還元すれば、可能な最も特別な偶発事象は、その偶発事象についての十分な関連訓練データを維持して、関連する確率の見積もりを意味のあるものとしつつ、特定の患者に利用可能な IV を仮定して創製することができる。

40

【0387】

また、本開示を読んだ後に、病気 - 遺伝子関連、遺伝子 - 遺伝子関連、および / または集団における遺伝子頻度についてのデータを用いるために同様な技術をどのようにして適用して、多変数線形および非線形回帰および論理回帰モデルの精度を改良することができることは当業者に明らかであろうことを注記する。さらに、本開示を読んだ後に、病気 - 遺伝子関連、遺伝子 - 遺伝子関連、および / または集団における遺伝子頻度についてのデ

50

ータを用いるための同様な技術を適用して、どのように適用して、結果データの利用を可能として、モデルに関連するその全ての独立変数とその結果データにつき測定されるものではないモデルを訓練することによって、多変数線形および非線形回帰および論理回帰モデルの精度を改良することができることは当業者に明らかであろう。さらに、本開示を読んだ後には、病気 - 遺伝子関連、遺伝子 - 遺伝子関連、および/または集団における遺伝子頻度についてのデータを用いるための同様な技術をどのようにして適用して、当該分野で良く理解される期待値最大化 (EM) アルゴリズムのような他の技術を用いて形成された偶発事象表モデルの精度を改良することができるかは当業者に明らかであろう。これらの技術は、HapMap Projectからの活用データ、およびNational Center for Biotechnology Information (NCBI) Online Mendelian Inheritance in Man (OMIM) およびdbSNPデータベースのような公のデータベースに含まれる他のデータに特に関連する。

10

【0388】

また、当該特許を通じて、我々が個体または対象に関連するデータに言及する場合、これは、該対象に感染したかもしれないいずれの病原体または該対象に感染しつつあるいずれの癌の該データは言及できるとも仮定する。該個体または対象データは、ヒト胚、ヒト胚盤胞、ヒト胎児、いくつかの他の細胞または細胞のセットについてのデータ、あるいはいずれかの種類の動物または植物にも言及することができる。

【0389】

20

明日のデータ：回帰モデルでの多因子表現型のモデル化

より多くのデータが多因子表現型での遺伝子型に関連して蓄積されるにつれ、支配的なシナリオは前記した(iii)となり、すなわち、表現型を正確に予測するためには遺伝子マーカーの複雑な組合せを考慮するのが望ましく、多次元線形または非線形回帰モデルが導かれる。典型的には、このシナリオについてのモデルを訓練するにおいて、潜在的プレディクターの数は、測定された結果の数と比較して大きいであろう。本明細書中に記載されたシステムおよび方法の例は、未決定の、または悪い条件の遺伝子型 - 表現型データセットについての疎なパラメーターモデルを創製する新規な技術を含む。該技術は、それについて多くのモデリング業績が比較のために利用でき、およびそれについてデータが多くの潜在的遺伝子プレディクターに関連して入手可能な抗 - レトロウイルス療法 (ART) に対するHIV/AIDSの応答のモデル化に焦点を当てることによって説明される。現実の実験室測定で交差 - 確証によってテストする場合、これらのモデルは、文献中で以前に議論されたモデル、および本明細書中に記載された他のカノニカル技術よりも正確に薬物応答表現型を予測する。

30

【0390】

2つの回帰技術を、遺伝子配列データからの抗 - レトロウイルス療法に対する応答においてウイルス表現型を予測する関係で記載し、説明する。双方の技術は、モデルパラメーターの粗なセットの連続的サブセット選択のために凸最適化を使用する。最初の技術は、最小絶対収縮および選択オペレーター (LASSO) を用い、これは l_1 ノルム喪失関数を適用して、疎な線形モデルを作り出し；第二の技術は径ベースの核関数と共にサポートベクトルマシーン (SVM) を用い、これは、 l_2 - 非感受性喪失関数を適用して、疎な非線形モデルを創製する。該技術は、 l_0 の逆転写酵素阻害剤 (RTI) および7つのプロテアーゼ阻害剤薬物 (PI) に対するHIV-1ウイルスの応答の予測に適用される。遺伝子データは、逆転写酵素およびプロテアーゼ酵素についてのHIVコーディング配列に由来する。この性能を可能とするこれらのモデルの鍵となる特徴は、喪失関数が、パラメーターの多くがゼロである単純モデルを創製する傾向があり、およびコスト関数の凸性が、モデルパラメーターを見出して、特定の訓練データセットについてのコスト関数を全体的に最小化することができることを確実にすることである。

40

【0391】

LASSOおよび l_1 選択関数

50

プレディクター M の数が訓練試料の数 N を超える場合、モデル化の問題は過剰決定系、または不適切である。というのは、N のプレディクターのいずれかの任意のサブセット、X 行列における関連行が直線的に独立している限り、訓練データについてのゼロ誤差を持つ線形モデルを生じるのに十分だからである。その結果、線形回帰方法によって戻された N - プレディクターモデルに信頼を置く気がしない。しかしながら、N よりもかなり少数の変数が低い訓練誤差を有するモデルを仮定する。モデルがより疎であれば、低い訓練誤差は偶然人工物である確率は低く；よって、プレディクターが独立した変数に因果的に関連するのがよりありそうである。これは、RTI データの場合のように、過剰決定系の問題における疎な解の重要性の基礎となる。同様な議論を、PI データに当てはまるように、行列 $X^T X$ での大きな条件数によって特徴付けられる悪条件の問題に適用することができる。この場合、見積もられたパラメーター

10

【 0 3 9 2 】

【数 9 5】

$$\hat{b}$$

はモデル誤差に対して、ならびに測定ノイズに対して高度に感受性であり、結果として、正確に一般化されないようである。過剰決定系および悪条件の問題は、可能なプレディクター - 遺伝子、蛋白質、または我々の場合には、突然変異部位の数が、測定された結果の数に対して大きな遺伝子データに典型的である。

【 0 3 9 3 】

20

そのような場合に対する 1 つのカノニカルアプローチはサブセットの選択である。例えば、段階的選択にて、各工程において、その変数が予測誤差に相関する優位性のレベルを示す最高 F - 検定統計学を有することに基づいて、単一プレディクターをモデルに加える。各変数を加えた後、残りの変数を全てチェックして、モデルのプレディクター誤差とのそれらの関連性において統計学的有意性の閾値未満までそれらのいずれも降下しないことを確実にする。この技術は、薬物応答予測の問題に成功して適用されてきた。しかしながら、選択プロセスの区別される性質のため、データの小さな変化はプレディクターの選択されたセットをかなり改変することができる。1 つの変数の存在または不存在は、もう 1 つの変数と関連する統計学的有意性、およびその変数がモデルに含まれ、またはそこから拒絶されるかに影響し得る。これは、特に、悪条件の問題について一般化での精度に影響する。

30

【 0 3 9 4 】

もう 1 つのアプローチは、収縮関数によって拘束されるべき見積もられたパラメーター

【 0 3 9 5 】

【数 9 6】

$$\hat{b}$$

の値についてである。カノリカル収縮関数は該パラメーターの平方の合計であり、これは：

【 0 3 9 6 】

40

【数 9 7】

$$\hat{b} = \arg \min_b \|y - Xb\|^2 + \lambda \|b\|^2 \tag{17}$$

[式中、 λ は、典型的には、交差 - 確証によって決定されるチューニングパラメーターである]

に従ってパラメーターを見出す ridge 回帰において適用される。この方法は非疎であって、パラメーターを 0 に設定しない。これは、一般化における精度を低める傾向があり、解を解釈するのが困難とする。

【 0 3 9 7 】

これらの問題は LASSO 技術によって取り組まれる。サブセット選択とは対照的に、

50

LASSOはプレディクター変数の離散的許容または拒絶を行わず；むしろ、それは連続的サブセット最適化を介して、一緒になって最も効果的なプレディクターとなる変数のセットを一斉に選択することを可能とする。それは l_1 ノルム収縮関数：

【0398】

【数98】

$$\hat{b} = \arg \min_b \|y - Xb\|^2 + \lambda \sum_{i=1, \dots, M} |b_i| \quad (18)$$

[式中、 λ は典型的には交差 - 確証によって設定される]

を用いる。LASSOはパラメーターの多くを0に設定する傾向がある。図20は、選択性と言及されるLASSOのこの特徴に対する洞察を供する。丁度2つの突然変異に基づくモデルは訓練データ $X = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}^T$, $y = \begin{bmatrix} 2 & 1 \end{bmatrix}^T$ で創製され、 x -軸および y -軸は、各々、2つのパラメーター b_1 および b_2 を表す。 l_1 および l_2 収縮関数の使用を比較し、ここで、双方の場合において、 $\|y - Xb\|^2 = 2$ となるように訓練データを同等によくフィットさせる解が見出される。大きな円(2001)、小さな円(2002)、および(2003)は、各々、コスト関数 $\|y - Xb\|^2$ 、 l_2 ノルム $\|b\|^2$ 、および l_1 ノルム $|b_1| + |b_2|$ についてのレベル曲線を表す。ridge回帰に対する解(l_2)が見出され、ここで、2つの円が交差し(2004)；LASSOについての解(l_1)が見出され、ここで、四角および大きな円が交わる(2005)。 l_1 ノルムについてのレベル曲線の「尖性」のため、軸 b_1 上にある解が見出され、これは、従って、疎である。より高次元へ拡大されたこの議論は、疎な解を生じるLASSOの傾向を説明し、なぜ達成された結果が文献に報告されたものよりも測定可能に良好であるかを示唆する。

【0399】

l_1 ノルムは、凸でありつつ、最も選択的収縮関数として見ることができる。凸性は、所与のデータセットに対して1つの全体的解を見出すことができることを保証する。最小角回帰と言及されるかなり有効な最近のアルゴリズムは、 M 工程においてLASSOの全体的解に収束することが保証されている。

【0400】

本開示を読んだ後は、 l_1 ノルムをどのようにして論理回帰の関係で用いて、カテゴリ変数の各状態の確率をモデル化することもできることは当業者に明らかであろうことを注記する。論理回帰において、測定のセットの事後確率の逆数に対応する凸コスト関数を形成することができる。事後確率は、各結果の尤度のモデル見積もりを仮定する観測された訓練データの確率である。 l_1 ノルムを凸コスト関数に加えることによって、得られた凸コスト関数を最小化して、特定の結果の確率をモデル化するための疎パラメーターモデルを見出すことができる。論理回帰についての l_1 ノルムの使用は、測定された結果の数がプレディクターの数に対して小さい場合に、特に関連し得る。

【0401】

サポートベクトルマシンおよび l_1 -ノルム

SVMは、特に、モデルが独立変数の間の複雑な相互作用を含む場合に、薬物応答および他の表現型の良好なモデル化を達成するように構成することができる。SVMについての訓練アルゴリズムは、 l_1 ノルム選択関数の使用を黙示的とする。SVMは、現実の価値の関数近似を行うことができ、かつ見積もり問題がHadamardの意味において不適切である場合でさえ、試料データの正確な一般化を達成することができる学習アルゴリズムである。正確に一般化されるSVMの能力は、SVMモデルおよび訓練アルゴリズムにおける2つの選択可能な特徴によって典型的には影響される。第一のものはコスト関数、または訓練において最小化されるべき関数の選択である。第二のものは、SVMの核、または線形回帰パラメーターの比較的小さなセットを用いて、SVMが、独立変数の間の相互作用を含む複雑な非線形関数をマッピングするのを可能とする関数の選択である。これらの特徴は以下に議論する。

10

20

30

40

50

【 0 4 0 2 】

線形関数近似：

【 0 4 0 3 】

【数 9 9 】

$$\hat{y}_i = f(x_i, b) = b^T x_i$$

を持つ対象 i y_i についての表現型をモデル化することを考える。まず、いくらか > 0 未満の誤差にペナルティを与えない「 - 非感受性喪失」関数と共に、パラメーターでの l_2 収縮関数よりなるコスト関数を最小化することによって b を見積もる。S V 回帰を拘束：

【 0 4 0 4 】

【数 1 0 0 】

$$y_i - b^T x_i \leq \varepsilon + \xi_i^+, i = 1 \dots N \quad (20)$$

$$b^T x_i - y_i \leq \varepsilon + \xi_i^-, i = 1 \dots N \quad (21)$$

$$\xi_i^+ \geq 0, \xi_i^- \geq 0, i = 1 \dots N \quad (22)$$

を条件として、以下の最適化：

【 0 4 0 5 】

【数 1 0 1 】

$$\hat{b} = \arg \min_{b, \xi_i^-, \xi_i^+} \frac{\|b\|^2}{2} + C \sum_{i=1}^N (\xi_i^- + \xi_i^+) \quad (19)$$

