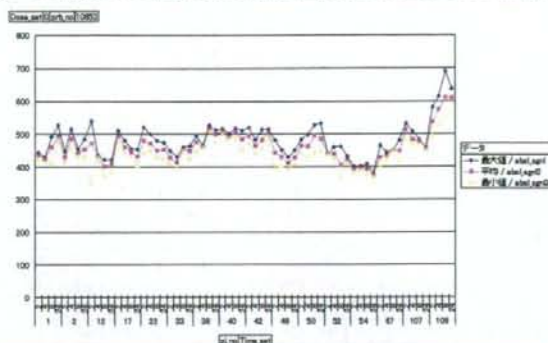


化合物間クラスタリングの課題

- 化合物間(Dose0同士)の距離を決定付けている遺伝子群の中からGc (Prb_No=10853)をグラフ化してみた。発現量が大きく、飽和していると考えられる



このようなプローブセットは排除する必要がある

まとめ(1)

- 課題
 - 遺伝子チップのプローブが飽和に近い状況になっていることが考えられる。飽和に近い状況に関して、絶対値変換の際の対策が必要である。
- 本クラスタリングの可能性
 - 今回は、階層型クラスタリングを行った。前段階として、全化合物同士の距離を計算している。これを用いることにより、MADICなどの距離を使用したクラスタリング手法を用いることも可能である。
 - また、距離計算の前段階として各遺伝子間の距離を計算している。これを用いることにより、化合物間の距離に対する各遺伝子の寄与を計算することができ、遺伝子に順位をつけることができる。逆に遺伝子リストを与えることで、飽和等の影響を排除した距離を計算することもできる。

化合物間クラスタリング(山脈構造)

◆同期率計算の山脈構造を用いる。ただし、計算量が多い

各遺伝子で

化合物aのクラスタリング結果を山脈化

化合物bのクラスタリング結果を山脈化



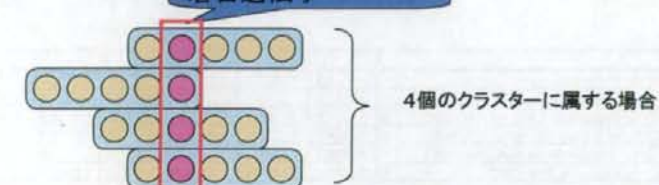
課題

1. クラスタリング対象外の取扱
2. ユニークとなった遺伝子の影響度
3. 計算量

同期率計算の山脈作成方法

着目する遺伝子に関わるクラスターを全部

着目遺伝子



属する個数の割合から、
山脈が描ける

この山脈の形状が似ていれば、近くの遺伝子の存在状況が似ているといえる

クラスタリング対象実験

prj_no	prj_name	chemical	sfc_no	sim_no	clst_no	山脈定義用データ数
12	TTG030-L	N-Methylaniline	2	1	1	99,957,189
17	TTG040-L	cisplatin	1	2	1	49,920,803
33	TTG044-L	Clofibrate	1	1	1	6,386,957
42	TTG055-L	N-ethyl-N-nitrosourea	1	1	1	9,802,013
50	TTG052-L	all trans retinoic acid	1	1	1	6,418,489
67	TTG061-L	Paraquat 2nd	1	1	1	54,690,275

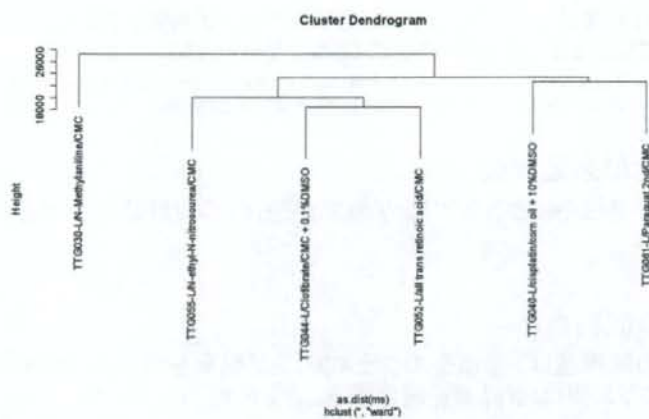
クラスタリング結果

- 距離マトリクス

	1	2	3	4	5	6	
TTG030-L/N-Methylaniline/CMC	1	0.0	25037.6	24231.9	24014.0	24928.7	25780.9
TTG040-L/cisplatin/corn oil + 10%DMSO	2	25037.6	0.0	20701.1	22379.2	22281.1	22690.0
TTG044-L/Clofibrate/CMC + 0.1%DMSO	3	24231.9	20701.1	0.0	18535.3	19138.9	20669.2
TTG052-L/all trans retinoic acid/CMC	4	24014.0	22379.2	18535.3	0.0	20175.7	22762.7
TTG055-L/N-ethyl-N-nitrosourea/CMC	5	24928.7	22281.1	19138.9	20175.7	0.0	22461.4
TTG061-L/Paraquat 2nd/CMC	6	25780.9	22690.0	20669.2	22762.7	22461.4	0.0

