

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4359805号  
(P4359805)

(45) 発行日 平成21年11月11日(2009.11.11)

(24) 登録日 平成21年8月21日(2009.8.21)

(51) Int. Cl.	F I		
G 0 6 F 19/00	(2006.01)	G 0 6 F 19/00	6 0 0
G 0 6 F 17/18	(2006.01)	G 0 6 F 17/18	D
C 1 2 Q 1/68	(2006.01)	C 1 2 Q 1/68	A

請求項の数 12 (全 23 頁)

(21) 出願番号	特願2000-620133 (P2000-620133)	(73) 特許権者	309024066
(86) (22) 出願日	平成12年5月25日 (2000.5.25)		アベンティス・ホールディングス・インコーポレイテッド
(65) 公表番号	特表2003-500715 (P2003-500715A)		アメリカ合衆国デラウェア州19807.
(43) 公表日	平成15年1月7日 (2003.1.7)		グリーンビル. スウィート200. ケネット・パイク3711
(86) 国際出願番号	PCT/US2000/014674	(74) 代理人	100091731
(87) 国際公開番号	W02000/071756		弁理士 高木 千嘉
(87) 国際公開日	平成12年11月30日 (2000.11.30)	(74) 代理人	100105290
審査請求日	平成19年5月17日 (2007.5.17)		弁理士 三輪 昭次
(31) 優先権主張番号	60/135,853	(74) 代理人	100140132
(32) 優先日	平成11年5月25日 (1999.5.25)		弁理士 竹林 則幸
(33) 優先権主張国	米国 (US)		
(31) 優先権主張番号	09/577,634		
(32) 優先日	平成12年5月24日 (2000.5.24)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 遺伝子発現のレベルの変化を推定するための数学的解析

(57) 【特許請求の範囲】

【請求項1】

複数のアレイ・ハイブリダイゼーションにおける遺伝子発現レベルの差分を示す値を計算するための方法であって、

(a) それぞれのアレイ・ハイブリダイゼーションにおける遺伝子のハイブリダイゼーション・シグナルの強度に関連する実験ノイズに関する値を決定すること、

(b) その決定した実験ノイズ値を用いて、それぞれのアレイ・ハイブリダイゼーションにおける遺伝子それぞれの強度の分布値を示す第1の解析確率分布関数を定義することであって、ノイズはガウスのものであると仮定し、ベイズの定理を用いること、

(c) 前記第1の解析確率分布関数を用いて、差次的に発現される遺伝子の遺伝子発現の相違を示す第2の解析確率分布関数を導くこと、および

(d) 前記差次的に発現される遺伝子の前記第2の解析確率分布関数を適用して、前記アレイ・ハイブリダイゼーションから実験的に導かれた強度及びノイズ値を用い、遺伝子発現の相違に関連する値を決定することを含む方法。

【請求項2】

前記決定された値が遺伝子転写物の濃縮物中の推定倍率変化を含む、請求項1に記載の方法。

【請求項3】

前記決定された値が少なくとも1つの前記推定倍率変化と関連する品質メトリックを含む、請求項2に記載の方法。

10

20

## 【請求項 4】

前記品質メトリックが、推定倍率変化が1より大きいときに倍率変化が1未満である確率、および推定倍率変化が1未満であるときに倍率変化が1より大きい確率の少なくとも1つを表す、請求項3に記載の方法。

## 【請求項 5】

前記決定された値が特定の信頼区間で与えられる前記倍率変化の信頼限界を含む、請求項1に記載の方法。

## 【請求項 6】

実験ノイズに関する値が決定され、第1の解析確率分布関数が定義され、第2の解析確率分布関数が導かれ、そして、それぞれのアレイ・ハイブリダイゼーションのそれぞれの遺伝子に適用する、請求項1に記載の方法。

10

## 【請求項 7】

コンピュータに下記手順を実行させることで、複数のアレイ・ハイブリダイゼーションにおける遺伝子発現レベルの差分を示す値を計算させるコンピュータ・プログラムであって、

(a) それぞれのアレイ・ハイブリダイゼーションにおける遺伝子のハイブリダイゼーション・シグナルの強度に関連する実験ノイズに関する値を決定すること、

(b) その決定した実験ノイズ値を用いて、それぞれのアレイ・ハイブリダイゼーションにおける遺伝子それぞれの強度の分布値を示す第1の解析確率分布関数を定義することであって、ノイズはガウスのものであると仮定し、ベイズの定理を用いること、

20

(c) 前記第1の解析確率分布関数を用いて、差次的に発現される遺伝子の遺伝子発現の相違を示す第2の解析確率分布関数を導くこと、および

(d) 前記差次的に発現される遺伝子の前記第2の解析確率分布関数を適用して、前記アレイ・ハイブリダイゼーションから実験的に導かれた強度及びノイズ値を用い、遺伝子発現の相違に関連する値を決定することをコンピュータに行わせるためのコンピュータ・プログラム。

## 【請求項 8】

前記決定された値が遺伝子転写物の濃縮物中の推定倍率変化を含む、請求項7に記載のコンピュータ・プログラム。

## 【請求項 9】

前記決定された値が少なくとも1つの前記推定倍率変化と関連する品質メトリックを含む、請求項8に記載のコンピュータ・プログラム。

30

## 【請求項 10】

前記品質メトリックが、推定倍率変化が1より大きいときに倍率変化が1未満である確率、および推定倍率変化が1未満であるときに倍率変化が1より大きい確率の少なくとも1つを表す、請求項9に記載のコンピュータ・プログラム。

## 【請求項 11】

前記決定された値が特定の信頼区間で与えられる前記倍率変化の信頼限界を含む、請求項7に記載のコンピュータ・プログラム。

## 【請求項 12】

コンピュータに、実験ノイズに関する値を決定させ、第1の解析確率分布関数を定義させ、第2の解析確率分布関数を導かせ、そして、それぞれのアレイ・ハイブリダイゼーションのそれぞれの遺伝子に適用させるように作用する、請求項7に記載のコンピュータ・プログラム。

40

## 【発明の詳細な説明】

## 【0001】

## 【発明の分野】

本発明は、差次的遺伝子発現のレベルの定量的推定のための数学的解析に関する。より詳細には本発明は、所与の実験測定から推測することができる遺伝子発現のレベルのすべての倍率変化の事後分布の数学的導関数に関する。

50

## 【 0 0 0 2 】

## 【 発明の背景 】

細胞は、その数多くのタンパク質成分を広範囲にわたる機能に利用する。これらの機能には、たとえばエネルギーの生成、高分子成分すべての生合成、細胞構造の維持、細胞内および細胞外の刺激に対して作用する能力などがある。生体内のそれぞれの細胞は、生物が発現することができるタンパク質のレパートリーを生成するために必要な情報を含んでいる。この情報は生物のゲノム中に遺伝子として保存されている。特有なヒト遺伝子の数は、30,000~100,000であると推定されている。

## 【 0 0 0 3 】

いかなる所与の細胞でも、遺伝子セットの一部分だけがタンパク質の形で発現される。すべての細胞中に存在すると考えられるタンパク質もいくつかある（つまり至る所で発現される）。なぜならそれらは、細胞のそれぞれのタイプで必要とされる生物学的機能を果たし、「ハウスキーピング」タンパク質であると考えられるからである。対照的に他のタンパク質は、特定の細胞タイプにのみ必要とされる特化した機能を果たす。たとえば筋肉細胞は、筋肉の高密度収縮性繊維を形成する特化したタンパク質を含む。細胞の特異的な機能の大部分が発現している遺伝子によって決定されると仮定すれば、生物のゲノム内に保存されている遺伝情報をタンパク質に転換するプロセスの最初のステップである転写が、細胞の活性を調整し指示する制御ネットワークによって高度に制御されることは理にかなっている。

## 【 0 0 0 4 】

遺伝子発現の制御は、特定の機能用に形成されている細胞（たとえば筋肉細胞への特化）において明らかな活性、または状態（たとえば活発に増殖しているか休止しているか）を調べる実験において容易に観察される。したがって、細胞がその状態を変えると、この分子生物学的/生理学的「状態」のために必要とされる、タンパク質の調整されている転写を観察することができる。この非常に詳細な、細胞の転写状態の全体的な知識によって、細胞の状態およびこの状態を制御する生物学的システムに関する情報が提供される。たとえば、いつどんなタイプの細胞において未知の機能の遺伝子のタンパク質産物が発現されるかという知識によって、その遺伝子の考えられる機能に関する有用な手がかりが提供されるはずである。正常な細胞における遺伝子発現のパターンを決定することによって、ただ1つの受精卵からの成熟した生物の生成および分化に必要なとされる、制御システムが高度に調整された活性化および不活性化を実現する方法についての詳細な知識が与えられる可能性がある。正常および病原性細胞の遺伝子発現パターンを比較することによって、有用な診断用の「指紋」が提供され、治療による介入の妥当な標的となる異常な機能を同定する助けとなる可能性がある。

## 【 0 0 0 5 】

残念ながら、数多くの遺伝子の転写状態を決定する実験を行うための能力は、1つの実験中での数多くの遺伝子転写物の存在および量について調査する能力に限界があるために最近まで阻止されてきた。1つの限界は、同定される遺伝子の数が少ないことである可能性がある。ヒトの場合、ヒトゲノム中にエンコードされているタンパク質のうち、ある程度まで物理的に精製され定量的に特徴付けられているのはわずか数千である。他の限界は、転写分析の方法にある可能性がある。