として公式化することができる。

【 0 4 0 6 】

コスト関数の第二項は、「非感受性」閾値 ε を超えてモデル化誤差の絶対値を最小化する。パラメーター C は、誤差 v s 重みに対する収縮の相対的重要性を見積もることを可能とする。この拘束された最適化を、ラグランジュの鞍点を見出す標準的技術を用いて解いて、K u h n - T u c k e r 拘束を満足させることができる。前記したコストおよび拘束を適合させるラグランジュは：

【 0 4 0 7 】

【数 1 0 2 】

$$\begin{aligned} L(b, \xi_i^+, \xi_i^-, \alpha_i^+, \alpha_i^-, \lambda_i^+, \lambda_i^-) = & \frac{\|b\|^2}{2} + C \sum_{i=1}^N (\xi_i^- + \xi_i^+) - \sum_{i=1}^N \alpha_i^- (y_i - b^T x_i + \varepsilon + \xi_i^-) \\ & - \sum_{i=1}^N \alpha_i^+ (y_i - b^T x_i + \varepsilon + \xi_i^+) - \sum_{i=1}^N (\lambda_i^- \xi_i^- + \lambda_i^+ \xi_i^+) \end{aligned} \quad (23)$$

である。パラメーター b 、 ξ_i^- 、 ξ_i^+ のベクトルに関して最小化し、ラグランジュ乗数 α_i^- 、 α_i^+ 、 λ_i^- 、 λ_i^+ のベクトルに関して最小化する。ラグランジュ乗数は K u h n - T u c k e r 拘束に従って望ましくは正であることを注記する。よって、パラメーターの最適なセットは、

【 0 4 0 8 】

【数 1 0 3 】

$$\alpha_i^+, \alpha_i^-, \lambda_i^+, \lambda_i^- \geq 0, i = 1 \dots N \quad (25)$$

を条件として、

10

20

30

40

50

【 0 4 0 9 】

【 数 1 0 4 】

$$(b_*, \xi_*^+, \xi_*^-) = \arg \min_{b, \xi^+, \xi^-} \max_{\alpha^+, \alpha^-, \beta^+, \beta^-} L(b, \xi^-, \xi^+, \alpha^+, \alpha^-, \lambda^+, \lambda^-) \quad (24)$$

に従って見出すことができる。最小化 / 最大化の順序は相互交換できるので、これらの変数に関する L の部分的導関数を 0 に設定することによって、変数 b 、 α_i^- 、 α_i^+ に関してまず最小化する。得られた方程式から、重みベクトルを

【 0 4 1 0 】

【 数 1 0 5 】

$$b = \sum_{i=1}^N (\alpha_i^+ - \alpha_i^-) x_i \quad (26)$$

の項で表すことができることが判明する。また、得られた方程式から、

【 0 4 1 1 】

【 数 1 0 6 】

$$\sum_{i=1}^N \alpha_i^+ = \sum_{i=1}^N \alpha_i^- \quad (28)$$

$$0 \leq \alpha_i^+ \leq C, i=1 \dots N \quad (29)$$

$$0 \leq \alpha_i^- \leq C, i=1 \dots N \quad (30)$$

を条件として、二次形式：

【 0 4 1 2 】

【 数 1 0 7 】

$$W(\alpha^+, \alpha^-) = -\sum_{i=1}^N \varepsilon(\alpha_i^- + \alpha_i^+) + \sum_{i=1}^N y_i (\alpha_i^+ - \alpha_i^-) - \frac{1}{2} \sum_{i,j=1}^N (\alpha_i^- - \alpha_i^+) (\alpha_j^- - \alpha_j^+) x_i^T x_j \quad (27)$$

を最大化することによって、係数 α_i^+ 、 α_i^- 、 $i = 1 \dots N$ を見出すことができるように、ラグランジュから変数を排除する。

【 0 4 1 3 】

これは、ベクトル b が計算されるのを可能とし、 ξ - 非感受性喪失関数に対する SVM モデルを十分に定義する。方程式 (1 1) から、モデルは、

【 0 4 1 4 】

【 数 1 0 8 】

$$f(x) = \sum_{i=1}^M \beta_i (x^T x_i) + b_0 \quad (31)$$

[式中、 $\beta_i = \alpha_i^+ - \alpha_i^-$]

として特徴付けることができることを注記する。得られたモデルは、セット { x_i , $i = 1 \dots M$ } 中のパラメーターの多くが 0 となる点で、疎となる傾向があろう。非ゼロの値 β_i に対応するベクトル x_i はモデルのサポートベクトルとして知られている。サポートベクトルの数は、チューナブルパラメーター C の値、訓練データ、およびモデルの適当性に依存する。以下の説明においては、今や、どのようにしてモデルを増加させて、核関数の使用でもって複雑な非線形関数を適合させることができるかを示す。次に、 ξ - 非感受性喪失関数は $1 - \exp(-\gamma \|x\|)$ ノルム収縮関数に関し、それは同じこと、すなわち、 $1 - \exp(-\gamma \|x\|)$ ノルムによる疎なパラメーターセットの一斉選択を実質的に達成することが示されるであろう。

【 0 4 1 5 】

10

20

30

40

50

変数の間で結合が可能な複雑な関数をモデル化するためには、方程式(17)の単純な内積を、ベクトルの間のより複雑な相互作用を計算する核関数で置き換える。核関数を挿入し、(17)中の我々の関数の近似は形態：

【0416】
【数109】

$$f(x) = \sum_{i=1}^N \beta_i K(x, x_i) + \beta_0 = \sum_{i=0}^N \beta_i K(x, x_i) \quad (32)$$

[式中、定義によると $K(x, x_0) = 1$ である]

を採る。これらのパラメータを見出すためには、前記したのと正確に同一の最適化方法を用い、全ての項 $x^T x_i$ を $K(x, x_i)$ で置き換える。前記したように、前記したのと同一の拘束に従い、

【0417】
【数110】

$$W(\alpha^+, \alpha^-) = -\sum_{i=1}^N \varepsilon (\alpha_i^- + \alpha_i^+) + \sum_{i=1}^N y_i (\alpha_i^+ - \alpha_i^-) - \frac{1}{2} \sum_{i,j=1}^N (\alpha_i^- - \alpha_i^+) (\alpha_j^- - \alpha_j^+) K(x_i^T x_j) \quad (33)$$

を最大化する独立変数を見出すことによって、 $\alpha_i = \alpha_i^+ - \alpha_i^-$ に従ってパラメータセットを計算する。前記したSVM結果では、径基礎核関数を選択した。

【0418】

さて、 l_1 ノルムの黙示的使用を説明するために：方程式(17)を最適化することを試みる代わりに、最適化：

【0419】
【数111】

$$\beta^* = \arg \min_{\beta} \int_{-\infty}^{\infty} \left(f(x) - \sum_{i=0}^N \beta_i K(x_i, x) \right)^2 dx + \varepsilon \sum_{i=0}^N |\beta_i| \quad (34)$$

で開始し、ここで、 l_1 収縮を明示的に用いて、 β_i の値を拘束させてあり、訓練データの離散的試料に対して定義される代わりに、データフィッティング誤差を、モデル化すべき仮定関数のドメインに対して定義する。さて、変数置換： $\alpha_i = \alpha_i^+ - \alpha_i^-$ ； $\alpha_i^+ \geq 0$ ， $\alpha_i^- \geq 0$ ， $\alpha_i^+ - \alpha_i^- = \alpha_i$ ， $i = 1 \dots N$ を行う。次いで、拘束

【0420】
【数112】

$$\alpha_i^+, \alpha_i^- \geq 0 \quad (36)$$

$$\alpha_i^+ \alpha_i^- = 0 \quad (37)$$

に従い、

【0421】
【数113】

$$W(\alpha^+, \alpha^-) = -\sum_{i=1}^N \varepsilon (\alpha_i^- + \alpha_i^+) + \sum_{i=1}^N y_i (\alpha_i^+ - \alpha_i^-) - \frac{1}{2} \sum_{i,j=1}^N (\alpha_i^- - \alpha_i^+) (\alpha_j^- - \alpha_j^+) K(x_i, x_j) \quad (35)$$

として、最適化を書き直すことができる。異なる拘束を有するこの解は、それにも拘わらず、もしSV方法についての値Cが、拘束 $0 \leq \alpha_i^+, \alpha_i^- \leq C$ が単純に拘束(21)および(22)となるのに十分に大きく選択されるならば、 l_1 -非感受性喪失関数に一致し、また、基礎関数の1つは、我々の場合についての方程式(17)におけるように定数

10

20

30

40

50

である。この場合、S V方法によって用いられるさらなる拘束

【0422】

【数114】

$$\sum_{i=1}^N \alpha_i^+ = \sum_{i=1}^N \alpha_i^-$$

を必要としない。拘束(25)は既に方程式(15)において黙示的であることを注記する。というのは、拘束(8)および(9)は同時に活性となることはできず、従って、ラグランジェ乗数 α_i^+ または α_i^- のうちの1つはスラックであるか、または0であるべきであるからである。

10

【0423】

これらの条件下で、 $1 - \alpha_i^-$ 収縮関数のアプローチを黙示的に用いて、 α_i^- - 非感受性喪失関数が疎な関数近似を達成することを見ることができる。

【0424】

多因子表現型予測：HIV-1薬物応答のモデリングの例

サルベージARTの表現型結果の予測に対する現在のアプローチは、主として、薬物養生法および遺伝子突然変異の多くの異なる順列と組み合わせた、主として、統計学的に有意な結果データの欠如のため、良好な予測パワーを示さない。この分野は、多数の不均一データセットの統合、および薬物応答予測の増強の双方についての逼迫した必要性を有する。

20

【0425】

本明細書中で示されたモデルは、訓練およびテスト目的のためのStanford HIVdb RTおよびプロテアーゼ薬物耐性データベースからのデータを用いた。このデータは、逆転写酵素(RT)またはプロテアーゼコーディングセグメントが配列決定されているHIV-1ウイルスの6644イン・ビトロ表現型テストよりなる。テストは、10の逆転写酵素阻害剤(RTI)および7つのプロテアーゼ阻害剤(PI)について行われた。RTIはラミブジン(3TC)、アバカビル(ABC)、ジドブジン(AZT)、スタビジン(D4T)、ザルシタピン(DDC)、ジダノシン(DDI)、デラビラジン(DLV)、エバピレンズ(EFV)、ネビラピン(NVP)およびテノフォビル(TDF)を含む。PIはアムプラナビル(APV)、アナザナビル(ATV)、ネルフィナビル(NFV)、リトナビル(RTV)、サキナビル(SQV)、ロピナビル(LPV)およびインジナビル(IDV)を含む。

30

【0426】

各薬物については、データは形式 (x_i, y_i) , $i = 1 \dots N$ の対に構造化されており、ここで、 N は訓練データを構成する試料の数であり、 y_i は測定された薬物の倍耐性(または表現型)であって、 x_i は突然変異のベクトル+定数、 $x_i = [1, x_{i1}, x_{i2}, \dots, x_{im}]^T$ であり、ここで、 M は関連酵素についての可能な突然変異の数である。もし m 番目の突然変異が i 番目の試料に存在すれば要素 $x_{im} = 1$ であり、その他の場合 $x_{im} = 0$ に設定する。各突然変異はコドン遺伝子座および置換されたアミノ酸の双方によって特徴付けられる。アミノ酸配列に影響しない突然変異は無視する。各薬物についての試料に1%を超えて存在する突然変異のみがモデルについての可能なプレディクターのセットに含まれることを注記する。というのは、耐性に関連する突然変異はあまり頻繁でなく起きるのはありそうもないからである。測定 y_i は野生型と比較した突然変異ウイルスについての薬物の倍耐性を表す。具体的には、 y_i は、野生型ウイルスの IC_{50} と比較した、突然変異したウイルスの IC_{50} (複製を50%だけ遅らせるのに必要な薬物の濃度)の比率の \log である。目標は、 x_i から y_i を正確に予測する各薬物についてのモデルを開発することである。データに対してバッチ最適化を行うためには、 $N \times M + 1$ 行列、 $X = [x_1, x_2, \dots, x_N]^T$ に独立変数をスタックし、ベクトル $y = [y_1, y_2, \dots, y_N]^T$ に全ての観察をスタックする。

40

【0427】

50

各アルゴリズムの性能は交差 - 検証を用いて測定する。各薬物については、一次相関係数 R を、モデルの予測された表現型応答、およびテストデータの現実の測定されたイン・ビトロ表現型応答の間で計算する。

【 0 4 2 8 】

【数 1 1 5 】

$$R = \frac{(\hat{y} - \bar{y}\vec{1})^T (y - \bar{y}\vec{1})}{\|\hat{y} - \bar{y}\vec{1}\|_2 \|y - \bar{y}\vec{1}\|_2} \quad (38)$$

ベクトル

【 0 4 2 9 】

【数 1 1 6 】

\hat{y}

10

が表現型 y の予測である場合、

【 0 4 3 0 】

【数 1 1 7 】

\bar{y}

はベクトル y における要素の平均を示し、

【 0 4 3 1 】

【数 1 1 8 】

$\vec{1}$

20

は全てのもののベクトルを示す。各薬物および各方法については、各々、訓練およびテストのためにデータを比率 9 : 1 にランダムに細分化する。1つの例において、10の異なる細分化を行って、訓練およびテストデータのいずれの重複もなくしてベクトル