クラスタリング結果

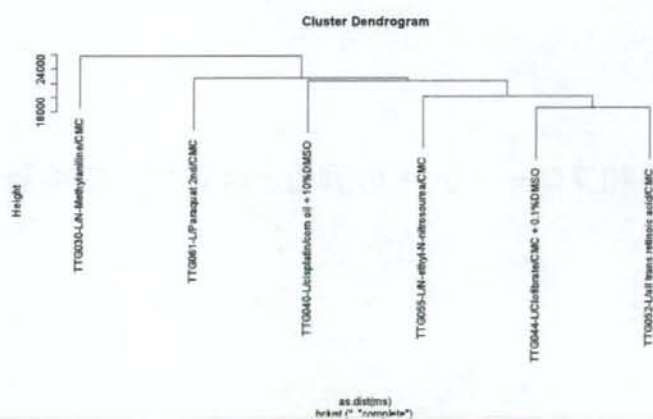
• Ward法



「Clofibrate」と「all trans retinoic acid」が近く、「N-Methylaniline」が遠いという結果になった

クラスタリング結果

• Complete法



「Clofibrate」と「all trans retinoic acid」が近く、「N-Methylaniline」が遠いという結果になった

まとめ(2)

- 優れた点
 - 山脈構造クラスタリングは、クラスタリング結果を用いるため、プロジェクト(化合物)ごとの偏差を受けにくい
 - 用量を含めたクラスタリング結果を用いるので、単純な化合物の影響を表現してクラスターに分割されるわけではないので結果の判断に留意が必要
- 改良が必要な点
 - 各遺伝子にクラスタリング結果が紐付くので計算量が大きい
- その他利点
 - 山脈構造は、重複型のクラスタリング結果を用いる。複数のクラスタリング結果を統合することも可能

5. 飽和プローブセットの調査と対処方法の検討

背景と目的

• 背景

- 平成18年度(STEP5)のテーマとして、化合物間クラスタリングの検討を行った。この中で、高発現のプローブセットで飽和により、正常に計測されていないと考えられるプローブセットが見つかった。これらのプローブセットを特定し、適切な対処を行わないと計測結果を用いた研究にて誤った結論を導く可能性がある。

• 目的

- 飽和していると考えられるプローブセットの特定とその対処方法を検討すること

仮説

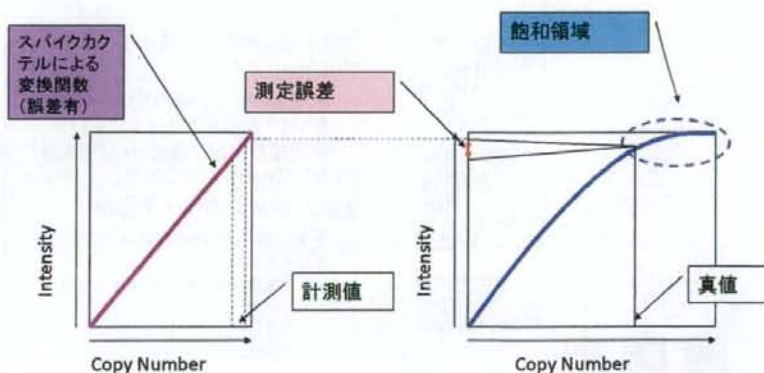


- 仮説1: 計測誤差は個体差に比べて小さい
 - 発現量が小さい領域では成り立たない(一般に認知されている)
 - 発現量が大きい領域でも成り立たない(Step5における発見)
- 仮説2: 観測誤差は不偏誤差
 - 不偏であるかのチェックがなされていない。
 - 発現量が大きい領域では、不偏ではない