## 【 0 0 0 6 】

近年の2つの技術的利点は、遺伝子転写の解析の助けとなっている。特定の組織中のmRNA転写物に由来する分子をクローニングし、次にこれらのライブラリーのメンバーのDNA端に高スループットの配列決定を適用することによって、発現される配列タグ(ESTs)のカatalogueが得られる。Boguski and Schuler, Nat. Genetics 10:369-370(1995)を参照のこと。これらの「シグナル配列」は、遺伝子の大きな一団の明確な同定手法を提供することができる。

## 【 0 0 0 7 】

さらにこれらの配列が由来するクローンは、生物学的サンプルからの転写物の定量化に

10

20

30

40

50

において使用することができる分析用の被験物となる。核酸ポリマー、DNAおよびRNAはコピー反応で生物学的に合成される。この反応では、1つのポリマーが対向する鎖の合成用の鋳型として働き、これがその相補体と名付けられる。他方からの鎖の分離（つまり変性）に続いて、これらの鎖を、詳細にはハイブリダイゼーションと呼ばれるプロセスで、相補的な配列を有する他の核酸の鎖とペアにすることができる。この特異的な結合が、特定のタンパク質の遺伝子産物を指定するmRNAなど、特定の核酸種の量を測定するための分析手順の基礎となる可能性がある。

#### 【0008】

第2の進展はマイクロアレイ/マイクロアッセイ技術に関する。これはハイブリダイゼーションをベースとするプロセスであり、多くの核酸種を同時に定量化することが可能である。たとえばDeRisi他のNat. Genetics 14:457-460(1996); Schena他のProc. Natl. Acad. Sci. USA 93:10614-10619(1996)を参照のこと。この技法は、個々の純粋な核酸種を少量ガラス表面上にロボットで配置（すなわちスポッティング）すること、このアレイに蛍光標識した核酸を複数個ハイブリダイゼーションすること、および結果として生じる蛍光標識されたハイブリッドをたとえば走査共焦点顕微鏡によって検出し定量化することの組合せである。転写物を検出するために使用すると、特定のRNA転写物（つまりmRNA）をDNA（つまりcDNA）にコピーすることができ、転写物のこのコピーされた形を、その後たとえばガラス表面上に固定する。

#### 【0009】

遺伝子発現データの解析における1つの問題は、他の実験における遺伝子の発現と比較したある実験中における遺伝子の発現レベルの全体的な倍率変化の推定である。遺伝子発現レベルにおける倍率変化のこの2つの生の測定値が与えられた場合、以前の方法で使用される最も単純な手法は、値の演算比を全体的な倍率変化の推定とみなすものであった。非常に強いシグナルの場合、これは基礎となるmRNA濃縮物における倍率変化の意味のある推定をもたらすが、弱いシグナルの場合は、使用される特定の実験システムに固有な「ノイズ」による汚染のために結果ははるかにあいまいになる。遺伝子発現レベルにおける倍率変化の推定のために以前に使用されていた他の技術は、特異的なシグナル強度に基づくものである（たとえばthe Affymetrix（登録商標）チップ）。しかしながら、前述の方法の使用によって発現レベルに割り当てられる値はマイナスであり、したがって遺伝子発現率がマイナスまたは不定であるという厄介な状況になる可能性がある。

#### 【0010】

##### 【発明の概要】

本発明は、マイクロアッセイ・プロトコールからの差次的遺伝子発現のレベルを定量化するための、非常に精度が高く再現性のある数学ベースの方法を提供する。

本発明の方法を使用して、遺伝子群の1つまたは複数のアレイ中の遺伝子発現のレベルの違いを計算することができる。この方法は、アレイ中のそれぞれの遺伝子のハイブリダイゼーション・シグナルの強度に関する実験ノイズを定義するものである。実験ノイズは、生物学系において見られる発現レベルの変化である生物学的ノイズではなく、観察レベルにおけるチップまたは他のマイクロアレイの変化である。遺伝子の検出は、必ずそうではないが、蛍光性に基づくことが多い。他の検出システムも使用されており、これに適合させることができる。このようなシステムには発光または放射性標識、標識されたプローブの容易な検出を可能にするビオチン化、ハプテン化、または他の化学タグがある。

#### 【0011】

数学的な記載については、以下のセクションI - ノイズ・モデルの公式化を参照のこと。ノイズには、ガウスおよびベイズの定理が適用されると考えられる。次いで明確な実験ノイズ項、シグマを使用して、解析（つまり、それが連続関数であるという数学的な意味で解析的である）確率分布関数（「pdf」）を定義し、これは遺伝子それぞれの強度の分布値を示す。これらのpdfを使用して解析同時pdfを導関数し、これはアレイ中で特異的に発現する任意の遺伝子または遺伝子産物の予想比率または倍率変化を示す。アレイ上の遺伝子群から実験的に導かれた強度およびノイズ値を使用して、同時pdfを適用す

る。これは(1) 遺伝子転写物中の濃縮物の倍率変化を推定するため、(2)  $j p d f$  を使用して特異的な信頼区間で与えられる倍率変化の信頼限界を確立するため、および(3) 倍率変化の評価に関するP-値、または倍率の変化の推定に関係のある品質メトリック(推定値が1より大きいときは倍率の変化が1未満である確率、または推定値が1未満であるときは倍率の変化が1より大きい確率)を導くためである。本発明の方法によって決定される推定の倍率変化は、観察される遺伝子発現のレベルの違いを示す。転写物の濃縮物が0に向かうときでも、全体の分散(つまりノイズ)は依然として大きい可能性がある。本発明の方法は数学的な式を使用して、アレイ中にある1つまたは複数の細胞または組織タイプの遺伝子発現のレベルの得られた測定値から導くことができる、遺伝子発現のすべてのレベルの事後統計分布について記載する。

10

#### 【0012】

マイクロアレイは、空間的に分離している組織内の担体材料上に位置する二本鎖または一本鎖DNA分子の整然としたアレイである。典型的にはニトロセルロースの大きなシートであるフィルタ「マイクロアレイ」とは対照的に、マイクロアレイでは10000個までのDNA分子が典型的には1~4平方センチメートルの領域内に適合することができるように、DNAがより高密度に詰め込まれた組織に配置されている。フィルタアレイのニトロセルロースベースの材料とは対照的に、マイクロアレイは典型的には被覆されたガラスを固形担体として使用する。DNAサンプルの整然としたアレイを有することによって、サンプルそれぞれの位置を突き止め、アレイ上のDNAが生成した元のサンプルに連結させることができる。マイクロアレイを調製するための方法および装置が記載されている。米国特許第5,445,934号および5,800,992号を参照のこと。これらは両方とも参照によって本明細書に取り込まれている。

20

#### 【0013】

マイクロアレイ上のDNAサンプルは、プローブ・サンプルがマイクロアレイ上のDNAサンプルと同等または同一の分子を含んでいるかどうかを同定するために蛍光的に標識されている、RNAまたはDNAプローブとハイブリダイズされる。適切な条件下では、プローブ分子はマイクロアレイ上のDNA分子にハイブリダイズする。一般に、生産性ハイブリッドからの同一またはほぼ同一な配列である。DNA-プローブ・ハイブリッド分子の存在は、蛍光検出装置によって検出される。ハイブリダイゼーション・シグナルが弱いまたは特定のDNA部位において存在しない場合、したがってプローブ中に対応するDNAまたはRNA分子はない。現在のマイクロアレイ用装置は、4つのまでの異なる蛍光プローブ・サンプルを1度にハイブリダイズすることができる。技術を改良することによって、より多くのプローブを1度にハイブリダイズすることができる。

30

#### 【0014】

最近まで、DNAのハイブリダイゼーションはニトロセルロース・フィルタ上で行われていた。マイクロアレイ上にDNAが直接スポットされるマイクロアレイとは対照的に、フィルタ上に細菌コロニーをスポットし、寒天成長培地上にフィルタを置き、そのフィルタを細菌コロニーの増殖を助長する条件下でインキュベートすることによってフィルタアレイは生成する。コロニーを溶解しフィルタを処理しDNAをフィルタ材料に固定することによって、細菌コロニー内部のDNAは切り離される。細菌のフィルタアレイを生成するプロセスは、典型的には2~4日かかる可能性がある。マイクロアレイには、フィルタアレイの方法に勝るいくつかの利点がある。たとえば一般にフィルタ方法では、内部にクローン化されたcDNAが含まれる細菌のコロニーを配列する。コロニーは数日にわたって生育され、溶解されてDNAを切り離しフィルタ上にDNAを固定しなければならない。細菌の残骸およびコロニーから切り離されるDNAの量が少ないために、コロニーのフィルタアレイのハイブリダイゼーションは確かなものではない。第2の利点は、マイクロアレイに関する相互作用がフィルタよりも速いことである。これは、フィルタ上でコロニーを生育しハイブリダイゼーションの次のラウンド用にそれらを調製するために必要とされる時間のためである。対照的に後のマイクロアレイの探索は、アレイの分析が終了した24時間後以内に始めることができる。マイクロアレイの他の利点は、蛍光的に標識され