【 0 4 3 2 】

【数 1 1 9 】

\hat{y}

30

および R を得る。次いで、この全プロセスを 10 回反復して、 R の 10 の異なる値を得ることができる。 R の 10 の異なる値を平均して、報告された R を得る。また、10 の異なる実験にわたって測定されたモデルの各々について R の標準偏差を決定して、モデルが統計学的に有意な方法で比較されることを確実にする。

【 0 4 3 3 】

表 1 1 は P I 薬物についての前記したモデルの結果を示し ; 表 1 2 は 1 0 の R T I 薬物についての結果を示す。結果は、訓練およびテストデータの 1 0 細分化にわたって平均した、修正係数 R の形式で示す。試料偏差から計算した R の平均値の見積もった標準偏差も示す。各薬物についての利用可能な試料の数を最後の列に示す。平均性能を増加させるためにテストした方法は : i) R R - R i d g e 回帰、i i) D T - 検出ツリー、i i i) N N - 神経ネットワーク、i v) P C A - 主成分分析、v) S S - 段階的選択、v i) 直線核での S V M _ L - サポートベクトルマシン、v i i) L A S S O - 最小絶対収縮および選択オペレーター、および v i i i) 径基礎核での S V M - サポートベクトルマシンである。表 1 1 および 1 2 の最後の行中の情報を図 2 1 に示す。図 2 1 中の円は、各 P I についての 1 0 の異なる実験にわたって平均し、かつ 7 の異なる P I にわたって平均した相関係数 R を示す。図 2 1 中の菱形は、各 R T I についての 1 0 の異なる実験にわたって平均し、かつ 1 0 の異なる R T I にわたって平均した相関係数 R を呈する。1 標準偏差誤差棒も示す。

40

【 0 4 3 4 】

50

モデル化技術がチューニングパラメーターを含む場合は常に、これらは、グリッドサーチアプローチを用い、交差 - 確証によって測定されたように技術の最適性能のために調整されている。全ての場合において、グリッド量子化は、グリッドからの最良の実行パラメーターが所与のデータについての最適パラメーターから現実的には識別可能であるのに十分良好であった。というのは、グリッド量子化による予測の差は実験ノイズを低めるからである。

【0435】

データには強い傾向があるが、試料の数の差のため、基礎となる遺伝子プレディクター、および薬物間で変化するデータ中の他の特異性、各アルゴリズムによって達成されるRの相互作用は薬物間で変化する得ることは注意すべきである。この変動は、表11(3ないし9行)および表12(3ないし12行)の個々の薬物行を調べることによって見ることができる。

10

【0436】

全ての方法のうち、SVMは良好に実行され、LASSOを僅かに凌ぐ(RTIについて $P < 0.001$; PIについて $P = 0.18$)。 - 非感受性喪失関数で訓練したSVMの性能は、サポートベクトルマシンに基づいた従前に報告された方法のそれよりもかなり良好である。非線形核関数を用いるSVMは、線形核関数を用い、および - 非感受性喪失関数を用いても訓練されるSVM Lを凌ぐ(RTIについて $P = 0.003$; PIについて $P < 0.001$)。SVMは、神経ネットワークを用い、かつ凸コスト関数および連続的サブセット選択を創製しない他の非線形技術をかなり凌ぐ(RTIおよびPI双方について $P < 0.001$)。凸コスト関数を用いて線形回帰モデルを訓練し、LOSSO技術は、SS技術をかなり凌ぐ(PIおよびRTI双方について $P < 0.001$)。トップの5つの方法、すなわち、SS、PCA、SVM_L、LASSO、SVM_Rは、全て、疎であるモデルを創製する傾向があるか、または限定された数の非ゼロパラメーターを有する。

20

【0437】

プレディクターとして選択された突然変異のサブセットを説明するために、本明細書中に開示されたある実施形態は第二の最良の実行モデル、すなわち、SVMとは異なり、プレディクターの間の非線形または論理的結合を模倣することを試みない線形回帰モデルを創製するLASSOに焦点を当てる。結果として、どのようにして多くのプレディクターを選択するかを示すのは直接的である。表13は、各モデルを訓練するにおいて用いられる、突然変異の数(表13、3列)、および試料の合計数(表13、2列)と共に、各PI薬物についてのプレディクターとしてのLASSOによって選択された突然変異の数(表13、4列)を示す。同一の表が、RTIについて示される(表14、同一列は同一事項に対応する)。

30

【0438】

選択された突然変異もまた薬物耐性の原因の理解を高めることができる。図22、23および24は、各々、PI、ヌクレオチドRTI(NRTI)、および非ヌクレオチドRTI(NNRTI)に対する応答を予測するためにLASSOによって選択されたパラメーターの値を示す。図面中の各列は薬物を表し;各行は突然変異を表す。関連突然変異はPI薬物についてはプロテアーゼ酵素に対する、およびRNTIおよびNNRTI薬物についてはRT酵素に対するものである。各四角の陰影は、その薬物についてのその突然変異に関連するパラメーターの値を示す。右側の色付き棒線(各々、2201、2301および2401)によって示されるように、陰影を付したダーカーであるプレディクターは増大した耐性に関連し;陰影を付したライターであるパラメーターは増大した感受性に関連する。突然変異は、関連パラメーターの平均の大きさを減少させる順序で左側から右側の順序とする。関連パラメーターをクラスにおいて全ての列または薬物にわたって平均する。40の最大のパラメーターの大きさに関連する突然変異を示す。特定の突然変異、または行については、パラメーターの値は、列、または同一クラスにおける異なる薬物にわたってかなり変化する。

40

50

【0439】

アルゴリズムRR、DT、NN、およびSSについては、モデルは、全ての遺伝子突然変異についてではなく、むしろDepartment of Health and Human Services (DHHS)によって耐性に影響すると考えられる部位で起こる突然変異のサブセットについて訓練した。独立変数の数の低下は、これらのアルゴリズムの性能を改良することが判明した。SVM_Lアルゴリズムの場合には、全ての突然変異についてモデルを訓練することによってPIに対する最良の性能を達成しつつ、DHHS突然変異サブセットのみを用いてRTIに対する最良の性能を達成した。全ての他のアルゴリズムについては、最良の全性能は、全ての突然変異についてモデルを訓練することによって達成された。

10

【0440】

プレディクターとしてのLASSOによって選択されたが、現在、耐性に影響するとDHHSによって判断された遺伝子座と関連付けられていない図22、23、および24に示された突然変異のセットは：PIについては-19P、91S、67F、4S、37C、11I、14Z；NRTIについては-68G、203D、245T、208Y、218E、208H、35I、11K、40F、281K；およびNNRTIについては-139R、317A、35M、102R、241L、322T、379G、292I、294T、211T、142Vである。LASSOおよびSVMのようないくつかの場合においては、LPVのような特定の薬物についての性能は、DHHSによって耐性に影響すると認識された遺伝子座のみが含まれた場合($R = 81.72$, $Std. dev. = 0.18$)と比較して、全ての突然変異がモデルに含まれた場合($R = 86.78$, $Std. dev. = 0.17$)、有意に改良された($P < 0.001$)ことを注記する。これは、DHHSによって認識されたものを超えた他の突然変異が薬物耐性において役割を演じることができることを説明する。

20

【0441】

凸最適化技術の使用は、本明細書中において、疎なパラメータセットの連続的サブセット選択を達成して、正確に一般化される表現型予測モデルを訓練することが示された。LASSOは、 l_1 ノルム収縮関数に適用して線形回帰パラメータの疎なセットを生じる基底関数でのかつ - 非感受性喪失関数で訓練したSVNは疎な非線形モデルを創製する。これらの技術の優れた性能は、それらのコスト関数の凸性、および疎なモデルを生じるそれらの傾向の点で、説明することができる。凸性は多くの潜在的プレディクターがある場合に、特定の訓練データセットについて全体的に最適なパラメータを見出すことができるのを確実とする。疎なモデルは、遺伝子データに典型的なように、特に劣決定または悪条件データの関連でよく一般化される傾向がある。 l_1 ノルムは、最も選択的な凸関数として見るることができる。選択的収縮関数を用いる疎なパラメータセットの選択は、Occam's Razorと同様な原理で：多くの可能な理論が観察されたデータを説明できる場合、最も単純なものは最も正しいようである：を発揮する。 - 非感受性喪失関数と共に l_2 収縮関数を用いるSVMは、サポートベクトルと関連するパラメータに適用された収縮関数として l_1 ノルムの明示的な使用と同様な効果を生じる傾向がある。

30

40

【0442】

l_1 収縮関数を用いる技術は、しばしば、IVの数が大きくて、データが未決定または悪条件である場合、正確に一般化することができる。結果として、独立変数の非線形または論理的組合せをモデルに加え、良好なプレディクターである組合せを訓練で選択されると予測することが可能である。SMVは、線形核関数よりも有意に良好に実行される、基底関数のような非線形核関数の使用と独立変数との相互作用をモデル化することが可能である。結果的に、本明細書中に開示した基本的な概念を変えないことなく、独立変数の論理的組合せをモデルに加えることによって、LASSOの性能を高めることができる。論理的項は、決定ツリーによって生じたものから、専門家則によって記載された論理的相互作用から、論理的回帰の技術から、または論理的項のランダム順列のセットさえから由来

50

することができる。LASSOの利点は、パラメーターが、サポートベクトルよりはむしろ独立変数、または独立変数を含む表現を直接的に組み合わせるので得られるモデルが解釈するのが容易であることである。モデルにおける多数の独立変数に対するLASSOの頑強性は、 l_1 ノルムの選択的性質およびその凸性双方によるものである。

【0443】

l_1 ノルムよりも収縮関数をより選択的に使用する他の技術が存在する。例えば、log-収縮回帰は、モデルパラメータセットに存在する情報の量を測定する暗号理論に由来する収縮関数を用いる。この技術は l_1 -ノルムの代わりに収縮関数としてlog関数を用い、その結果、非凸である。パラメーターの疎なセットを求めるための理論的に興味があるアプローチを供しつつ、ペナルティ関数の非凸性は、対応する回帰を解くことが、LASSOよりも依然として計算の、扱いやすくなく、プレディクターの大きなセットについては所与のデータについての全体的最小よりはむしろ局所的な最小のみを生じさせることができることを意味する。

10

【0444】

本明細書中に記載された技術は、広い範囲の表現型予測問題についての線形および非線形回帰モデルの創製に適用することができる。それらは潜在的遺伝子プレディクターの数が測定された結果の数と比較して大きい場合に特に関連する。

【0445】

遺伝子独立変数を異なる空間へマッピングすることによる回帰モデルの単純化

前記したように、遺伝子マーカーの複雑な組合せを考える場合、SNP変数をもう1つの変数空間に投影して、分析を単純化することが可能であることを注記する。この変数空間は、HapMap Projectによって記載されたクラスターまたはピンのような、突然変異の公知のパターンを表すことができる。言い換えれば、前記した特定のSNP突然変異を表すベクトル x_i よりはむしろ、それは、個体が特定のHapMapクラスターまたはピンに入るか否かを表すことができる。例えば、前記した表記方法に従い、Bが関連HapMapピンの数であるベクトル $x_i = [x_{i1}, x_{i2}, \dots, x_{ib}]^T$ があると想像する。もし個体のSNPパターンがb番目のピンに入るならば、要素 $x_{ib} = 1$ を設定し、そうでなければ0を設定することができる。別法として、もし個体SNPおよび特定のピンの間の重複が不完全であって、カテゴリー「他の」において単純に個体を置き換えるのが望ましくないのであれば、各 x_{ib} を、SNPのパターンおよびピンbのその間の重複の割合と等しく設定することができる。本明細書中に開示された概念を変えることなく多くの他の技術は回帰問題を公式化することが可能である。

20

30

【0446】

結果予測についての交差確認によるモデルの選択

この議論を進めた中で、専門家則、偶発事象表、線形および非線形回帰を含む異なる表現型予測技術を記載した。さて、訓練データの使用に基づき、特定の対象についての特定の Kategorie または非 Kategorie 結果をモデル化するのが最良であるモデル化技術のセットから選択する一般的アプローチを記載する。図25は、システムについての説明的フローダイアグラムを供する。図25に記載されたプロセスは、特定の患者、モデル化すべき表現型、およびデータをテストし訓練する所与のセットで利用できるデータを仮定して最良のモデルを選択する一般的アプローチであり、該プロセスは特定のモデル化技術から独立している。好ましい実施形態において、用いることができるモデル化技術のセットは、専門家則、偶発事象表、LASSOで、またはデータが劣決定されていない場合は単純な最小二乗で訓練された線形回帰モデル、およびサポートベクトルマシンの非線形回帰モデルを含む。