現象

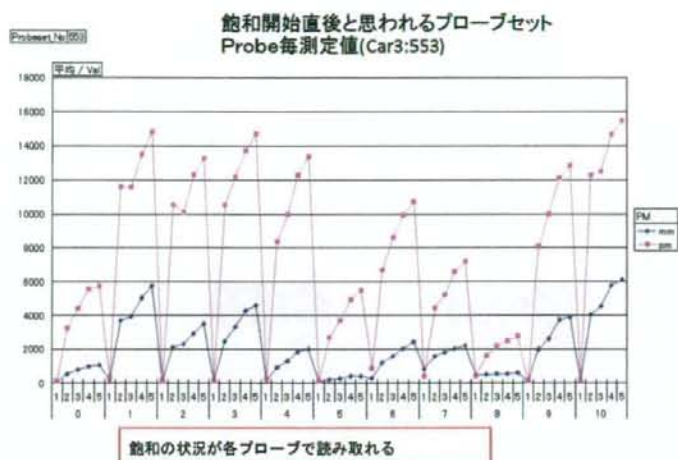
- 高発現域において正常に計測されていないと考えられるプローブセットが見つかった
- 飽和の可能性
 - 完全に飽和しているのではなく、飽和への漸近状態
 - 不偏の条件を満たさない
 - 特に飽和近辺では不偏の条件が明らかに崩れている
- 個体差と計測誤差を別個に取り扱うことが可能か？
 - N=3ならば、正規分布とみなすことが可能
 - 計測誤差が個体差に比べて小さければ、計測誤差が0であるとみなしても影響は少ない
 - 既存の統計手法で分析可能なはず
- 計測誤差が個体差よりも大きく、不偏の条件を満たさないならば、特別の取扱が必要である
 - 低発現域では、計測誤差(チップ偏差)は個体差よりも大きい。
 - 既知の課題として研究対象になっている
 - 中発現域では、計測誤差が不偏の条件を満たすと仮定できる？
 - 既存の分析手法で可能と考えられる
 - 高発現域は、飽和しつつある領域で、計測誤差が大きく、かつ、不偏の条件を満たしていない
 - Percellomeにより量を計測するためには補正が必要と考えられる。

Percellome変換の基礎

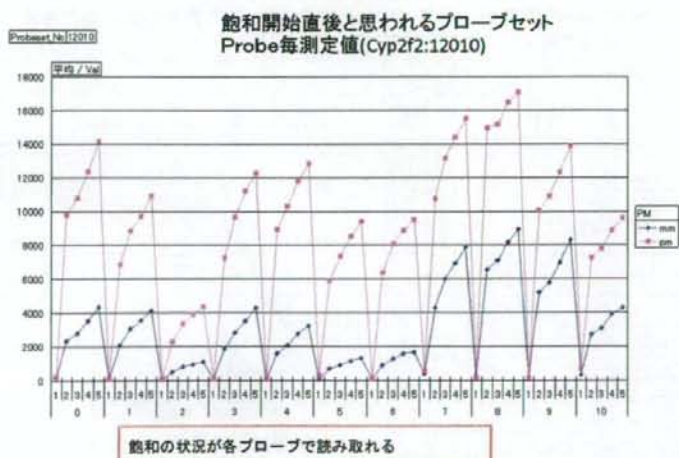


- 初期のデータでは、メーカーSOPをそのまま採用していたため、スパイク自身が飽和していた可能性がある
- 飽和領域はプローブごとorプローブセットごとの検討を行う

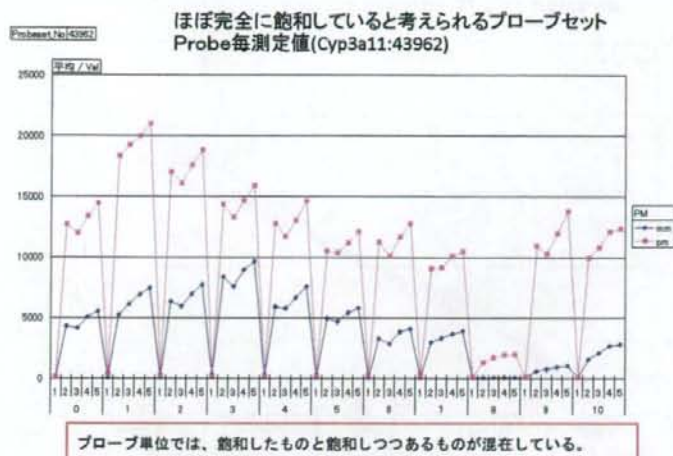
データ状況把握(飽和プローブセット)