40

50

たプローブを使用するための能力である。これによって、ハイブリダイゼーション検出用の非放射性の方法が提供される。対照的に、一般にフィルタのハイブリダイゼーションでは、放射性リン酸またはイオウで標識されたプローブを使用する。マイクロアレイを複数のプローブと同時にハイブリダイズすることができる。対照的にフィルタアレイは、1度に1つのプローブとしかハイブリダイズすることができない。マイクロアレイの最も重要な利点の1つは、ハイブリダイゼーション・シグナルの再現性および感受性である。典型的には、マイクロアレイのハイブリダイゼーション・シグナルは、フィルタアレイと比べて高く感受性が大きい。さらにフィルタアレイは、プローブとフィルタ上のDNAの間の生産的なハイブリダイゼーションとは関係のない偽のバックグラウンド・シグナルをしばしば示す。

10

## 【0015】

核酸断片のランダムなサンプルがひとたびマイクロアレイの固体表面(たとえばガラス)に固定されると、次いでその核酸断片のランダムなサンプルを、当該の遺伝子または配列に相補的な1つまたは複数の標識されたプローブにハイブリダイズすることができる。一般に、ハイブリダイズされないプローブは取り除かれる。次いでその標識されたプローブは、当分野で知られている方法(たとえば共焦点顕微鏡)によって検出される。たとえばスライド像を、スポット発見分析、局所的な背景の決定、スポット中のシグナル強度の分布、およびノイズ比へのシグナル用のアレイビジョン像分析ソフトウェア(Imaging Research)によって分析することができる。次いで統計的な評価を以下に記載するように行う。

20

## 【0016】

本発明は数学に基づく方法を使用して、差次的に発現される遺伝子のレベルにおける倍率変化を定量化する。詳細には本発明は、ベイズのフレームワークに根拠を有する単純な演繹手法を使用して、差次的遺伝子発現の数学的解析において使用された以前の方法の発見的方法に基づく制限を回避する。本発明は、至急に遺伝子発現のレベルにおける倍率変化の点的な推定値を探索するのではなく、所与の測定値から推測することができる差次的遺伝子発現のすべての倍率変化の事後分布の数学的な式を導く。この事後分布から以下の情報を得ることができる。(i) 遺伝子発現のレベルの倍率変化の推定量、(ii) 任意の所与の信頼度レベルにおける、倍率変化の信頼限界、および(iii) 変化の統計的な有意性を評価するためのP-値。本発明のさらなる利点は、両方のシグナルがゼロまたはマイナスの場合、発見的方法の閾値に頼ることなく、倍率変化の推定値および信頼限界がシグナルペアに割り当てられる可能性もあるということである。ゆえに、本明細書で開示される数学的フレームワークは、所与のサンプル内のすべてのシグナルの推定値を1つにする。

30

## 【0017】

## 【発明の詳述】

図10を参照すると、コンピュータ・システム102はプロセッサ104、メモリ106、ディスク・ドライブ108、ディスプレイ110、キーボード112、およびマウス114を含む。プロセッサ104は、Intel(登録商標)Corporationによって作成されるPentium(登録商標)IIIプロセッサなどのパーソナル・コンピュータの中央演算処理装置(CPU)であってよい。メモリ106は、ランダム・アクセス・メモリ(RAM)および読み取り専用メモリ(ROM)を含む。ディスク・ドライブ108はハードディスク・ドライブを含み、フロッピーディスク・ドライブ、CD-ROMドライブ、および/またはジップ・ドライブを含んでもよい。ディスプレイ110はブラウン管(CRT)であるが、他の形のディスプレイ、たとえばTF Tディスプレイを含めた液晶ディスプレイ(LCD)も許容される。キーボード112およびマウス114は、ユーザ(図示せず)のためのデータ入力機構となる。コンポーネント104、106、108、110、112および114は母線116によって接続されている。コンピュータ・システム102は、たとえばメモリ106中に、プロセッサ102を制御するための命令を含むソフトウェア・コードを保存して以下に記載する機能を行うことができる。

40

## 【0018】

50

## 1. ノイズ・モデルの公式化

所与の遺伝子の発現レベルの測定値は、異なる実験または繰り返しの実験において  $x$  に基づいており、以下のように書くことができる。

$$x = C n + \quad (1)$$

上式で  $n$  は溶液中の遺伝子の mRNA の絶対的物理濃度（モル濃度）であり、 $C$  はモル濃度を記録される強度と関係付ける一定の比例定数であり、 $\quad$  はノイズ項である。以下の等式では、絶対的な mRNA の濃度を決定することは求められず、したがって簡潔性のために以下では  $C = 1$  に値を設定する。

$$x = n + \quad (2)$$

等式 (2) では、ノイズ項は 3 つの異なる因子に分解することができる。

$$= \quad + \quad + \quad \quad (3)$$

上式で  $\quad$  はバックグラウンド強度の変動からの変数であり、 $\quad$  は他の mRNA（特異的または非特異的）のクロス・ハイブリダイゼーションから生じる項であり、 $\quad$  はオリゴヌクレオチドまたは cDNA 密度、および他の同様の因子のチップ間の変動から生じる「比例変数」項である。たとえば、最終シグナル  $x$  がいくつかのディファレンシャル・シグナルの平均化から得られる Affymetrix（登録商標）チップの場合、それぞれのノイズ項はプラスまたはマイナスであってよく、平均して約 0 である。したがって、バックグラウンドおよびクロス・ハイブリダイゼーション項をただ 1 つのノイズ項に分類することができる。

$$c = \quad + \quad \quad (4)$$

### 【0019】

ノイズ全体の平均および分散は以下のように書かれる。

$$\langle \quad \rangle = 0 \quad (5)$$

$$c^2 = \text{var}(\quad) + \text{var}(\quad) = \quad^2 + (\quad n)^2 \quad (6)$$

上式で  $\text{var}(\quad) = \quad^2 = (\quad n)^2$  となるような比例変数パラメータである。比例変数項  $\quad$  は変数  $c$  の係数と同等である（Chen 他 J. Biomed. Optics 2: 364 (1997) によって最初に定義された）。他のノイズ項もこの方法に介入するので、ノイズの分散の合計は  $n = 0$  のときでも大きいままである可能性がある。以下の等式では、 $\quad$  は正規分布していると考えられる。

### 【0020】

一例として Affymetrix（登録商標）チップを使用して、不在 / 存在決定アルゴリズムによってシグナルされる遺伝子すべての測定値  $x$  の分散が不在であると考えることにより、組み合わせさせたバックグラウンドおよびクロス・ハイブリダイゼーションノイズの分散  $\quad^2$  を計算する。繰り返しの実験における強度の最も高い四分位数を比較することによって、比例定数項は推定されている。

### 【0021】

典型的な Affymetrix（登録商標）チップベースの実験用の等式 (6) 内の項の規模を例示するために、発現レベルの平均値  $x$  は  $\text{Med}(x) = 80$ 、 $\quad = 2.5$  および  $\quad = 0.25$  である。したがって、ノイズ比  $\text{Med}(x) / \quad$  を超えるメジアン・シグナルは、約 3 つしかない。バックグラウンド・ノイズの標準偏差のみで  $\quad = 3 \sim 4$  であるので、バックグラウンド・ノイズを約 1 桁上回るクロス・ハイブリダイゼーション・ノイズでは  $\quad = \quad$  である。

### 【0022】

可変性の  $\quad$  は比較的大きく  $\quad$  は小さいので、等式の右手側に  $n = x$  と書くことによって、等式 (6) をわずかに簡略化することができる。

$$\quad^2 = \quad^2 + (\quad x)^2 \quad (7)$$

したがって、ノイズの分散を定量的に推定するために、前もって基本濃度を覚えておく必要はない。

### 【0023】

## 2. 濃度の事後分布

10

20

30

40

50

等式(2)は濃度に関する測定値を与えるが、本発明では測定値の関数として濃度を得ている。変数  $n$  および  $x$  についてベイズの定理を書くことによって、確率の点でこれを公式化することができる。

【数1】

$$P(n|x) = \frac{P(x|n)P(n)}{P(x)} \quad (8)$$

等式(8)において、 $P(x|n)$ は $x$ についての条件付確率分布関数(pdf)であり、 $n$ に関する条件付で、 $P(n)$ は $n$ の事前分布であり(したがって実際に測定が行われる前の $n$ の知識の状態を反映する)、 $P(x)$ 、 $x$ についてのpdfは本質的に正規化項である。それゆえ、等式(2)およびガウスのノイズの仮定から、以下の等式を導くことができる。

【0024】

【数2】

$$P(x|n) = \frac{1}{(2\pi\sigma_c^2)^{1/2}} \exp\left(-\frac{(x-n)^2}{2\sigma_c^2}\right) \quad (9)$$

上式で  $\sigma_c(n)$ 、等式(6)。

分布 $P(n)$ には、事前知識として、濃度は必ずプラスであるという事実のみが使用される。

【0025】

【数3】

$$P(n) = \begin{cases} 0 & n < 0, \\ \mu^{-1} n^{\mu-1} & n \geq 0, \end{cases} \quad (10)$$