40

【0447】

該プロセスは、モデル化されるであろう、あるいはもしそれが Kategorie 変数であれば、それについて確率をモデル化することができる、特定の対象および特定の従属変数(DV)を選択で開始する2501。次いで、該システムは、対象の記録に関連し、かつDVの結果のモデル化に関連し得る独立変数(IV)のセットを決定する2502。システム

50

のヒトユーザーは、ユーザーがモデルに関連して可能と考えるIVのそのサブセットを選択することもできる。次いで、システムはチェックして2503a、モデルが既に訓練され、独立変数の所与の組合せ、およびモデル化すべき所与の従属変数について選択されているか否かをみる。もしこれが当てはまり、かつ出来合いのモデルを訓練し、テストするのに用いるデータが旧式でなければ、システムは、そのモデルを用いる予測の創製に直接的に向かう2519。そうでなければ、システムは、注目する特定のDVを有し、かつ注目する特定の対象と同一のIVのセットを有しても有しなくてもよい。全ての他の記録をデータベースから抽出するであろう。そうすることにおいて、システムは、データがモデルを訓練しテストするのに利用できるか否かを決定する2503b。もし答えが否であれば、システムは、いずれかの利用可能な専門家則があるかをみるためにチェックして2515、対象で利用可能なIVのサブセットに基づいて結果を予測する。もし専門家則が利用できなければ、システムは出て2504、それが有効な予測をできないと示す。もし1以上の専門家則ができれば、システムは、特定の対象のデータに最良に適する専門家則のサブセットを選択する2505。好ましい実施形態において、対象にいずれの専門家則を適用するかを選択は、その専門家則見積もりにおける信頼性のレベルに基づくであろう。もしそのような信頼性見積もりが利用できなければ、それらの特異性のレベルに基づいて、すなわち、注目する対象で利用できるどれくらい多くのIVを専門家則が予測で用いるかに基づいてランク付けすることができる。次いで、専門家則の選択されたサブセットを用いて予測を生じさせる2506。

【0448】

もしデータが利用できると判断されたならば2503b、システムはチェックして2516、テストおよび訓練データで失われたいずれかのデータがあるか否かを決定する。言い換えれば、関連DVを含む全ての記録について、システムはチェックして、全ての記録が、注目する患者について利用できるのと正確に同一のIVのセットを有するか、およびいずれがモデルにおいて潜在的予測であり得るかをチェックする。典型的には、答えは「否」であろう。というのは、異なる情報が異なる患者で利用可能だからである。もし失われたデータがあればシステムは四方を進んで対象にとって最良の可能な予測をなすのに用いるべきIVのセットを見出す。この手法は時間を消費するものである。というのは、それは多数ラウンドのモデル訓練および交差 - 確証を含むからである。その結果、この手法における最初の工程は、考えられるIVのセットを、利用可能な計算時間に基づいて管理可能なサイズに低下させることである2507。好ましい実施形態においてIVのセットは、やはり利用可能なDVを有する対象のあるパーセンテージについてのそのIVに関するデータがあることに基づいて低下させる。単純な線形回帰モデルを仮定し、それらはモデル化誤差に関連する程度に基づいてIVを選択する段階的選択のような当該分野で知られた他の技術を用いて、IVのセットをさらに低下させることができる。次いで、システムはループに入り、そこでは、残りのIVの各組合せが調べられる。好ましい実施形態において、各IVおよびDVについても以下の状態を考慮する：各IVはモデルに含めることができるか、または含めることができず、全ての対象について陽性であるIVまたはDVについての数値データでは、該データはその対数を取ることによって進行させても、させなくてもよい。IVの包含/排除および前処理の各特定の組合せについてモデル化技術のセットを適用する2510。

【0449】

ほとんどのモデル化技術は、テストデータでの交差 - 確証を用いるグリッド - サーチアプローチに基づいて最適化し、またはチューニングすることができるいくつかのチューニングパラメーターを有するであろう。例えば、先に議論したLASSO技術については、多くの値が変数パラメーター について調べられる。 の各値について、回帰パラメーターを訓練することができ、モデルの予測をテストデータの測定された値と比較することができる。同様に、先に議論したサポートベクトルマシンアプローチでは、グリッド - サーチアプローチを用いて最適化すべきチューニングパラメーターはC、 、および、おそらくは核関数の特徴を記載するパラメーターを含む。偶発事象表に基づいた技術ではチュ

10

20

30

40

50

ーナブルパラメーターは、先に議論したように、偶発事象を所与の対象についてできるだけ特異的としつつ、偶発事象表モデルから許容できる最高の標準偏差と比較することができる。

【0450】

多くの異なる行列を用いて、モデル予測をテストデータと比較して、チューナブルパラメーターを最適化し、モデルを選択することができる。好ましい実施形態において、誤差の標準偏差を用いる。他の実施形態において、予測されたおよび測定された結果の間の相関係数Rを用いることができる。論理的回帰または偶発事象表の関係で、事後確率、すなわち、各テスト結果の尤度のモデルの予測を仮定するテストデータの所与のセットの確率を用いることもできる。いずれの測定基準を用いようとも、もし予測誤差の標準偏差をテスト測定基準として用いるならば、予測誤差の標準偏差の最小化のような、測定基準の値を最適化するチューニングパラメーターのその値を選択する。モデル訓練および交差 - 確証はゆっくりとしたプロセスであるので、この段階2510において、異なるチューニングパラメーターが調べられるように規定するグリッドは、最良のモデルおよび最良のチューニングパラメーターの粗いアイデアのみを得ることができるように、利用可能な時間の量に基づいておおまかセットされる。

【0451】

一旦、全ての異なるIV/DV組合せがこのようにして調べられたならば2511、システムが、テスト測定基準の最良の値を達成した、IV/DVの組合せ、モデルおよびチューニングパラメーターを選択する。もし失われたデータがなければ、システムはIV/DVのすべての組合せをチェックする工程をとばすことを注記する。代わりに、システムは、異なるモデル化技術およびチューニングパラメーターを調べ2508、テスト測定基準を最大化するモデル化方法およびチューニングパラメーターのセットを選択する。次いで、より細かく間隔を設けられたグリッドを用い、システムは最良の回帰モデルの洗練されたチューニングを行い、チューニングパラメーター値の各セットについて、テストデータとの相関を決定する。テスト測定基準の最良な値を生じるチューニングパラメーターのセットを選択する。次いで、システムは、予測誤差の標準偏差のようなテスト測定基準が、予測が有効と考えられるように、選択された閾値未満であるか否かを決定する2518。例えば、1つの実施形態において、 $R > 0.5$ の相関係数は予測が有効とみなされるのに望ましい。もし得られたテスト測定基準が閾値を満足しないならば、予測を行うことができない2517。もしテスト測定基準が必要な閾値を満足するならば、予測で用いたIVおよびモデルがテストデータで達成した相関係数の組合せと共に、表現型予測を生じさせることができる。

【0452】

失われたデータでの癌コホートにおける交差確証によるモデル選択の説明

この態様を示すためには、National Institute of Health's Pharmacogenomic Research Networkの一部であって、どのようにして個々の遺伝子変異が異なる薬物応答に寄与するかを発見する使命を有するPharmGKBで見出すことができる結腸癌に関連する遺伝子および表現型データの利用に焦点を当てた。このデータベースについては、鍵となる挑戦は失われた情報であった。理想的には、前記した回帰技術を適用して、特定の患者に利用できるすべてのIVからのモデルについてのIVサブセットを自動的に選択したいであろう。しかしながら、これは、モデルを訓練し、テストするために他の患者から入手できるデータの量を制限する。その結果、あまり十分でないIVを含有するデータベースについては、独立変数の全ての可能なサブセットを通じてサーチすることが可能である。各々について、前記したように、必要な結果が測定され、および独立変数の関連セットが利用できる患者のセットを抽出することができる。前記したように、可能な方法の空間をサーチして、陽性数値の独立変数のlogを取ることのような含まれた独立変数を前処理することもできる。含まれた独立変数の各組合せ、および独立変数前処理技術については、テストデータでの交差 - 確証によってモデルを訓練し、テストする。テストデータでの最良の交差 - 確証を有す

るモデルを選択する。一旦、IVについての所与のセットのためにモデルを創製したならば、網羅的なモデルサーチを必要とすることなくIVの同一セットが供給された新しい患者データにそのモデルを適用する。

【0453】

この技術は、結直腸癌薬物イリノテカンについての臨床的副作用を予測するのに用いられてきた。ひどい毒性がイリノテカンを受ける癌患者で共通して観察される。イリノテカン薬物動態学および副作用等、イリノテカン代謝酵素および推定関連性のトランスポーターをコードする遺伝子の対立遺伝子変種との間の関係を記載するデータが含まれた。患者を、MDR1 P-糖蛋白質(ABC B1)、多薬物耐性-関連蛋白質MRP-1(ABC C1)およびMRP-2(ABC C2)、乳癌耐性蛋白質(ABC G2)、チトクロームP450イソ酵素(CYP3A4, CYP3A5)、カルボキシエステラーゼ(CES1, CES2)、UDPグルクロノシル-トランスフェラーゼ(UGT1A1, UGT1A9)、および肝臓転写因子TCF1をコードする遺伝子における変異について遺伝子タイプ分けした。この研究のための遺伝子配列データに関連する表現型データを表15に記載する。

10

【0454】

図26は、ファルマコゲノミック移動エンジンを用いて供給された利用可能なPharmGKBデータが所与の、イリノテカンでの結腸癌治療のための予測結果のモデルを説明する。図26において、モデルは、関連遺伝子座(2601)、用いるインジケーター、この場合、0ないし24時間からのCPT-11の濃度曲線下面積(AUC)のlog(2602)、および12日ないし14日に絶対好中球カウントのNadirのlog(2604)を予測するための0ないし24時間のSN-38 AUCのlog(2603)を示す。テストデータでモデルを交差-確認し、R=64%の相関係数が達成された(2605)。モデル予測の経験的標準偏差はモデルを訓練するのに用いられた(2607)結果のヒストグラムに重ねて示す(2606)。これらの統計学を用いて、イリノテカン治療を完全に差し控えるような通知した治療決定を行い、あるいは顆粒球コロニー刺激因子のような第2の薬物を投与して低いANCおよび得られた感染を妨げることができる。

20

【0455】

高められた診断報告

病気治療の関係では創製された遺伝子データは、データを用いて治療用療法を選択するのに助けることができる臨床家にとって最も用いられるものである。1つの態様において、表現型予測を状況に当てはめ、臨床家または患者に対する報告に組織化する。もう1つの態様において、本明細書中に開示されたシステムおよび方法は、診断lab2703がlabテスト2701および医療報告2702からのデータを検証し、それをデータセンター2704に送り、そこで、それは開示された方法を用いて分析された標準的腫瘍学に一体化されるより大きなシステム(図27参照)の一部として用いることができ、高められた診断報告2705が創製され、医師2706に送られる。

30

【0456】

報告を創製することができる1つの可能な状況は、イリノテカンで治療される結腸癌患者についての予測臨床結果に関するであろう。それは、治療のための禁忌の概念、投与スケジュール、副作用プロフィールをコードすることができる。そのような副作用の例は、2つとも普通である骨髄抑制および後期-開始下痢、緊急の医療的看護を必要とするイリノテカン治療の用量-律速副作用を含む。加えて、ひどい好中球減少症およびひどい下痢は、各々、患者の28%および21%に影響する。あるUGT1A1対立遺伝子、肝臓帰納テスト、ギルバート症候群の過去の医療的履歴、および抗-痙攣薬およびいくつかの抗-催吐薬のようなチトクロームp450を誘導する患者投薬の同定は、イリノテカン用量調整を警告するインジケーターである。

40

【0457】

図28は、表現型予測を用いるイリノテカンでの結直腸癌治療についての高められた報告のモック-アップである。治療に先立ち、報告は患者の癌の段階、過去の医療履歴、現

50

在の投薬および薬物用量を推奨するためのUGT1A1遺伝子型を考慮する。最初の薬物投与からほぼ1日後に報告は、UGT1A1遺伝子中の突然変異、および患者の血液から測定された代謝産物（例えば、SN-38、CPT-11）に基づいた、ほぼ2週間の時間における患者の絶対好中球カウントの予測されたNadirの予測を含む。この予測に基づき、医師は、患者にコロニー刺激因子薬物を与え、またはイリノテカン用量を変更するか否かを決定することができる。また、患者を血液カウント、下痢のグレードについてモニターする。データ源および推奨の正当性を供する。

【0458】

態様の組合せ

先に述べたように、本開示の利点を仮定すれば、他の態様、特徴および実施形態は本明細書中に開示された方法およびシステムの1以上を実行することができる。以下に、開示された発明の種々の態様を複数の方法で組み合わせることができる状況を説明する例の短いリストを掲げる。このリストは包括的であることを意図せず、本発明の態様、特徴および実施形態の多くの他の組合せが可能であることに注意するのは重要である。