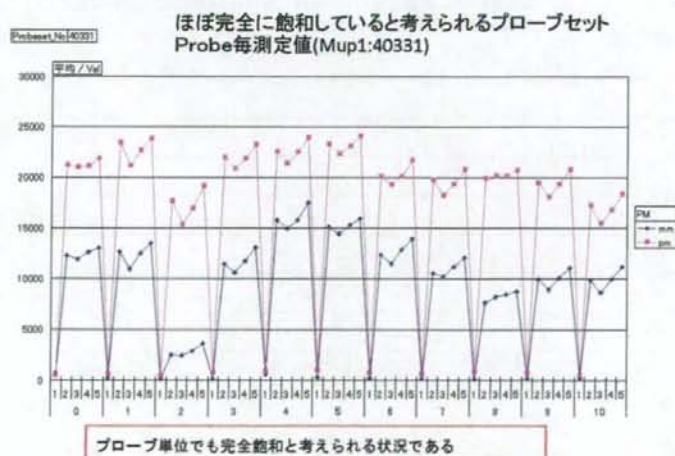
データ状況把握(飽和プローブセット)



データ状況把握(飽和プローブセット)

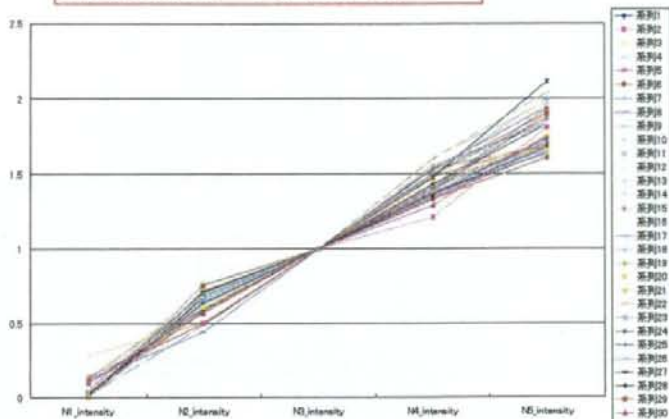


データ状況把握(飽和プローブセット)



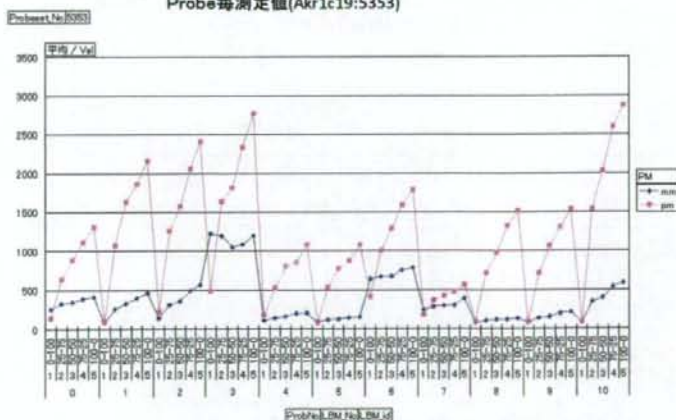
データ状況把握(線形的増加)

線形性が認められるプロブセットを見つけ出す



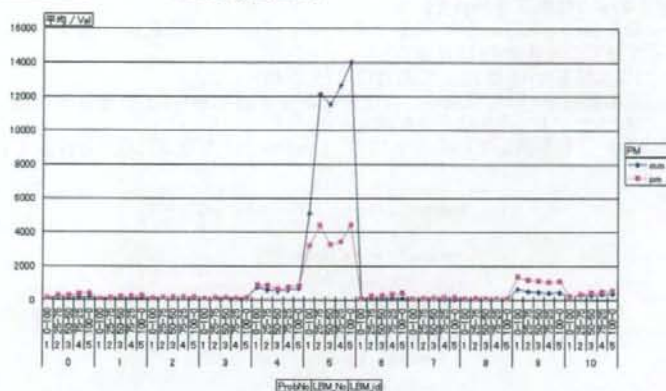
データ状況把握(線形的増加)

最もきれいな直線を描くProbesetでも、Probe単位ではブレがある。
Probe毎測定値(Akr1c19:5353)