上式で極限 $\mu \rightarrow 0$ は非常に短い。(これは、極限 $\mu \rightarrow 0$ におけるステップ関数分布を得るための単なる細工であり、この間 $P(n)$ を常時積分可能に保つ)。最後に、 $P(x)$ は以下の積分式によって得られる。

【数4】

$$P(x) = \int_{-\infty}^{\infty} dn P(n) P(x|n) \quad (11)$$

【0026】

極限 $\mu \rightarrow 0$ の場合、等式(8)を以下の形式に書き直すことができる。

【数5】

$$P(n|x) = \frac{P(x|n)}{\hat{P}(x)}, \quad n \geq 0 \quad (12)$$

この場合 $P(x|n)$ は等式(9)によって与えられ、ここで分母は以下のものである。

【数6】

$$\hat{P}(x) = \int_0^{\infty} dn P(x|n) \quad (13)$$

エラー関数を使用して、等式(13)を容易に評価することができる。濃度の推定値に関する

る等式(12)の結果を直接探索するのではなく、倍率の変化の分布を定量化するために、それを以下で使用する。

【 0 0 2 7 】

### 3 . 倍率の変化の事後分布

所与の遺伝子に関して、2つの所与の実験(たとえば実験1および2)の間の遺伝子発現のレベルにおける倍率の変化を評価することが望まれていると考えられる。たとえば、実験中のmRNA濃度がそれぞれ $n_1$ および $n_2$ である場合、実験1と比較した実験2の濃縮物の倍率の変化 $R$ はしたがって以下の式で与えられる。

【数7】

$$R = \frac{n_2}{n_1} \quad (14)$$

10

等式(14)では $n_1$ および $n_2$ への直接のアクセスはないが、ベイズの項における $R$ の推定値を、 $R$ の事後分布を以下のように書くことによって即座に公式化することができる。

【 0 0 2 8 】

【数8】

$$f_R(R|x_1, x_2) = \int_0^\infty dn_1 \int_0^\infty dn_2 \delta\left(\frac{n_2}{n_1} - R\right) P(n_1|x_1) P(n_2|x_2) \quad (15)$$

20

上式で $x_1$ および $x_2$ はそれぞれ実験1および2における強度の測定値であり、 $\delta$ はDiracデルタ関数のことであり、 $P(x|n)$ は前述の等式(12)によって与えられる。

等式(15)に示される積分を実施することは非常に直接的であり、どちらかというやや退屈な作業である。 $R$ の分布関数(つまり $f_R(R|x_1, x_2)$ )の $x_1$ および $x_2$ に直接従属している)は以下の式

【数9】

$$f_R(R) = \frac{C(x_1)C(x_2)}{2\pi\sigma_1\sigma_2} \exp\left(-\frac{x_1^2(R-R_0)}{2(\sigma_2^2 + R^2\sigma_1^2)}\right) I(x_1, x_2) \quad (16)$$

30

によって得られ、上式で $\sigma_i^2 = \sigma_i^2(x_i)$ 、 $I = 1, 2$ 、ここで $\sigma_i(x)$ は等式(7)によって与えられ以下の正規化項を有する。

【数10】

$$C(x) = \frac{2}{1 + \operatorname{erf}(x/\sqrt{2}\sigma_i(x))} \quad (17)$$

40

上式で $\operatorname{erf}$ はエラー関数であり、(Abramowitz, M. and Stegun, I.A., p.297 Handbook of Mathematical Functions (Dover, New York, 1972)を参照のこと)以下の定義式を有する。

【 0 0 2 9 】

【数11】

$$I = \sigma_{12}^{-1} \exp\left(-\frac{\alpha_{12}^2}{2\alpha_{12}^2}\right) + \alpha_{12}(2\pi\sigma_{12}^2)^{1/2} \frac{1}{2} \left(1 + \operatorname{erf}\left(\alpha_{12}/\sqrt{2}\sigma_{12}\right)\right) \quad (18)$$

50

上式で  
【数 1 2】

$$\frac{1}{\sigma_{12}^2} = \frac{1}{\sigma_1^2} + \frac{R^2}{\sigma_2^2} \quad (19)$$

$$\alpha_{12} = \left( \frac{x_1}{\sigma_1^2} + \frac{Rx_2}{\sigma_2^2} \right) / \left( \frac{1}{\sigma_1^2} + \frac{R^2}{\sigma_2^2} \right) \quad (20)$$

10

である。

見た目がやや複雑ではあるが、等式(16)にはわずかに2つの単純な極限しかなく、それらは2つのシナリオを使用することによって以下で論じられるはずである。

【0030】

3.1 ケース1 - 濃度が高い場合

両方の実験において、ノイズの標準偏差と比較してRNA濃度が大きい場合、結果として  $x_i > > (x_i)$ 、 $i = 1, 2$ 、 $R$  はほぼ正規な分布を有する。

【数 1 3】

$$f_R(R) \approx \frac{1}{(2\pi\sigma_R^2)^{1/2}} \exp \left( -\frac{(R - R_0)^2}{2\sigma_R^2} \right) \quad (21)$$

20

この極限では、 $R$  の平均は単なる測定値の比である。

【数 1 4】

$$\langle R \rangle = R_0 = \frac{x_2}{x_1} \quad (22)$$

したがって、 $R$  の分散  $\sigma_R^2$  は以下の式で与えられる。

【数 1 5】

$$\sigma_R^2 = \frac{\sigma_2^2 + x_2^2 \sigma_1^2 / x_1^2}{x_2^2} \quad (23)$$

30

等式(7)を代わりに使用し、 $R$  の標準偏差の単純な近似値を得ることができる。

【0031】

【数 1 6】

$$\sigma_R = \sqrt{2\alpha} R_0 \quad (24)$$

40

したがって、高濃度の極限では(つまりケース1)、測定値の比と比較した実際の倍率変化の標準偏差は以下の定数によって与えられる。

【数 1 7】

$$\frac{\sigma_R}{R_0} = \sqrt{2\alpha} \quad (25)$$

調べてみると、シグナルがどんなに大きくても、測定される全体の倍率の変化の約 2 (  $\alpha = 0.25$  の場合  $\pm 35\%$  ) という倍率の変化の推定において既約の変数が残ることを等式(24)は示す。

50

【 0 0 3 2 】

## 3.2 ケース2 - 濃度が非常に低い場合

両方の実験において、ノイズの標準偏差と比較してRNA濃度が非常に低く  $x_1$  (  $x_2$  )、  $I = 1, 2$  である場合、したがって分布は以下の「ユニバーサル」な形をとる。

【数18】

$$f_R(R) \approx \frac{1}{\pi} \frac{1}{1+R^2} \quad (26)$$

上式では簡潔性のために  $\sigma_1 = \sigma_2$  であると仮定する。

10

【 0 0 3 3 】

この極限では、Rの分布は濃度とは完全に独立しており、その影響はノイズによって支配されている。等式(26)はいわゆるコーシー分布を定義し (Keeping, E.S., Introduction to Statistical Inference (Dover, New York, 1995) を参照のこと)、それは非常に広域であり明確な平均値を有していない。コーシー分布の1つの固有の「病的な」性質は、使用されるサンプルの合計数に関係なく、多くの独立サンプルの平均値はいずれの意味でも1つの数字ではなく、等式(26)に従って分布したままであるということである。これとは逆に、コーシー分布のメジアンはちょうど1であり、サンプルのメジアン範囲は1までであり、メジアンに関係のある前述の病原体はない。

【 0 0 3 4 】

20

Rの累積分布関数は以下の式によって与えられる。

【数19】

$$P(R \leq \rho) = \frac{2}{\pi} \tan^{-1} \rho \quad (27)$$

たとえば、90%の信頼限界は(0.16, 6.3)であり、等式(26)の分布が非常に広いことを示す。なぜなら、ノイズと比較してシグナルが充分弱いとすれば、これらの範囲は  $R_0 = 1$  のときでも得られるからである。

最後に、変換  $\mu = \log R$  の下では、等式(26)の分布関数は完全に対称的になる。

30

【数20】

$$f_u(u) = \frac{1}{\pi} \frac{1}{\cosh(u)} \quad (28)$$

したがって、いくつかの例では対数表示が有用である可能性があるが、変換のさらなる使用が本明細書においてさらに求められるべきではない。

【 0 0 3 5 】

一連の対 ( $x_1, x_2$ )、両方のノイズ項  $\sigma_1 = \sigma_2 = 20$  の一定の標準偏差の事後分布  $f_R(R)$  を図1は示す。この図中では、比  $x_2 / x_1$  は常に4であるが (両方のシグナルが0の場合を除いて)、シグナルとノイズの比は非常に変動的である。最も高いシグナル・レベルでは、( $x_1, x_2$ ) = (100, 400) であり、 $f_R(R)$  は  $R = 4$  付近で強烈にピークに達する。しかしながらこの極限においてさえも、視覚的な観察によって、68%の信頼区間 (正規分布の2つの標準偏差の幅に対応する) が約 (3, 5) であることが示される。このことは、最も小さいシグナルとノイズの比が  $100 / 20 = 5$  であるときでも、実際の倍率変化を  $3 < R < 5$  を超える値であると推測することはできないことを示す。