【0459】

1つの例は、各々の値を最適化する方法での種々のゲノタイピング測定技術を利用することができる。例えば、labは、Applied Bioscience Taqmanアッセイのような低シグナルの場合において、高価であるが、高い品質のデータを与えることができる技術を用いて、標的DNAを測定し、およびAffymetrix's 500K Genechip、またはMIPSのような高価であるが、多量の遺伝物質を必要とする技術を用いて、良好な質のデータを与え、親DNAを測定することができる。

【0460】

もうひとつの例は、IVF治療を受けているカップルが婦人から収穫された卵を有し、男性からの精子で受精させ、8つの生きた胚を生じる状況であろう。胚盤胞を各胚から収穫し、胚盤胞からのゲノムデータを、Taqmanゲノタイピングアッセイを用いて測定する。他方、Molecular Inversion Probes（分子逆転プローブ）を用い、双方の親から採った組織からジブロイドデータを測定する。男性の精子の1つからの、および婦人の卵の1つからのハプロイドデータもMIPを用いて測定する。親の遺伝子データを用いて8つの胚盤胞のSNPデータを清澄化する。次いで、清澄化された遺伝子データを用いて胚の潜在的表現型に関して予測を行う。最も有望なプロフィールを有する2つの胚を選択し、婦人の子宮に着床させる。

【0461】

もう1つの例は、その夫がテイ・サックス病の家族履歴を有する妊娠した婦人が、彼女が担う胎児が遺伝的に罹患性であるかを知りたがっているが、羊水穿刺は流産のかなりの危険性があるのでそれを受けることを望まない状況であろう。彼女は血液を吸い取り、幾らかの胎児DNAを彼女の血液から単離し、MIPを用いてそのDNAを分析する。彼女および彼女の夫は、従前に分析された彼らの十分なゲノムデータを既に有しており、それはイン・シリコで利用可能である。医師は、親ゲノムのイン・シリコ知識および本明細書中に開示した方法を用いて胎児DNAデータを清澄化し、テイ・サックス病の原因である臨界的遺伝子が胎児のゲノムに存在するかをチェックすることもできる。

【0462】

もう1つの例は、44歳の妊娠した婦人が、彼女が担う胎児がダウン症候群を有し得るかに関心がある状況であろう。彼女は、流産の個人的履歴を仮定すれば、出生前診断で用いる煩わしい技術を有することを警戒しており、従って、彼女は自分の血液を分析することを選択する。健康ケア実践者は、母体血液試料中の胎児細胞を見出すことができ、婦人自身も遺伝子データの知識とともに本明細書中に開示した方法を用い、異数性について診断することができる。

【0463】

もう1つの例は、カップルがIVF治療を受けており；彼らは婦人から収穫した卵を有し、男性からの精子で受精させ、9つの生きた胚を生じる状況であろう。胚盤胞が各胚か

10

20

30

40

50

ら収穫され、胚盤胞からのゲノムデータをIlluminaビーズアッセイを用いて測定する。他方、分子逆転プローブを用いて双方の親から採取された組織からジプロイドデータを測定する。同一方法を用い、父親の精子からのハプロイドデータを測定する。母親から入手できる過剰な卵はなく、従って、バルクジプロイド組織試料は彼女自身の父親および母親から採取され、精子試料は彼女の父親から採取される。それらはすべてMIPを用い分析され、本明細書中で開示された方法を用いて、母親のゲノムについての遺伝子分析を供する。次いで、父親のジプロイドおよびハプロイドデータとともにそのデータを用いて胚盤胞の各々の遺伝子データの高度に正確な分析を行う。表現型予測に基づき、カップルは3つの胚を着床させることを選択する。

【0464】

もう1つの例は、競走馬飼育者が、彼の優勝競走馬によって種付けされた子馬がそれ自体が優勝馬となる尤度を増加させることを望む状況である。彼は所望の雌馬がIVFによって妊娠されるように手配し、雄馬および雌馬からの遺伝子データを用いて、生きた胚から測定された遺伝子データを清澄化する。清澄化された胚遺伝子データは、育種者が関連遺伝子型 - 表現型相関を見出し、望ましい競走馬を最も生産するような着床用の胚を選択することを可能とする。

【0465】

もう1つの例は、妊娠した婦人が彼女が担う胎児がいずれかの深刻な病気に対する素因があるか否かを知りたい状況であろう。父親は既に亡くなっており、従って父親の兄弟および父親の父親から創製されたハプロイドおよびジプロイドデータを用いて、胎児血液サンプリングの間に集められた胎児細胞から測定された胎児の遺伝子データを清澄化することを助ける。健康ケア実践者によって契約された会社は清澄化された遺伝子データを用いて各予測の信頼性ととも、胎児が呈するような表現型のリストを提供する。

【0466】

もう1つの例は、乏しい研究室技術のため、汚染された胎児遺伝子データと場合によっては闘わなければならない羊水穿刺であろう。開示された方法を用いて、母性および父性遺伝子データを用いて汚染された胎児遺伝子データを清澄化することができる。開示された方法が汚染DNAの増大した速度を補うことができることを知って、稔性手法を緩和させることによって研究室がコストを切り詰めることができる状況を想像することができる。

【0467】

もう1つの例は、40代の婦人が妊娠を得ようとしてIVFを受けている状況であろう。彼女は、胚をスクリーニングして、遺伝病を最も有しないようであり、最も着床し、妊娠まで持っていきそうなものを選択することを望む。彼女が用いているIVFクリニックは生きた胚の各々から胚盤胞を収穫し、標準的な手法を用いてDNAを増幅し、鍵となるSNPを測定する。次いで、技術者は本明細書中に開示された方法を用いて、染色体アンバランスについてスクリーニングし、また、胚の遺伝子データを見出し、それを清澄化して、各胚の表現型素因について予測を行う。

【0468】

もう1つの例は、妊娠した婦人が羊水穿刺を有し、本明細書中に開示された方法とともに、血液試料中の胎児細胞における遺伝物質を用いて異数性および他の染色体異常についてスクリーニングする。

【0469】

1つの例は、径基礎核関数およびノルム喪失関数とともにサポートベクトルマシンを用いる非線形モデルがヒト成人の遺伝子型および表現型データを利用して、早期開始アルツハイマー病の尤度を予測し、該病気の開始を遅らせることができる可能なライフスタイルの変化および運動養生法を提案する。

【0470】

もう1つの例は、LASSO技術を用いる線形モデルが、癌の遺伝子データとともに、肺癌に罹った成人婦人の遺伝子型および表現型データを利用して、いずれの医薬が該病気

10

20

30

40

50

の進行を遅らせるのに最も効果的であるかを予測する婦人の医師についての医師用の報告を作成する。

【0471】

もう1つの例は、複数のモデルを、クローン病患者の遺伝子、表現型および臨床データより成る集合データについてテストし、次いで、最も正確であることが判明する非線形回帰モデルが成人男性の表現型および臨床データを利用して、彼のクローン病の徴候を緩和するようであるある種の栄養サプリメントを提案する報告を作成する状況であろう。

【0472】

もう1つの例は、Hapmap Projectを通じて獲得されたデータから形成された偶発事象表を利用し、かつ胚からの胚盤胞から集めた遺伝子情報を利用するモデルを用いて、もし胚が着床すれば、結果をもたらす子供のありそうな表現型に関して予測を行う状況であろう。

【0473】

もう1つの例は、新生児に感染するHIVの株の遺伝子情報を利用する線形回帰モデルを用いて、いずれの抗ウイルス薬物が、もし投与されたならば、成人に達する最大のチャンスに彼女に与えるかを示唆する赤ん坊の医師用の報告を作成する状況であろう。

【0474】

もう1つの例は、新しい研究が公表され、中年婦人における心筋梗塞の罹患率、およびある遺伝子および表現型マーカーの間のある相関を示唆する状況であろう。次いで、これは非線形回帰モデルの使用を促進して、中年データの集合データならびにそのデータがシステムに知られている個体の遺伝子および表現型データを再度調べ、次いで、該モデルは、心筋梗塞の危険性が最もある婦人を同定し、予測される危険性を彼らに通知する各医師に送られる報告書を作成する。

【0475】

もう1つの例は、複数のモデルを試みられた種々の薬物介入を含めた、結腸癌に罹った人々の集合データについてテストされる状況であろう。最良の予測を可能とすることが判明するモデルを用いて、実験的新しい医薬から最も利益を受けるであろう患者、および該新しい医薬に対する権利を所有する会社によってそれらの結果が用いられ、臨床試験を行うにおいて彼らを助ける。

【0476】

定義

SNP（単一ヌクレオチド多形）：個体間変異を示す傾向がある染色体上の特別な遺伝子座。

SNPを要求すること：直接および間接的証拠を考慮し、特定の塩基対の同一性を質問すること。

対立遺伝子を要求すること：SNPを要求すること遺伝子データを清浄化すること：関連する個人の遺伝子データ、および本明細書中に記載された方法を用いて不完全な遺伝子データを取り、誤差のいくつかまたは全てを修正すること。

不完全な遺伝子データ：以下の：対立遺伝子ドロップアウト、不明瞭な塩基対測定、正しくない塩基対測定、偽シグナル、または失われた測定のいずれかを持つ遺伝子データ

信頼性：要求されたSNP、対立遺伝子、または対立遺伝子のセットが個体の現実の遺伝子状態を表す統計学的尤度

多重遺伝子：複数の遺伝子または対立遺伝子によって影響される

ノイズな遺伝子データ：不完全な遺伝子データとも呼ばれる不完全遺伝子データ；

未清浄化遺伝子データ：測定された遺伝子データ、すなわち、生の遺伝子データにおいてノイズの存在について修正するのにいずれの方法も用いられたことがない；また、粗遺伝子データとも呼ばれる

直接的関係：母親、父親、息子、または娘

染色体領域：染色体のセグメント、または全染色体

親サポート：遺伝子データを清浄化する開示された方法で時々用いられる名称

10

20

30

40

50

染色体のセクション：1塩基対ないし全染色体のサイズの範囲とすることができる染色体のセクション。

【0477】

(表)

【0478】

【表1】

	常染色体	X-連鎖	Y-連鎖	ミトコンドリア	合計
公知の配列を持つ遺伝子	9918	157	18	37	10153
公知の配列および表現型を持つ遺伝子	382	30	0	0	389
表現型の記載、分子基礎既知	1652	147	2	26	1827
メンデル表現型または遺伝子座、分子基礎未知	132	15	1	0	150
疑われるメンデル基礎を持つ表現型	2092	147	2	0	2241
合計	15381	496	21	63	16419

10

表1

【0479】

【表2】

技術の記載	計算に適用された式 N=	アクセイ名称	E(ml)	E(m)	Sd(ml)	Sd(m)	コピー数の比較	モザイク現象	SNP当たりの価格 N	さらなるコスト	因子シグマ	1-Cor	セグメント	ウェル当たりの細胞	総細胞 WoPCR	コスト	
平均の比較 参照試料なし	$2^2s^2(2s^2/(m-nm))^2$	Taqman	32.26	30.75	0.597	0.597	2vs1	100.0%	16	\$0.20	30	5	2.87E-07	5	50	800	\$46.00
平均の比較 参照試料あり	$2^2s^2(2s^2/(m1-m0))^2$	Taqman	32.26	30.75	0.597	0.597	2vs1	100.0%	32	\$0.20	30	5	2.87E-07	5	50	1600	\$62.00
平均の比較 参照試料なし	$4(2^2)(0.95^s2+0.05^s2)(0.05^2/(m1-m0)^2)$	Taqman	32.26	30.75	0.597	0.597	2vs1	5.0%	1001	\$0.20	30	2	2.30E-02	5	50	5000	\$1,031.00
平均の比較 参照試料あり	$4(2^2)(1.95^s2+0.05^s2)(0.05^2/(m1-m0)^2)$	Taqman	32.26	30.75	0.597	0.597	2vs1	5.0%	2002	\$0.20	30	2	2.30E-02	5	50	100100	\$2,032.00
平均の比較 参照試料なし	$4(2^2)(0.95^s2+0.05^s2)(0.05^2/(m1-m0)^2)$	MPs	32.26	30.75	0.597	0.597	2vs1	5.0%	1001	\$0.01	200	2	2.30E-02	5	50	5000	\$250.00
平均の比較 参照試料あり	$4(2^2)(1.95^s2+0.05^s2)(0.05^2/(m1-m0)^2)$	MPs	32.26	30.75	0.597	0.597	2vs1	5.0%	2002	\$0.01	200	2	2.30E-02	5	50	100100	\$300.10
平均の比較 参照試料なし	$4(1^2)((1-f)^s2+ff^s2)(f^2/(m1-m0)^2); f=0.081$	Taqman	32.26	30.75	0.597	0.597	2vs1	8.1%	384	\$0.20	30	2	2.30E-02	5	50	19200	\$414.00
平均の比較 参照試料なし	$4(1^2)((1-f)^s2+ff^s2)(f^2/(m1-m0)^2); f=0.053$	Taqman	32.26	30.75	0.597	0.597	2vs1	4.1%	384	\$0.20	30	1	1.56E-01	5	50	19200	\$414.00
平均の比較 参照試料なし	$4(2^2)((1-f)^s2+ff^s2)(f^2/(m1-m0)^2); f=0.19$	Taqman	32.26	30.75	0.597	0.597	2vs1	19.0%	70	\$0.20	30	2	2.30E-02	5	50	3500	\$100.00
平均の比較 参照試料なし	$2^2s^2(2s^2/(m1-m0))^2$	qPCR	27.65	26.64	0.53	0.53	2vs1	100.0%	27	\$0.15	30	5	2.87E-07	5	50	1350	\$50.25
平均の比較 参照試料あり	$2^2s^2(2s^2/(m1-m0))^2$	qPCR	27.65	26.64	0.53	0.53	2vs1	100.0%	55	\$0.15	30	5	2.87E-07	5	50	2700	\$71.25
平均の比較 参照試料なし	$4(2^2)(0.95^s2+0.05^s2)(0.05^2/(m1-m0)^2)$	qPCR	27.65	26.64	0.53	0.53	2vs1	5.0%	1754	\$0.15	30	2	2.30E-02	5	50	87700	\$1,345.00
平均の比較 参照試料あり	$4(2^2)(1.95^s2+0.05^s2)(0.05^2/(m1-m0)^2)$	qPCR	27.65	26.64	0.53	0.53	2vs1	5.0%	3508	\$0.15	30	2	2.30E-02	5	50	175400	\$2,661.00
平均の比較 参照試料なし	$4(1^2)((1-f)^s2+ff^s2)(f^2/(m1-m0)^2); f=0.107$	qPCR	27.65	26.64	0.53	0.53	2vs1	10.7%	384	\$0.15	30	2	2.30E-02	5	50	19200	\$318.00
平均の比較 参照試料なし	$4(1^2)((1-f)^s2+ff^s2)(f^2/(m1-m0)^2); f=0.054$	qPCR	27.65	26.64	0.53	0.53	2vs1	5.4%	384	\$0.15	30	1	1.56E-01	5	50	19200	\$318.00
平均の比較 参照試料なし	$4(2^2)((1-f)^s2+ff^s2)(f^2/(m1-m0)^2); f=0.19$	qPCR	27.65	26.64	0.53	0.53	2vs1	19.0%	121	\$0.15	30	2	2.30E-02	5	50	6050	\$120.75