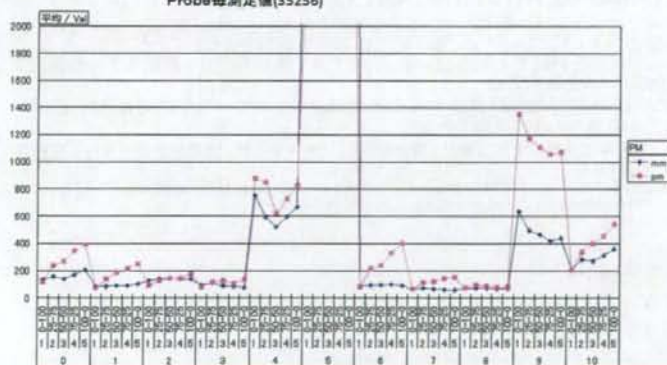
データ状況把握(線形的増加)

MMの方が大きいProbeが存在する。他の遺伝子に合致した可能性がある。
Probe毎測定値(35256)



データ状況把握(線形的増加)

幾つかのProbeで、直線性が見られる
Probe毎測定値(35256)



MMの大きなプローブは無視され、小さい値の線形性がみられるプローブから、線形性の強い遺伝子として検出された。

基本的アイデア①:

Langmuir isotherm equation

- Irving Langmuirによって1918年に導出された理論的な吸着等温式である。以下のような仮定を持っている。
 - 吸着媒には有限な数Nの吸着サイトがあり、そこだけで吸着質分子と結合する。
 - すべての吸着サイトは等価である。
 - 1つの吸着サイトは1つの吸着質分子としか結合しない。
 - 空の吸着サイトM、気相中の吸着質S、吸着サイトに結合した吸着質M-Sの間に $M + S \rightleftharpoons M-S$ の化学平衡が成立する。
- 各プローブが飽和している場合には、Langmuirの方程式に従うと仮定する

$$\log(P_{M_{ps}}) = \log\left(I_p \frac{k_p c_j}{1 + k_p c_j} + bg_p\right) + \varepsilon_{ps}$$

$p = 1, \dots, P$: probe_index

j : concentration_index

c : concentration

$l = 1, \dots, L$: replicate_index

I : Saturation_Intensity

k : equilibrium_constant

bg : background_component

基本的アイデア②

AIC(Akaike Information Criteria)

- すべてのプローブが飽和しているとは考えられない。また、飽和していたとしても、一定の値を示しているため、飽和しているかを判別できないことが考えられる。これらの状況を判断するため、Langmuirモデル、線形モデル、定数モデルを選択するためにAICを使用する
- AICとは、元統計数理研究所所長の赤池弘次によって1971年に考案された統計モデルの良さを示す指標である
 - モデルは複雑にすれば細かく適合させることができるが、偶然の影響を受ける可能性も高くなる
 - Akaike, Hirotugu (1974). "A new look at the statistical model identification". IEEE Transactions on Automatic Control 19 (6): 716-723.

$$AIC = -2 \log(L) + 2k$$

L : 尤度

k : パラメータ数

$$\text{確率密度関数: } p_{r,s}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\text{尤度関数: } L_r(\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$L_r(\mu, \sigma) = \prod_{i=1}^n L_r(\mu, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

$$\log L_r(\mu, \sigma) = \log\left(\prod_{i=1}^n L_r(\mu, \sigma)\right) = \log\left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}\right)$$

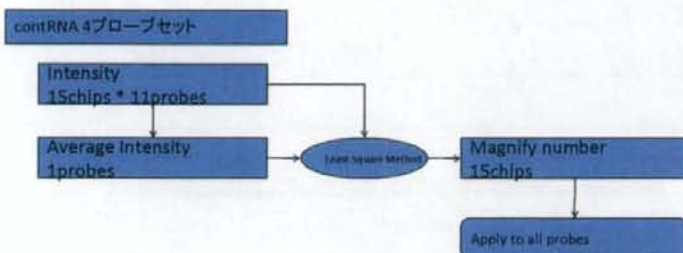
$$= \sum_{i=1}^n \log\left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}\right)$$

$$= \sum_{i=1}^n \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) + \sum_{i=1}^n \log\left(e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}\right)$$

$$= -n \log \sqrt{2\pi\sigma^2} - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

チップごとの標準化 (contRNA)