40

【 0 0 3 6 】

シグナルとノイズの比が減少すると、分布  $f_R(R)$  が広がるだけでなく、そのピーク・シフトが低下する。したがって図1では、測定値対 (40, 10) については、分布のメジアンは約2.2であり、発生する実際の最大値は1という値に非常に近い。分布関数の

50

この広がりおよびシフトによって、弱まっているシグナルでは、測定値の比は実際の倍率変化の非常に当てにならない指標となることが示される。最後に、両方ともゼロ(0, 0)である測定値の極限では、等式(26)の復元によって分布が非常に広いこと、メジアン  $R = 1$  および  $R = 0$  でピークであることが示される。

【0037】

図2は、以下の構成を使用することによって、前に図1において定量化された分布  $f_R(R)$  の状態を定性的に示す。(i) 値  $(x_1, x_2)$  のそれぞれの対について、マイナス軸に沿った領域を除き、各次元において  $\pm \epsilon$  の程度で平面状のポイント  $(x_1, x_2)$  付近にボックスを描く。(ii) 次にボックス中のすべてのポイントに原点から1組の線を引き、これらの線の傾きの分布は、事後分布  $f_R(R)$  を表す。

10

【0038】

4. ベイズによる倍率変化の推定

ベイズによる倍率変化  $R$  の推定は、等式(16)および測定値  $x_1$  および  $x_2$  の知識に基づいて行うことができる。最初に、累積分布関数を定義する。

【数21】

$$F(R) = P(R \leq R) = \int_0^R f_R(R) dR \quad (29)$$

$F(R)$  は数値の積分の使用によって評価されることが好ましい。 $F(R)$  の数値の値に基づいて、以下の情報を容易に得ることができる。

20

【0039】

4.1 倍率の変化の推定量  $R$

倍率の変化の推定量  $R$  として、メジアンの推定量を選択した。

【数22】

$$\hat{R} = \text{Med}(R) \quad (30)$$

つまり、 $F(R)$  の  $R$  の値 =  $1/2$ 。

たとえばMAP(事後確率の最大値)または平均推定量などの、他の推定量も考えられる。たとえばVan Trees, H.L., Detection, Estimation and Modulation Theory. Part I (John Wiley and Sons), New York, 1998)を参照のこと。しかしながら、平均推定量はここでは使用しない。なぜなら、 $f_R(R)$  は明確な平均値を有していないからである(つまり、等式(21)の正規に近い極限においてさえも、これは常に従属率  $1/R^2$  の「テール」を有しているはずである)。したがって、メジアンの推定量には変換  $(R \rightarrow 1/R)$  下で確固性および対称性という2つの利点があり、本明細書に取り入れられているものである。形式的にはメジアンの推定量は、(推定 - 実際値)エラー項の絶対値を減少させる、たとえば最小限にするものである。たとえばVan Trees, H.L., Detection, Estimation and Modulation Theory. Part I (John Wiley and Sons), New York, 1998)を参照のこと。

30

40

【0040】

4.2 信頼限界  $R_p$  および  $R_{1-p}$

$p < 1$  であるとする、信頼限界  $R_p$  および  $R_{1-p}$  は対応する百分位数の値として定義される。

$$F(R_p) = p \quad (31)$$

$$F(R_{1-p}) = 1 - p \quad (32)$$

【0041】

4.3 変化の有意性の  $P$ -値

$R > 1$  (「1と比較した実験2において発生する有意なプラスの倍率変化」という仮説を、 $R < 1$  でありこれを有意な変化についての仮説の  $P$ -値  $P$  として定義する、相補的

50

な仮説の確率を評価することによって試験することができる。これは以下のように単純に表せる。

$$P = F(R = 1) \quad (33)$$

図1と関連して論じられるすべての測定値対に関する結果は、 $p = 0.16$ によって決定される信頼限界で以下の図4においても例示される。P-値を有することは、有意であるとみなされるこれらの測定値対のみを保持するための強力な選択基準となることに留意されたい。したがって、図4に示される測定値の比はすべて4に等しい(ただし(0, 0のコースは除く))が、最初の3つの記載事項((100, 400)、(50, 200)、(25, 100))のみが有意な変化(つまり、0.05の信頼レベルで)を示すことが分かっている。代わりに、これら前述の表形式の記載事項それぞれについて、倍率変化の信頼限界が知られている。したがって、測定値対(25, 100)では、 $R = 3.6$ という推定値が値(2.0, 8.8)によって括弧でくくられ、この実施例においてはここで報告される区間より正確には、倍率変化を確認することができないことが示される(つまり、2ほど小さく、8.8ほど大きい実際の倍率変化はデータと一致する)。

10

【0042】

5. 強度の対( $x_1, x_2$ )の( $R, P$ )平面への写像

等式(30)および等式(33)によって、強度の対( $x_1, x_2$ )を数値の対( $R, P$ )に写像する。倍率変化の有意性が重み付けされた表像を提供する、この写像を図3に例示する。簡潔性のために、 $\alpha_1 = \alpha_2$ を選択した。

図3は、( $\log(R), P$ )平面における定数 $x_1$ および定数 $x_2$ のラインを示す。 $R$ ( $1/R$ )(つまり、 $\log(R) - \log(R)$ )の相互変換の下では、その形状は対照的である。なぜなら、この特定の実施例のために値 $\alpha_1 = \alpha_2$ を選択したからである。所与のRに関してPの範囲は明確であり、上限 $P_u(R)$ は以下のように表される。

20

【数23】

$$0 \leq P \leq P_u(\hat{R}) \quad (34)$$

$R$ 、または $R = 0$ のとき $P_u(R) = 0$ なので、明確な0でない範囲に常に存在するとしても、大きな倍率変化は小さなP-値と必ず相関関係がある。

30

$P_u(R)$ に関する式は以下の式によって与えられる(項 $P_u(R)$ の導関数に関する付録Aの等式(65)を参照のこと)。

【0043】

【数24】

$$P_u(\hat{R}) = \begin{cases} \text{erfc} \left( t_m \frac{(\hat{R}^2 + \sigma_2^2 / \sigma_1^2)^{1/2}}{(1 + \sigma_2^2 / \sigma_1^2)^{1/2}} \right), & \hat{R} \geq 1 \\ \text{erfc} \left( t_m \frac{(1 / \hat{R}^2 + \sigma_1^2 / \sigma_2^2)^{1/2}}{(1 + \sigma_1^2 / \sigma_2^2)^{1/2}} \right), & \hat{R} \leq 1 \end{cases} \quad (35)$$

40

上式で $\text{erfc}$ は相補的なエラー関数であり(Abramowitz, M. Stegun, I.A., p. 297 Handbook of Mathematical Functions (Dover, New York, 1972)を参照のこと)、 $t_m = 0.477$ である。厳密には等式(35)は、

**R □ 1 または R ■ 1**

について漸近的にのみ正当性があるが、実際には図3において見ることができるよう、すべてのR値のための優れた近似値を提供する。

50

## 【 0 0 4 4 】

$R > 1$  の場合、境界  $P = P_u(R)$  は  $(x_1, x_2)$  平面のライン  $x_1 = 0$  (つまり  $x_2$  軸) に対応する。この境界上の点は、所与の倍率変化  $R$  が少なくとも有意であるための点である (つまり  $P$  の最大値を有する)。定数  $R$  のラインは  $(x_1, x_2)$  平面内の弧に対応しており、これらはすべて  $x_2$  軸を基準とし ( $R > 1$  の場合)、ここで点  $P$  は最大値であり、したがってライン  $x_2 = R x_1$  への漸近線であり、 $P$  は急速に 0 になる傾向がある。

## 【 0 0 4 5 】

## 6. P F O L D アルゴリズムの実施

前に記載した推定スキームは、P F O L D と呼ばれる C + + プログラムにおいて実施されている。ノイズ項の 2 つの強度および対応する標準偏差を特定する、所与の入力パラメータのセット  $(x_1, x_2, \sigma_1, \sigma_2)$  に関して、P F O L D は最初に分布関数  $f(R)$  を数学的に評価する (明確な範囲  $R_{min} \leq R \leq R_{max}$  の通常の組み合わせ  $R_i = R_{min} + i \cdot (R_{max} - R_{min}) / N$ ,  $i = 0, 1, \dots, N$  の点における等式(16)を参照のこと、ここで  $R_{min}$ 、 $R_{max}$  および  $R$  は自動的に選択されて関数の変数をすべて得る (図 1))。次いで  $f_R(R)$  の数値積分法によって累積分布関数  $F(R)$  (等式(29)を参照のこと) を見つけ、等式(30)、(31)、(32) および(33)をそれぞれ順に解くことによって、次にセクション 5 の推定値をすべて (つまり倍率変化  $R$ 、信頼限界  $(R_p, R_{1-p})$ 、および  $P$ -値  $P$  を容易に評価することができる。これら前述の等式の根を見つける際には、単純な二分法を使用した。たとえば Press, W., 他の Numerical Recipes in C, 2nd Edition, p.353 (Cambridge University Press, Cambridge, 1997) を参照のこと。