20

表2

【0480】

30

40

【表 3】

snp_id	e1	e2	p1	p2	m1	m2	pe1	pe2	pp1	pp2	pm1	pm2
101100940	C	T	T	C	C	T	0.9538	0.8902	0.8626	0.8580	0.8654	0.9101
101164838	T	C	T	C	T	C	0.9359	0.9521	0.9406	0.9253	0.9957	0.8770
rs1463589	C	C	T	C	C	C	0.9428	0.9928	0.9841	0.9266	0.8661	0.9798
101028396	C	G	C	G	C	C	0.9252	0.8792	0.9246	0.9856	0.9819	0.8631
101204217	A	G	G	A	G	G	0.9799	0.9843	0.9194	0.9478	0.9438	0.9709
101214313	A	G	G	A	G	A	0.8513	0.9863	0.9521	0.9707	0.8570	0.9639
101231593	G	A	G	G	A	A	0.9857	0.9653	0.8908	0.9036	0.9431	0.9832
rs1426442	G	C	C	G	C	G	0.9338	0.9278	0.9469	0.9514	0.8766	0.9017
rs7486852	C	C	C	T	T	C	0.9566	0.9616	0.9390	0.8673	0.8785	0.8889
101266729	A	G	A	G	A	G	0.9238	0.9500	0.9026	0.9855	0.8760	0.9381

10

表3

【 0 4 8 1 】

【表 4】

snp_id	e1	e2	p1	p2	m1	m2	pe1	pe2	pp1	pp2	pm1	pm2
101019515	G	G	G	G	G	G	0.9134	0.8768	0.8666	0.9690	0.8679	0.8599
101100940	C	T	T	C	C	T	0.9538	0.8902	0.8626	0.8580	0.8654	0.9101
101160854	A	A	A	A	A	A	0.8705	0.9769	0.8763	0.8870	0.9311	0.9553
rs4980809	A	G	G	A	A	G	0.9638	0.9951	0.9582	0.9621	0.9197	0.9199
101058479	G	A	G	A	G	A	0.9003	0.9885	0.8906	0.9235	0.9787	0.8792
101236938	G	G	G	G	G	A	0.8528	0.9710	0.8810	0.9249	0.9274	0.9891
rs7137405	T	T	T	T	T	A	0.9360	0.9918	0.9148	0.9558	0.9135	0.9388
101251161	G	G	G	G	G	G	0.9802	0.8620	0.9372	0.8501	0.9891	0.8679
101270051	G	G	G	G	G	A	0.9004	0.9643	0.9778	0.9060	0.9943	0.8962
rs215227	G	G	G	G	G	A	0.9244	0.9236	0.9629	0.8575	0.9019	0.9362
101245075	G	G	G	G	G	G	0.9958	0.8593	0.9129	0.8504	0.8534	0.9866
101158538	A	G	A	G	G	G	0.9471	0.8909	0.8710	0.9581	0.8961	0.9046
rs2535386	A	A	A	A	A	A	0.9273	0.9479	0.9867	0.8918	0.9264	0.9750
rs6489653	T	T	T	T	T	T	0.9453	0.9776	0.9051	0.8547	0.9636	0.9532
101137205	C	G	C	C	G	G	0.8619	0.9503	0.9029	0.9426	0.8845	0.9282
101089311	T	C	C	C	C	T	0.8844	0.9381	0.9719	0.8636	0.9186	0.9652
101205712	A	A	A	A	A	A	0.8513	0.9226	0.8755	0.8999	0.9193	0.8535
101124605	G	G	G	G	G	G	0.8981	0.9093	0.9075	0.8676	0.8931	0.9258
101025989	G	T	T	G	G	T	0.9695	0.9016	0.8722	0.8821	0.9787	0.9273
rs4766370	T	A	A	T	T	A	0.8886	0.9166	0.8762	0.8767	0.9890	0.8536

20

30

表4

【 0 4 8 2 】

40

【表 5】

snp_id	true_value	true_hyp	ee	pp	mm	SnipProb	HypProb
101100940	CT	p2 m2	CT	TC	CT	0.8416	0.5206
101164838	CT	p2 m1	TC	TC	TC	0.9061	0.5206
rs1463589	CC	p2 m1	CC	TC	CC	0.9946	0.5206
101028396	GC	p2 m1	CG	CG	CC	0.9791	0.5206
101204217	AG	p2 m2	AG	GA	GG	0.9577	0.5206
101214313	GA	p1 m2	AG	GA	GA	0.9308	0.5206
101231593	GA	p1 m2	GA	GG	AA	1.0000	0.5206
rs1426442	CG	p1 m2	GC	CG	CG	0.9198	0.5206
rs7486852	CC	p1 m2	CC	CT	TC	0.9138	0.5206
101266729	AG	p1 m2	AG	AG	AG	0.9296	0.5206

10

表5

【 0 4 8 3 】

【表 6】

snp_id	true_value	true_hyp	ee	pp	mm	SnipProb	HypProb
101019515	GG	p1 m1	GG	GG	GG	1.0000	0.9890
101100940	TC	p1 m1	CT	TC	CT	0.9946	0.9890
101160854	AA	p1 m1	AA	AA	AA	1.0000	0.9890
rs4980809	GA	p1 m1	AG	GA	AG	0.9961	0.9890
101058479	GG	p1 m1	GA	GA	GA	0.9957	0.9890
101236938	GG	p1 m1	GG	GG	GA	1.0000	0.9890
rs7137405	TT	p1 m1	TT	TT	TA	1.0000	0.9890
101251161	GG	p1 m1	GG	GG	GG	1.0000	0.9890
101270051	GG	p1 m1	GG	GG	GA	1.0000	0.9890
rs215227	GG	p1 m1	GG	GG	GA	1.0000	0.9890
101245075	GG	p1 m1	GG	GG	GG	1.0000	0.9890
101158538	AG	p1 m1	AG	AG	GG	0.9977	0.9890
rs2535386	AA	p1 m1	AA	AA	AA	1.0000	0.9890
rs6489653	TT	p1 m1	TT	TT	TT	1.0000	0.9890
101137205	CG	p1 m1	CG	CC	GG	1.0000	0.9890
101089311	CC	p1 m1	TC	CC	CT	0.9940	0.9890
101205712	AA	p1 m1	AA	AA	AA	1.0000	0.9890
101124605	GG	p1 m1	GG	GG	GG	1.0000	0.9890
101025989	TG	p1 m1	GT	TG	GT	0.9973	0.9890
rs4766370	AT	p1 m1	TA	AT	TA	0.9973	0.9890

20

30

表6

【 0 4 8 4 】

40

【表7】

Pop.Freq	ph	pd	DHAlgorithm1		DHAlgorithm2	
			P1精度	P2精度	P1精度	P2精度
データ	0.95	0.95	0.982	0.951	0.95	0.906
データ	0.75	0.75	0.891	0.811	0.749	0.618
データ	0.25	0.25	0.71	0.71	0.253	0.25
データ	0.5	0.9	0.849	0.838	0.499	0.768
データ	0.9	0.5	0.942	0.734	0.898	0.347
データ	0.6	0.8	0.852	0.816	0.601	0.673
均一	0.95	0.95	0.95	0.906	0.949	0.905
均一	0.75	0.75	0.749	0.612	0.749	0.612
均一	0.25	0.25	0.25	0.248	0.25	0.25
均一	0.5	0.9	0.69	0.669	0.501	0.671
均一	0.9	0.5	0.901	0.412	0.901	0.413
均一	0.6	0.8	0.678	0.618	0.6	0.618

10

表7

【0485】

【表8】

Pop.Freq	ph	pd	pe	PSAlgorithm1		PSAlgorithm2	
				P1精度	P2精度	P1精度	P2精度
データ	0.95	0.95	0.95	0.834	0.815	0.928	0.931
データ	0.75	0.75	0.75	0.797	0.769	0.819	0.819
データ	0.25	0.25	0.25	0.711	0.682	0.703	0.687
データ	0.5	0.9	0.9	0.849	0.838	0.866	0.864
データ	0.9	0.5	0.5	0.792	0.809	0.756	0.752
データ	0.6	0.8	0.8	0.777	0.788	0.835	0.828
均一	0.95	0.95	0.95	0.673	0.631	0.898	0.901
均一	0.75	0.75	0.75	0.549	0.497	0.635	0.65
均一	0.25	0.25	0.25	0.239	0.249	0.252	0.25
均一	0.5	0.9	0.9	0.601	0.611	0.814	0.818
均一	0.9	0.5	0.5	0.459	0.391	0.449	0.468
均一	0.6	0.8	0.8	0.544	0.511	0.672	0.679

20

30

表8

【0486】

【表9】

Farrer(2005)			Alvarez(1999)		Labert(1998)	
	D+	D-	D+	D-	D+	D-
A+B+	58/151	25/206	161/350	176/400	189/350	224/400
A+B-	28/151	69/206				
A-B+	43/151	20/206	B+	334/573	104/509	
A-B-	28/151	20/206	B-	239/573	405/509	

40

表9

【0487】

【表10】

prob.									
P(A+B+D+)	0.28	0.024	0.24	0.32	0.37	0.039	0.29	0.45	
P(A+B+D-)	0.18	0.016	0.16	0.21	0.17	0.026	0.12	0.22	
P(A+B+D+)	0.22	0.027	0.18	0.27	0.28	0.037	0.21	0.35	
P(A+B+D-)	0.18	0.018	0.15	0.22	0.12	0.023	0.08	0.16	
P(A-B+D+)	0.31	0.028	0.28	0.36	0.39	0.038	0.28	0.4	
P(A-B+D-)	0.05	0.017	0.02	0.09	0.1	0.021	0.08	0.14	
P(A-B+D+)	0.2	0.032	0.13	0.28	0.02	0.011	0	0.04	
P(A-B+D-)	0.58	0.034	0.52	0.63	0.61	0.034	0.54	0.68	

表10

【0488】

【表11】

方法	平均	APV	ATV	NEV	RTV	SOV	LPV	IDV	平均
RR	平均	79.30	89.44	78.37	87.96	82.93	75.35	82.81	79.45
	標準偏差	0.26	0.90	1.18	0.49	0.60	1.36	0.65	0.86
NN	平均	83.16	68.45	87.59	90.88	89.10	73.67	89.07	83.13
	標準偏差	0.52	1.22	0.19	0.36	0.15	0.82	0.10	0.61
PCA	平均	86.49	73.52	87.65	91.98	88.72	84.79	87.86	86.86
	標準偏差	0.19	0.79	0.14	0.12	0.12	0.26	0.08	0.33
DT	平均	77.32	65.12	83.61	85.57	83.63	89.05	82.54	78.11
	標準偏差	0.33	0.93	0.15	0.21	0.31	0.87	0.34	0.54
SS	平均	85.25	73.27	89.18	91.65	89.81	77.15	87.79	84.87
	標準偏差	0.03	1.02	0.14	0.10	0.08	0.39	0.08	0.42
SVM-L	平均	87.48	75.24	88.95	92.24	89.53	87.13	88.12	86.96
	標準偏差	0.15	1.05	0.08	0.09	0.10	0.23	0.10	0.42
LASSO	平均	89.17	76.60	89.88	93.47	91.00	86.78	88.61	87.93
	標準偏差	0.08	0.62	0.07	0.16	0.07	0.17	0.07	0.26
SVM-R	平均	88.22	82.27	89.92	92.29	89.48	86.19	90.28	88.38
	標準偏差	0.03	0.46	0.06	0.07	0.06	0.32	0.03	0.22
試料		577	142	617	510	598	1253	579	

表11

【0489】

【表12】

方法	平均	STC	ABC	AZT	D4T	DDC	DDI	DLV	EFV	NYP	TDF	平均
RR	平均	87.28	74.56	78.3	79.76	70.66	69.65	80.66	78.12	73.47	71.85	76.48
	標準偏差	0.72	1.15	0.94	0.61	1.06	0.66	0.17	0.29	0.63	0.66	0.69
NN	平均	94.33	74.96	71.72	80.23	73.42	65.79	84.15	80.5	84.05	68.71	77.73
	標準偏差	0.26	0.50	0.45	0.82	0.69	1.09	0.43	0.60	0.60	1.23	1.18
PCA	平均	92.23	83.3	84.14	88.83	85.28	85.19	80.88	77.43	77.34	68.36	82.30
	標準偏差	0.11	0.28	0.16	0.17	0.12	0.23	0.21	0.25	0.32	0.88	0.34
DT	平均	94.68	81.68	81.32	80.87	73.66	78.35	81.92	78.78	85.93	88.91	79.48
	標準偏差	0.22	0.31	0.58	0.52	0.78	0.61	0.30	0.92	0.36	2.05	0.84
SS	平均	94.91	84.48	81.86	86.88	83.99	82.17	83.87	80.15	78.8	67.27	82.52
	標準偏差	0.07	0.22	0.42	0.32	0.17	0.67	0.21	0.28	0.68	1.39	0.59
SVM-L	平均	84.31	85.27	84.04	89.09	85.5	83.53	87.07	83.78	78.77	87.9	83.94
	標準偏差	0.12	0.21	0.21	0.15	0.17	0.32	0.16	0.24	0.46	1.87	0.68
LASSO	平均	95.48	86.32	85.54	91.37	87.54	85.74	85.16	82.31	82.19	70.68	86.22
	標準偏差	0.07	0.18	0.11	0.10	0.14	0.24	0.12	0.20	0.20	0.41	0.20
SVM-R	平均	95.94	86.84	89.99	89.30	85.29	85.40	85.33	84.61	85.39	78.76	86.62
	標準偏差	0.11	0.19	0.13	0.12	0.26	0.22	0.20	0.26	0.27	0.41	0.25
試料		356	354	353	399	318	354	366	384	401	86	

表12

【0490】

10

20

30

40

【表13】

薬物	APV	ATV	MFV	RFV	SCV	LPV	IDV
#試料	577	142	617	510	598	253	579
#突然変異	320	320	320	320	320	320	320
#プレディクター	120	50	90	120	120	100	110

表13

【0491】

【表14】

10

薬物	3TC	ABC	AZT	D4T	DDC	DDI	DLV	EFV	NVP	TDF
#試料	356	354	353	356	319	354	335	334	401	96
#突然変異	791	791	791	791	791	791	791	791	791	791
#プレディクター	80	70	100	110	80	100	170	120	120	60

表14

【0492】

【表15】

20

対象ID	Pharm_GKBにおける対象についてのPharm_GKBアクセションID
性別	男性、女性、または未知
人種民族性	カッコに入れた自己報告情報を伴うNIH 管理および予算オフィスによって規定
用量 (mg/m ²)	イリノテカン [®] の調整された容量
BSA (m ²)	平方メートルで表した身体表面
用量 (mg)	ミリグラムで表したイリノテカンの用量
CPT-11 AUC ₂₄	蛍光検出法と共に逆相HPLC法において測定した、1サイクルの1日目の、 0ないし24時間の濃度曲線下のイリノテカン、SN-38およびSN-38Gの面積 (hr*ng/ml)
SN38 AUC ₂₄	
SN38 G AUC ₂₄	
APC AU ₀₋₂₄	0ないし24時間の濃度曲線下の7-エチル-10-[4-N-(5-アミノペンタン酸)1-ヒペリジノ]-カルボニルオキシカンプトテンシン (APC) 面積
ANC Nadir	最初のイリノテカンに用量 (12ないし14日) 後の絶対抗中球カウント (ANC) のNadir値
下痢グレード	処理のサイクルの間に観察された通常 [®] の毒性基準バージョン2. 0下痢を用いて評価した