## 【 0 0 4 6 】

## 7. モンテカルロのシミュレーション

式のデータの解析における主な問題点は、有意ではない倍率変化から有意なものを分離することである。このプロセスにおける統計値  $(R, P)$  の有用性を評価するために、一連のモンテカルロのシミュレーション (たとえば Cowan, G., Statistical Data Analysis, p.41 (Clarendon Press, Oxford, 1998) を参照のこと) を実施した。これは実際の実験値を近似するためのものであった。対数正規分布に従って (たとえば Cowan, G., Statistical Data Analysis, p.34 (Clarendon Press, Oxford, 1998) を参照のこと) 以下の式によって計算することによって、濃度の値  $n$  を発生させた。

## 【 0 0 4 7 】

$$n = \exp(y) \quad (36)$$

上式で  $y$  はガウスのランダムな変数であり、以下のパラメータと共に発生する。

$$\langle y \rangle = 7.25 \quad (37)$$

$$\sigma_y = 1.22 \quad (38)$$

上式で  $\langle y \rangle$  および  $\sigma_y$  はそれぞれ、 $y$  の平均値および標準偏差である。等式(37)および(38)のパラメータは、それぞれ以下の値を有する第 25 百分位数、メジアンおよび第 75 百分位数に分布する結果となる。

## 【 0 0 4 8 】

$$n_{25} = 618 \quad (39)$$

$$n_{50} = 1408 \quad (40)$$

$$n_{75} = 3208 \quad (41)$$

濃度に関する対数正規分布の選択は事後観察によって決定され、実際の実験値において存在する遺伝子の強度の分布は、ほぼ対数正規分布である。Affymetrix (登録商標) チップについては、抗体染色手順次いでハイブリダイゼーションによって、等式(39)、(40) および(41)において百分位数によって示される強度が典型的である。

## 【 0 0 4 9 】

等式(36)によって発生する  $n$  のそれぞれの値については、 $b$  の真の倍率変化を、以下の 2 つの強度の値を計算することによってノイズと組み合わせてシミュレートした。

$$x_1 = n + \sigma_1 \quad (42)$$

$$x_2 = b n + \sigma_2 \quad (43)$$

上式でノイズ項  $\epsilon_1$  および  $\epsilon_2$  はガウスのランダムな変数とは相関関係がなく、平均値  $\langle \epsilon_1 \rangle = \langle \epsilon_2 \rangle = 0$  であり、等式(6)によって与えられる標準偏差は以下のパラメータを有する。

$$\sigma_c = 600, \quad \sigma = 0.25 \quad (44)$$

最後に、等式(42)および(43)によって計算した強度 ( $x_1, x_2$ ) から、対応する推定値 ( $R, F$ ) を等式(30)および(33)によって計算した。

【0050】

7.1 クラスの割り当て

明確さの程度を確実にするために、2組のシミュレーションを実施し、以下に与える倍率の変化をした遺伝子のクラスをそれぞれで定義した。

クラス0：変化なし、 $b = 1$ 。

クラス1：変化あり、 $b = 3$ 。

P F O L D を使用して2つのクラス間の遺伝子を区別することの有効性を評価して、クラス1に属する遺伝子を選択した。この評価を行うために、( $R, F$ ) 平面内の受容領域  $D$  を定義し(たとえばCowan, G., Statistical Data Analysis, p. 47 (Clarendon Press, Oxford, 1998) を参照のこと)、遺伝子のクラスのメンバーシップに関する以下の予測値  $\pi$  をさらに定義する。

【0051】

【数25】

$$\pi = \begin{cases} p & \text{if } (\hat{R}, P) \in D_1 \quad (\text{すなわち、クラス1に割り当てられた遺伝子}) \\ \alpha & \text{if } (\hat{R}, P) \in D_0 \quad (\text{すなわち、クラス0に割り当てられた遺伝子}) \end{cases} \quad (45)$$

上式で  $p$  および  $\alpha$  はそれぞれ、受容領域中のプレゼントおよびアセントを表す。受容領域  $D$  の一例は、以下の式によって定義される長方形を決定する表面である。

【数26】

$$D = \{ \hat{R} \geq R_c, P \leq P_c \} \quad (46)$$

しかしながら、より一般的な領域も考えられるはずである。

【0052】

$D$  のいかなる選択によっても、以下の確率の推定値を導くことが可能である。

$P(p | 0)$  = クラス0の遺伝子がクラス1に割り当てられる確率、

$P(\alpha | 0)$  = クラス1の遺伝子がクラス0に割り当てられる確率、

これはそれぞれの場合において、 $D$  に属するかあるいは属さないインスタンス変数 ( $R, F$ ) の数を単に計測することによって行われる。遺伝子の大きなセット中のいくつかは倍率の変化なしであり、他のものが倍率変化3である事前確率用の値が代わりに仮定される場合は(簡潔性のために、1または3以外の倍率変化は可能ではないと仮定する)、

$P_0$  = 遺伝子がクラス0中にある事前確率、

$P_1$  = 遺伝子がクラス1中にある事前確率である。

【0053】

次いでベイズの定理を使用して以下の事後確率を得ることができる。

$P(p | 0)$  = 有意な倍率変化を割り当てられた遺伝子が実際には変化しなかった確率、

$P(\alpha | 0)$  = 変化なしのカテゴリーに割り当てられた遺伝子が実際には変化した確率。

その結果は以下の式である。

$$P(0 | p) = P_0 P(p | 0) / P_p \quad (47)$$

10

20

30

40

50

$$P(1) = P_1 P(1) / P_0 \quad (48)$$

上式で  $P_0$  および  $P_1$  はそれぞれ、遺伝子がクラス 1 またはクラス 0 中にあるということが出来る事後確率の合計であり、これらは以下の式で与えられる。

【 0 0 5 4 】

$$P_p = P(p=0)P_0 + (1-P(p=1))P_1 \quad (49)$$

$$P = (1-P(p=0))P_0 + P(p=1)P_1 \quad (50)$$

最大の関心が持たれる 2 つの量には以下のものがある。

$$P(1|1) = \text{絶対的な偽陰性率} \quad (51)$$

$$P(0|p) = \text{相対的な偽陽性率} \quad (52)$$

等式(51)および(52)の定義は対称的なものではない。なぜなら、 $P(0|p)$ の計算には前述の  $P_1$  の値が必要であるが、 $P(1|1)$ の場合はそうではないからである。絶対的な偽陰性率は、発現されている遺伝子すべての断片の指標であり、これらの遺伝子は所与の厳密さの検出スキームによって失われるはずである。一方で相対的な偽陽性率とは、誤って分類されている検出される遺伝子の断片であり、実際これらは変化をしなかった。したがって、本明細書で定義される偽陰性率は、このように検出スキームの効率(たとえばCowan, G., Statistical Data Analysis, p.47 (Claredon Press, Oxford, 1998)を参照のこと)または感度の指標であり(つまり、より小さな値が好ましい)、一方で偽陽性率はスキームの純度(たとえばCowan, G., Statistical Data Analysis, p.47 (Claredon Press, Oxford, 1998)を参照のこと)または選択性の指標である(つまり、より小さな値が好ましい)。

【 0 0 5 5 】

## 7.2 実験的シミュレーションの結果

図 6 中には(パネル A および B)、クラス 0 (変化なしのクラス)からの 1000 個の遺伝子、およびクラス 1 (3つの倍率の変化を示した遺伝子)からの 1000 個の遺伝子それぞれによって発生した(R、P)平面内の分散したプロットの比較がある。この概念の理解を促進するために、以下の等式は遺伝子の変化の事前確率に関するシナリオに基づく。

$$P_i = 0.2 \quad (53)$$

4000 個の変化していない遺伝子のバックグラウンドに対する、3つの倍率が変化している 1000 個の遺伝子の形状をこれによって定義する。

【 0 0 5 6 】

エラー率の決定表面の位置への依存を示すための選択の方法は、決定スキームのいわゆるレシーバ動作特性(ROC)のグラフ形式の構築である。たとえばVan Trees, H.L., Detection, Estimation and Modulation Theory. Part I (John Wiley and Sons), New York, 1998)を参照のこと。ROCによって、偽陽性率を低下、たとえば最小限にすること(これによってスキームの選択性が增大する)偽陰性率を低下、たとえば最小限にすること(これによっても所与のスキームの全体の選択性が增大する)の間の関係をはっきりと視覚化すること可能になる。

【 0 0 5 7 】

図 7 では、得られるレシーバ動作特性(ROC)をプロットし、このとき倍率変化自体が  $t = T$  である決定境界用の統計値として使用される。したがって前述の例では、受容領域は単に以下のように定義される。

【 数 2 7 】

$$D = \{ \hat{R} \geq R_c \} \quad (54)$$

したがって決定表面は(R、P)平面内の縦軸である。図 7 では、偽陽性率  $P(0|p)$  および偽陰性率  $P(1|1)$  を  $R_c$  の関数としてプロットする。受容に関する試験の厳密性が高くなると(つまり  $R_c$  が增大すると)、相対的な偽陽性率は低下するが(つまり検出されるサンプルの純度が增大するが)、これに付随して絶対的な偽陰性率も増大すること