表15

30

【図1】

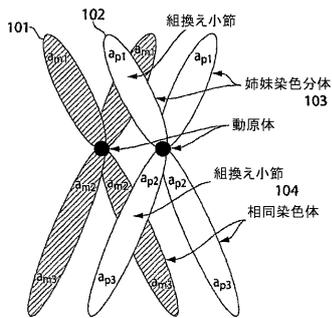


Fig.1

【図2】

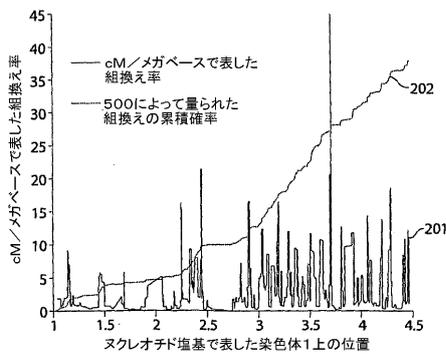


Fig.2

【図6】

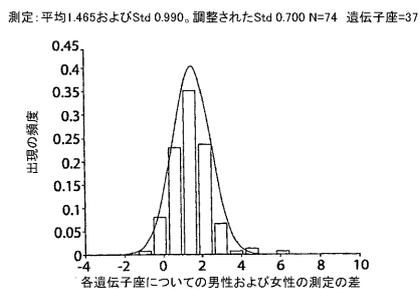


Fig.6

【図7】

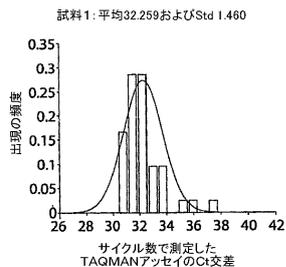


Fig.7

【図3】

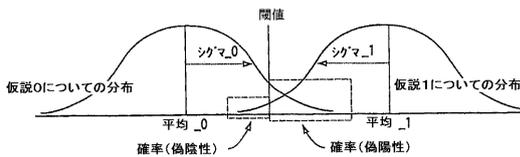


Fig.3

【図4】

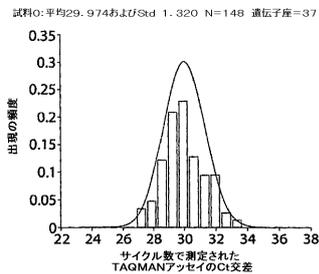


Fig.4

【図5】

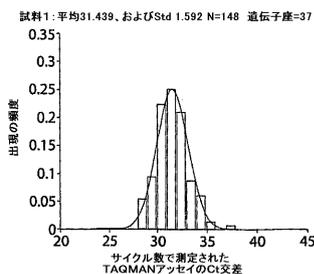


Fig.5

【図8】

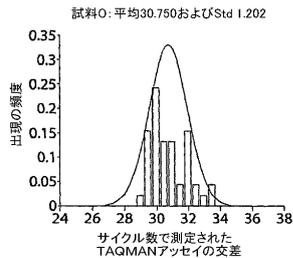


Fig.8

【図9】

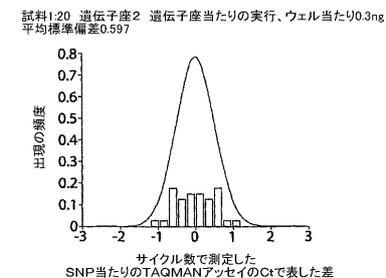


Fig.9

【 図 1 0 】

試料0: 平均26.640およびStd 1.148

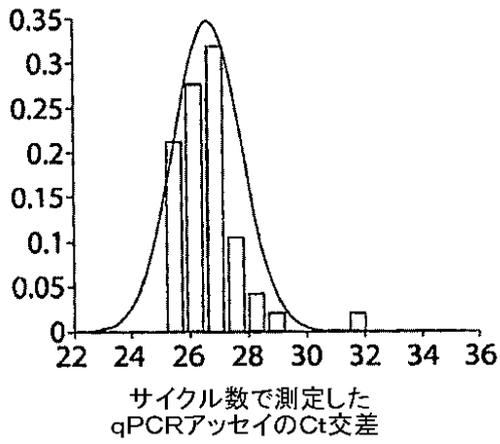


Fig. 10

【 図 1 1 】

試料1: 平均27.652およびStd 1.401

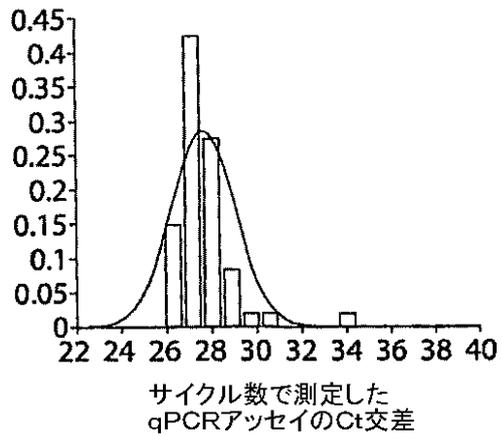


Fig. 11

【 図 1 2 】

測定: 平均1.012およびStd 0.750 調整されたStd 0.530

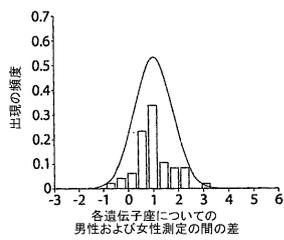


Fig. 12

【 図 1 5 】

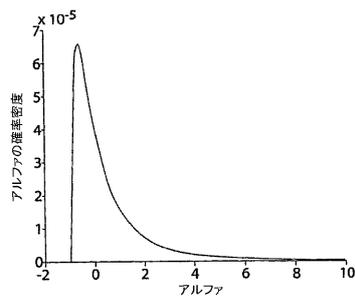


Fig. 15

【 図 1 3 】

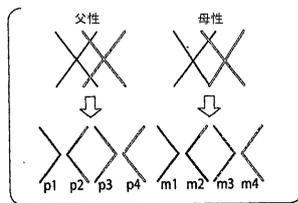


Fig. 13

【 図 1 4 】

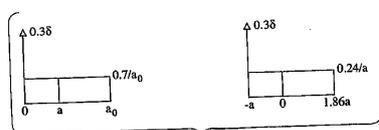


Fig. 14

【図16】

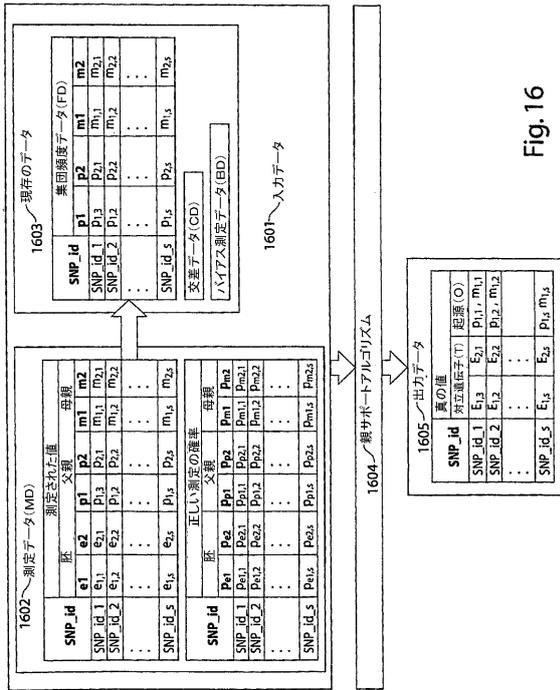


Fig. 16

【図17】

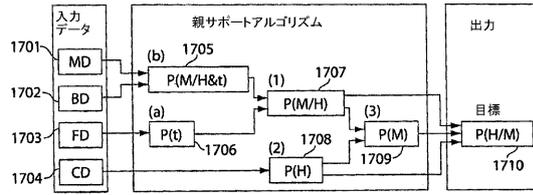


Fig. 17

【図18】

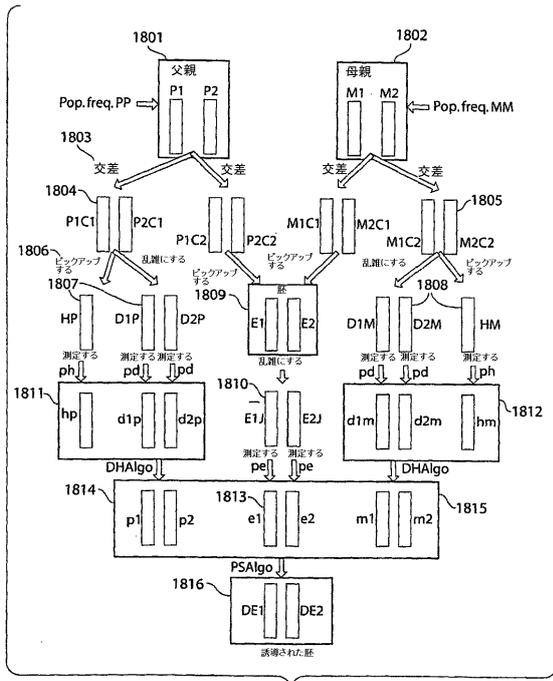


Fig. 18

【図19】

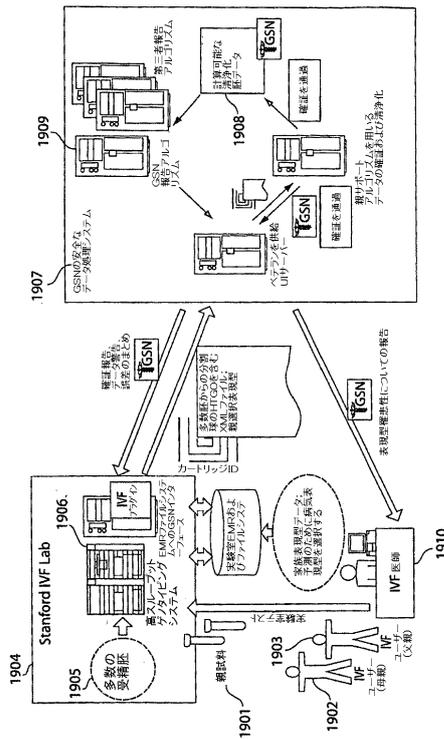


Fig. 19

【 図 2 0 】

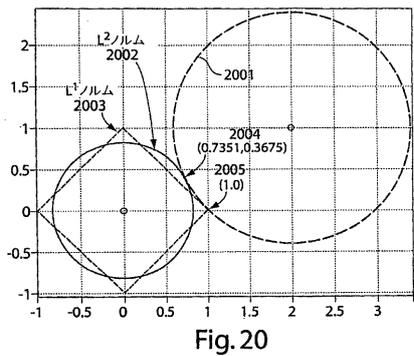


Fig. 20

【 図 2 1 】

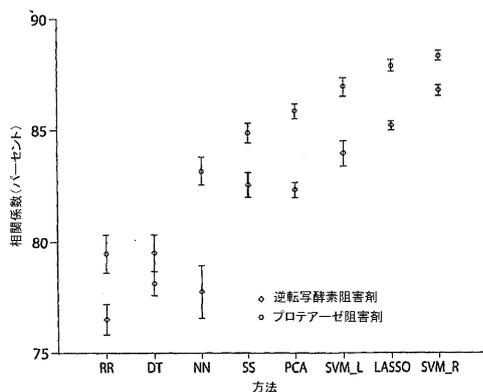


Fig. 21

【 図 2 3 】

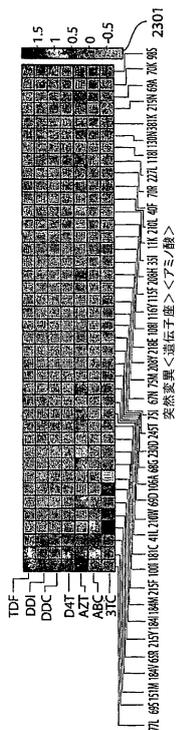


Fig. 23

【 図 2 2 】

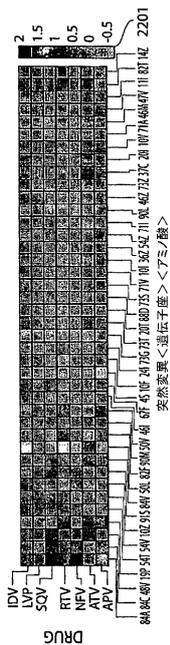


Fig. 22

【 図 2 4 】

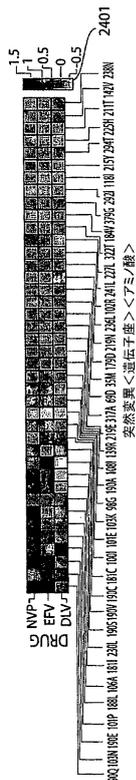


Fig. 24

【 図 2 5 】

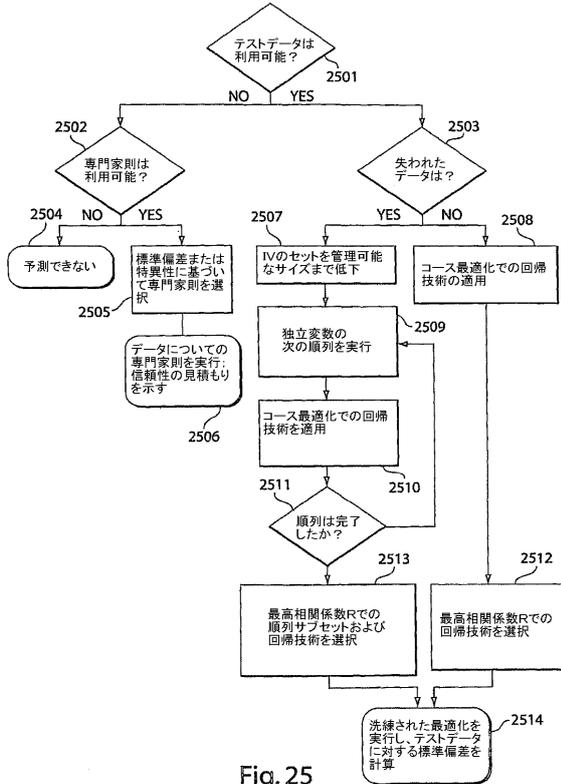


Fig. 25

【 図 2 6 】

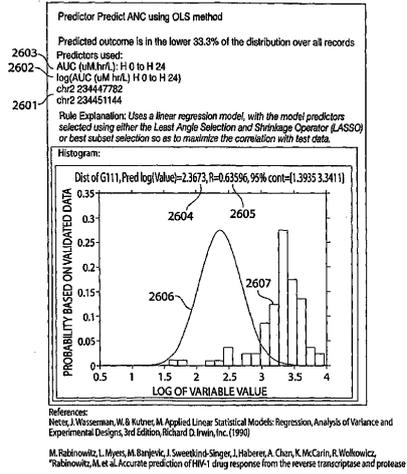


Fig. 26

【 図 2 7 】

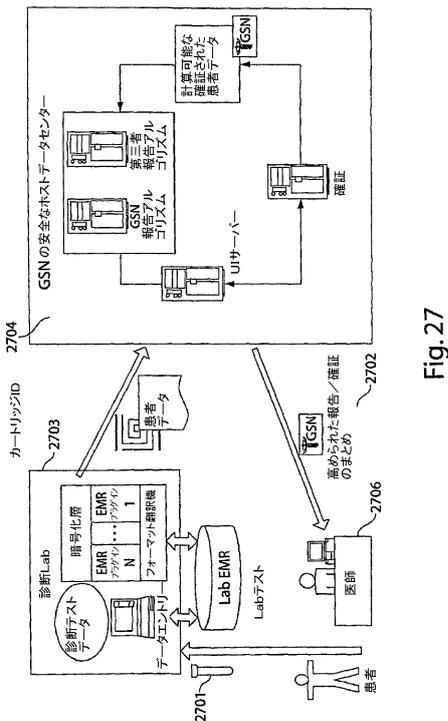


Fig. 27

【 図 2 8 】

Enhanced Colorectal Cancer Treatment Report		GSN CENTRORESEARCH	
NAME: R.E. Wiseman	Report Date: March 26, 2006		
MNH: 112234455	Clinical Indication: Malignant neoplasm, colon		
DOB: 08/08/1958	ICD-9 153.9		
Sex: Male	Staging of Colorectal Carcinoma	Dukes C; TNM Stage IIIA	
Therapy: Irinotecan		Q3 Week Dosing Regimen	
Date of therapy: March 25, 2006			
Other pertinent past medical history:			
Abnormal bilirubin glucuronidation (Gilbert's)			
Significant drug-drug interactions with patient current medications		Effect on Irinotecan levels	
Phenytoin			
Paroxetine			
Ketocoazole		CONTRAINDICATED	
Data Sources:			
Quest Labs	3/20/2006	UGT1A1 Assay - 3/26/2006	
EMR Data	3/25/2006	Heterozygous for UGT1A1 *28	
White blood cell count	3.0 K/uL	AUC _{0-12h} (ng · h/mL) 5582.9	
Absolute neutrophil count	3000/mm ³	AUC _{0-24h} (ng · h/mL) 1261.3	
Neutropenia grade	None	ANC Prediction:	
Hemoglobin	12 g/dL		
Platelet count	150 K/uL	Grade 2 Neutropenia	
Hypertibrinemia	1.2 mg/dL		
Diarrhea grade	grade 2 (4-6 stools/day)		
SUMMARY RECOMMENDATIONS		RATIONALE	
Pre-Treatment		Gilbert's Syndrome	
Consider reducing starting dose to 350 mg/m ²		Discontinue ketocoazole 1 week prior to therapy, resume 2 days after therapy if indicated	
Consider alternative regimens not including Irinotecan		Ketocoazole c-> Irinotecan interaction	
Day 1 Post Treatment Predictions		*28/*28 genotype	
Decrease dose by 50 mg/m ²		Grade2 neutropenia (1000-1499/mm ³)	
Consider addition of colony stimulating factor		Grade2 neutropenia (1000-1499/mm ³)	
Omit dose until resolved to baseline, then reduce to 250 mg/m ²		Grade2 diarrhea (4-6 stools/day)	

Fig. 28

フロントページの続き

- (31)優先権主張番号 60/754,396
(32)優先日 平成17年12月29日(2005.12.29)
(33)優先権主張国 米国(US)
- (31)優先権主張番号 60/774,976
(32)優先日 平成18年2月21日(2006.2.21)
(33)優先権主張国 米国(US)
- (31)優先権主張番号 60/789,506
(32)優先日 平成18年4月4日(2006.4.4)
(33)優先権主張国 米国(US)
- (31)優先権主張番号 60/817,741
(32)優先日 平成18年6月30日(2006.6.30)
(33)優先権主張国 米国(US)
- (31)優先権主張番号 11/496,982
(32)優先日 平成18年7月31日(2006.7.31)
(33)優先権主張国 米国(US)
- (31)優先権主張番号 60/846,610
(32)優先日 平成18年9月22日(2006.9.22)
(33)優先権主張国 米国(US)
- (74)代理人 100117019
弁理士 渡辺 陽一
- (74)代理人 100150810
弁理士 武居 良太郎
- (74)代理人 100141977
弁理士 中島 勝
- (72)発明者 マシュー ラビノビッツ
アメリカ合衆国 カリフォルニア 94028, ボルトラ バレー, ハイフィールドズ ロード
80
- (72)発明者 ミレナ バンジェビック
アメリカ合衆国 ニューヨーク 10009, ニューヨーク, イースト 14ティーエイチ
ストリート 445, アpartment 2エフ
- (72)発明者 ザカリー ポール デムコ
アメリカ合衆国 マサチューセッツ 02144, サマービル, セント ジェームス アベニ
ュー 31ビー
- (72)発明者 デイビッド スコット ジョンソン
アメリカ合衆国 カリフォルニア 94028, ボルトラ バレー, ゴールデン オーク ド
ライブ 425

審査官 松原 寛子

- (56)参考文献 特表2004-502466(JP,A)
特表2004-533243(JP,A)
厚生省精神・神経疾患研究委託費 筋ジストロフィーの臨床・疫学及び遺伝相談に関する研究
平成6・7年度
日本不妊学会雑誌, 日本, 2001年, 第46巻、第1号, p.43-46
遺伝学的検査に関するガイドライン, 日本, 遺伝医学関連学会(日本遺伝カウンセリング学会、
日本, 2003年
Nature Reviews Genetics, 2004年, Vol.5, p.251-261

Ecology Letters, 2004年, Vol.7, p.509-520

(58)調査した分野(Int.Cl., DB名)

C12N 15/11

C12Q 1/6837

C12Q 1/6874

G06F 19/22

JSTPlus/JMEDPlus/JST7580(JDreamIII)

CPlus/MEDLINE/EMBASE/BIOSIS(STN)