が分かる。さらに図8では、統計値  $t = -P$  に関して得られるレシーバ動作特性 (ROC) を示す。

【0058】

図5は、セクション7に記載されるモンテカルロの実験的シミュレーションに基づいて、変化のない ( $b = 1$ ) 4000個の遺伝子のバックグラウンドに対する実際の倍率変化が  $b = 3$  である、1000個の遺伝子の検出に使用される3つの統計的方法のパフォーマンスを表形式で示す (つまり、3つの倍率が変化している遺伝子の事前確率は  $P_1 = 0.2$  である)。固定された偽陽性率  $P(0 | p) = 0.3$  に関する詳細な結果を報告する。図9は、RおよびP統計値について、変化のあった ( $b = 3$ ) 遺伝子の断片  $P_1$  の関数として、感受性を比較したことを示す。固定された相対的な偽陽性率  $P(0 | p) = 0.3$  を全体を通して課す。

10

【0059】

図5および図7に示す結果は、倍率変化 R の代わりに、またはこれと共に P - 値 (統計値  $t = -t$ ) を使用することによって、さまざまな範囲のパラメータの感度が大幅に増大し、以前はこれを検出するのは非常に困難であった。

結論としては、式のデータに固有なシグナルとノイズの比が低いために、このノイズは慎重に考慮しなければならない。本明細書に記載される P F O L D アルゴリズムによって、ノイズを処理するための理論的および実際的なフレームワークが与えられる。たとえば P F O L D アルゴリズムによって、遺伝子発現レベルの倍率の変化に関する2つの重要な計量値を与える。(i) 比の全体的な「質」を反映する P - 値、および (ii) 遺伝子 (群) の発現における倍率の変化の「量」を反映する R。さらに、P F O L D の p - 統計値は、このような遺伝子の発現レベルの小群の変化する遺伝子および / または小さな倍率の変化を定量化するために不可欠である。

20

【0060】

付録 A :  $X_1 = 0$  に対する分布

所与の倍率変化  $R - 1$  に関して、 $x_1 = 0$  である強度の対から最も有意性の低い予測値が生じる。この従属関係を調べるために、 $x_1 = 0$ 、および  $x_2 > >_{1,2}$  であるとき、等式(16)から近似値を発生させる。この結果は以下の通りである。

【数28】

$$f_R(R) \approx \left(\frac{2}{\pi}\right)^{1/2} \frac{Ry}{(R^2 + \sigma_2^2 / \sigma_1^2)^{3/2}} \exp\left(-\frac{y^2}{2(R^2 + \sigma_2^2 / \sigma_1^2)}\right) \quad (55)$$

30

上式では  $y = x_2 / x_1$  である。以下の変換を伴なう。

【数29】

$$u = \frac{1}{(R^2 + \sigma_2^2 / \sigma_1^2)^{1/2}} \quad (56)$$

40

【0061】

区間  $0 < R < \infty$  を  $0 < u < 1 / x_2$  に写像し、uに関する分布は以下の式によって与えられる。

【数30】

$$f_u(u) = \left(\frac{2}{\pi}\right)^{1/2} y \exp(-y^2 u^2 / 2) \quad (57)$$

分布関数の正しい正規化が得られることを確かめるのは容易である。なぜなら以下の通

50

りだからである。

【数 3 1】

$$\int_0^{\infty} f_R(R) dR = \int_0^{\sigma_1/\sigma_2} f_u(u) du \approx \int_0^{\infty} f_u(u) du = 1 \quad (58)$$

【 0 0 6 2 】

等式(58)においては近似が保たれる。なぜなら  $y = x_2 / x_1 > 1$  だからである。以下の式によって  $y$  の所与の値について  $P$  - 値を計算する。

【数 3 2】

$$P = P(R \leq 1) = P(u \geq u_1) \quad (59)$$

上式で

【数 3 3】

$$u_1 = \left(1 + \sigma_2^2 / \sigma_1^2\right)^{1/2} \quad (60)$$

である。

【 0 0 6 3 】

等式(57)を使用し、 $u_1 < u < \infty$  の範囲で積分することによって以下のことが分かる。

【数 3 4】

$$P = \operatorname{erfc} \left( \frac{y}{2^{1/2}} \left(1 + \sigma_2^2 / \sigma_1^2\right)^{-1/2} \right) \quad (61)$$

$y$  の所与の値について予測される倍率変化  $R$  を計算するために、分布のメジアンとしての  $R$  の定義を使用する。

【数 3 5】

$$\frac{1}{2} = P(R \leq \hat{R}) = P(u \geq) \quad (62)$$

これは結果として以下の等式になり、

【数 3 6】

$$\frac{1}{2} = \operatorname{erf} \left( \frac{y}{2^{1/2}} \left(\hat{R}^2 + \sigma_2^2 / \sigma_1^2\right)^{-1/2} \right) \quad (63)$$

上式から以下のような  $y$  と  $R$  の間の関係が分かる。

【 0 0 6 4 】

【数 3 7】

$$y = 2^{1/2} t_M \left(\hat{R}^2 + \sigma_2^2 / \sigma_1^2\right)^{1/2} \quad (64)$$

上式で  $t_M = 0.477$  は等式  $\operatorname{erf}(t_M) = 1/2$  の根である。等式(64)を使用して等

10

20

30

40

50

式(61)から  $y$  を除去することができ、以下の最終的な等式を得る。

【数 3 8】

$$P_n(\hat{R}) = \begin{cases} \operatorname{erfc}\left(t_m \frac{(\hat{R}^2 + \sigma_2^2 / \sigma_1^2)^{1/2}}{(1 + \sigma_2^2 / \sigma_1^2)^{1/2}}\right), & \hat{R} \geq 1 \\ \operatorname{erfc}\left(t_m \frac{(1 / \hat{R}^2 + \sigma_1^2 / \sigma_2^2)^{1/2}}{(1 + \sigma_1^2 / \sigma_2^2)^{1/2}}\right), & \hat{R} < 1 \end{cases} \quad (65)$$

10

上式で  $t_m = 0.477$  である。

【0065】

図11を参照すると、遺伝子群の少なくとも1つのアレイ中の遺伝子発現のレベルの違いを計算するための、プロセス150が示されている。プロセス150では数学的な式を使用して、少なくとも1つのアレイ中の1つまたは複数の細胞または組織タイプにおける遺伝子発現のレベルの得られる測定値から導くことができる、さまざまなレベルの（たとえばすべてのレベル）遺伝子発現の事後統計分布を記載する。

【0066】

ステージ152では、少なくとも1つのアレイ中の遺伝子それぞれのハイブリダイゼーション・シグナルと関係のある実験ノイズを定義する。このノイズは実験ノイズであり、チップまたは他のマイクロアレイ上での観察レベルにおける変数であり、生物学的システムにおいて見られる発現レベルの変数である生物学的ノイズではない。

20

【0067】

ステージ154では、定義した実験ノイズを使用して、遺伝子それぞれの強度の分布値を示す解析確率分布関数（pdf）を定義する。このノイズはガウスのものであると仮定し、ベイズの定理を使用して解析pdfを定義する。この解析pdfは連続関数である。

ステージ156では、解析pdfを使用して、少なくとも1つのアレイ中で差次的に発現される遺伝子または遺伝子産物について、予想比または倍率変化を示す解析同時pdfを導く。差次的に発現される任意の遺伝子または遺伝子産物について、予想比または倍率変化を示す解析同時pdfを導くことができる。

30

【0068】

ステージ158では、少なくとも1つのアレイ上の遺伝子から実験的に導かれた強度およびノイズ値を使用することによって同時pdfを適用し、遺伝子に関する倍率の変化に関係のある値を決定する。遺伝子転写物の濃縮物中の倍率の変化を推定する。特定の信頼区間で与えられる倍率の変化の信頼限界を確立する。倍率の変化を推定に関係のあるp値、または品質メトリックも導く。この値は、推定値が1より大きいときは倍率変化が1未満である確率、または推定値が1未満であるときは倍率変化が1より大きい確率を表す。

推定倍率の変化は、観察される遺伝子発現のレベルにおける違いを表す。転写物の濃度がゼロに向かうときでも、分散の合計（ノイズ）は依然として高い可能性がある。

40

【0069】

均等物

本発明の具体的な実施形態の前述の詳細な記載から、遺伝子発現の変化の推定値についての数学的に導かれる独特な事後分布が記載されている。個々の実施形態は本明細書において詳細に開示されているが、これはたとえば例示のみの目的で行われており、以下に続く添付の特許請求の範囲の範囲に関して制限することを目的とするものではない。詳細には、前述の記載および添付の図面からさまざまな置換形態、変更形態、修正形態が当業者に明らかになり、特許請求の範囲によって定義されるように本発明の精神および範囲から逸脱することなく、これらを実験に対して作成することができる。このような修正形態は添付の特許請求の範囲内に入ると企図される。

50

## 【 0 0 7 0 】

ソフトウェアの性質により、ソフトウェア、ハードウェア、ファームウェア、配線、またはこれらのいずれかの組み合わせを使用して、前述の機能を実施することもできる。これらの機能を実施するアイテムは、異なる物理的位置において機能の一部が実施されるように分布されていることを含めて、さまざまな位置に物理的に位置していてもよい。

## 【 図面の簡単な説明 】

【 図 1 】 倍率変化の事後分布、一連の測定値対  $(x_1, x_2)$  に関する等式(16)を示す線グラフである。 $(0, 0)$  以外はいずれの場合も、測定値の比は 4 である。両方のノイズ項の標準偏差は、 $\sigma_1 = \sigma_2 = 20$  の一定値に保つ。

【 図 2 】 等式(16)の導関数を定性的に例示し、図 1 に示される分布の状態を示す図である。シグナル  $(x_1, x_2)$  の対ではそれぞれ、平面状の 1 ポイント周辺で  $\pm \sigma_c$  の程度で 1 つのボックスが描かれる。ボックスのポイントに原点から線が引かれる。これらの線の傾きの分布は、倍率変化の事後分布である。パネル A はシグナル  $(100, 400)$  の構築を示す。パネル B はシグナル  $(5, 20)$  の構築を示す。

【 図 3 】 等式(30)および(33)によって導かれる、 $(R, P)$  平面への強度  $(x_1, x_2)$  のグラフ図である。一定の  $x_1$  および  $x_2$  の線が示される。黒い線は P の範囲の上限  $P_y(R)$  である。

【 図 4 】 図 1 の測定値の対すべてに関する、倍率の変化の推定値 R、68% の信頼区間  $(R_p, R_{p-1})$ 、および有意なプラスの倍率変化に関する P - 値の結果を表形式で示す図である。

【 図 5 】 本明細書のセクション 6 に記載されているモンテカルロのシミュレーションに基づいて、変化のない  $(b = 1)$  4000 個の遺伝子のバックグラウンドに対して実際の倍率変化が  $b = 3$  である、1000 個の遺伝子の検出に使用される 3 つの統計的方法のパフォーマンスの結果を表形式で示す図である（つまり、3 つの倍率が変化している遺伝子の事後確率は  $P_1 = 0.2$  である）。固定された偽陽性率  $P(0 | p) = 0.3$  に関する詳細な結果を報告する。

【 図 6 】 (1) クラス 0、変化のないクラス（パネル A）からの 1000 個の遺伝子、および (2) クラス 1（パネル B）からの 1000 個の遺伝子によって発生した  $(R, P)$  平面における散乱したプロットのばらつきの図である。

【 図 7 】 統計値  $t = P$  に関するレシーバー動作特性 (ROC) を表す線グラフである。

【 図 8 】 統計値  $t = -P$  に関するレシーバー動作特性 (ROC) を表す線グラフである。

【 図 9 】 R および P 統計値について、変化のあった  $(b = 3)$  遺伝子の断片  $P_1$  の関数として、感受性を比較したことを示す線グラフである。固定された相対的な偽陽性率  $P(0 | p) = 0.3$  が全体を通して課される。

【 図 10 】 本発明の特徴を実施するためのコンピュータ・システムのブロック図である。

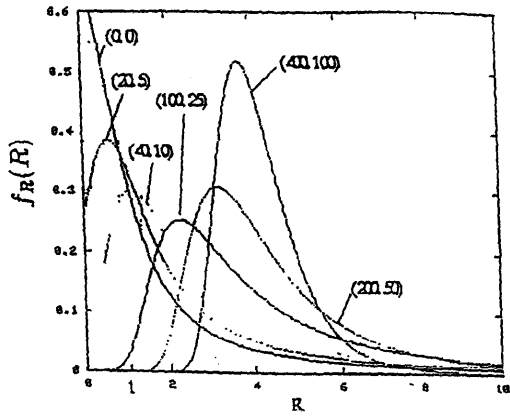
【 図 11 】 本発明のプロセスの流れ図である。

10

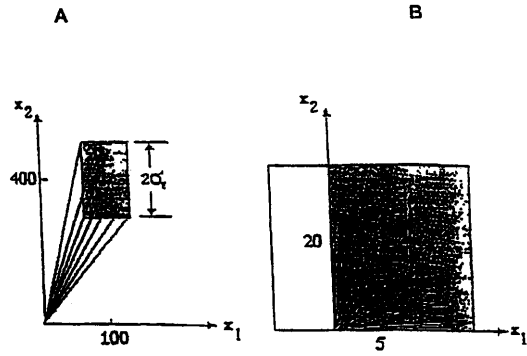
20

30

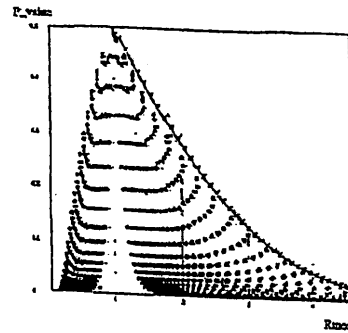
【 図 1 】



【 図 2 】



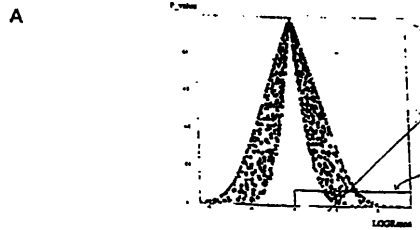
【 図 3 】



【 図 4 】

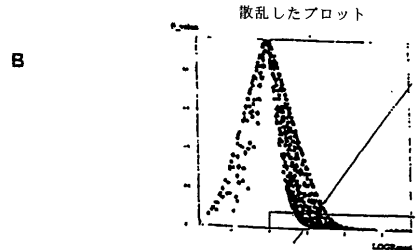
$x_1$	$x_2$	$R_p$	$R$	$R_{1p}$	P-値
100	400	3.3	4.0	5.0	0.0
50	200	2.8	3.9	6.5	$5.98 \times 10^{-1}$
25	100	2.0	3.6	8.8	$4.41 \times 10^{-1}$
10	40	0.93	2.2	7.3	0.18
5	20	0.51	1.5	5.7	0.34
1	4	0.32	1.1	4.6	0.46
0	0	0.23	1.0	4.4	0.5

【 図 6 】

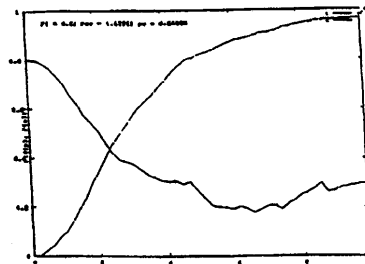


【 図 5 】

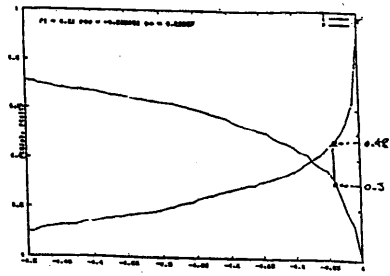
決定パラメータ		一定のP.P.率P(0 p)=0.3に関する結果				
メインの統計値t	領域	$N_{ror}$	$N_{rr}$	T.P.率 P(p 1)	$t_c$	Med(R)
フォールド変化R	$R \geq t_c$	371	260	0.96	4.37	6.1
P-値 P	$P < t_c$	814	570	0.57	0.037	4.2
判別式 $\log(R) - 5.48P$	$t \geq t_c$	713	499	0.50	0.85	4.4



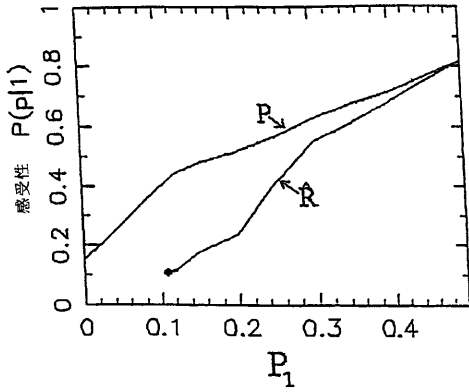
【 図 7 】



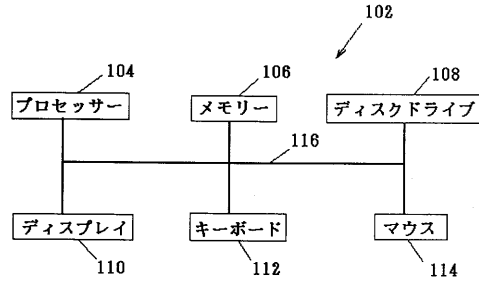
【図 8】



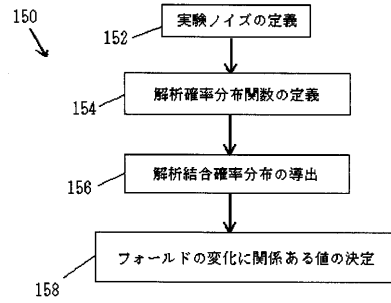
【図 9】



【図 10】



【図 11】



---

フロントページの続き

- (72)発明者 ジョウアキム・サイルハーバー  
アメリカ合衆国マサチューセッツ州02138・ケムブリッジ・サクラメントプレイス10
- (72)発明者 スティーヴン・ブシュネル  
アメリカ合衆国マサチューセッツ州02052・メドフィールド・サウスストリート41
- (72)発明者 ライナー・フックス  
アメリカ合衆国マサチューセッツ州02193・サドベリー・パウカードライブ40

審査官 宮久保 博幸

- (56)参考文献 特表平09-500199(JP,A)  
国際公開第97/048331(WO,A1)  
国際公開第99/054724(WO,A1)  
特表平09-507027(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 19/00

G06F 17/18