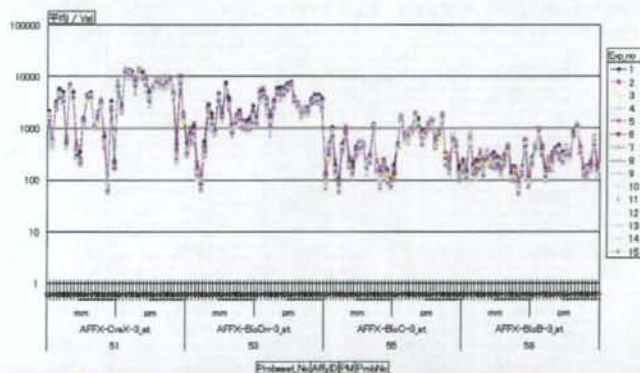
- チップごとの標準化のためにControl用RNAのデータ(contRNA)を使用する
 - 各プローブに対する平均値を求め、それらの平均値に近づくように(最小二乗法)にチップごとの変換倍率を求め、その値をチップの全プローブにかける。
 - 対数領域で計算する
 - ContRNAはハイブリ液に一定量を加えたE.Coli由来のmRNAの測定データを示す



チップごとの標準化 (contRNA)

チップごとの標準化 (contRNA)

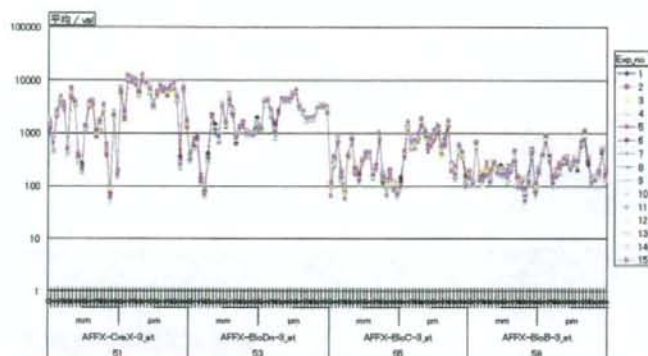
- 露光等の問題により、チップ毎にバイアス値が異なる。このような測定誤差の変動を補正し、標準化を行う。
- contRNAのプローブセットを使用する。



プローブ毎に変動しているが、チップ(Exp. no)毎に似たパターンを示している。これらの平均値への誤差が最小になるように、定数倍することで標準化を行う。

チップごとの標準化 (contRNA)

チップごとの標準化 (contRNA) の実施後のグラフ



わずかであるが、線の重なりが増え、チップ間のばらつきが小さくなったことが読み取れる

Langmuir式フィッティング検証

Langmuir式にフィッティング可能かいくつかのプロブセットで試してみる
実験対象プロブセット

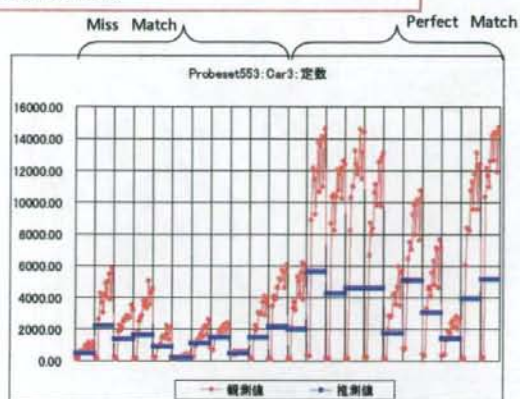
・ 計算を試みるプロブセット

- 飽和、線形と思われたプロブセット
- LBM, 50%-50%における、発現値順位、最上位、10%、50%、80%
- Percellome Spike Cocktail (内2プロブセットは、プロブ数が異なるので今回は省略)

	No.	PNNo	Avg	SD	AFIXID	GeneSymbol	CV	
飽和	511	46	992	170020	63810	460202.at	Cerb	26
	523	86	12010	862640	37024	1349732.at	Cyrb6	26
	513	87	43867	238426	37833	141883.at	Cyrb21	26
	508	11	43031	888620	37844	142095.at	Mafk	26
線形	678	5548	5703	77776	7776	1400404.at	Mafk128	26
	2613	13120	26242	8056	889	1420540.at	Sfr	111
発現値	558	11	43031	1107610	37224	141395.at	Ameo2	23
	558	2	43015	1090020	43838	141395.at	Fabp1	26
	578	3	29590	1190177	44044	1424110.at	Mafk	26
	558	4	4602	100017	3818	140952.at	Ameo2	26
	578	5	41510	1013204	37414	141910.at	Ameo2	26
中央値	1028	4632	46068	45025	4821	1415708.at	Taf1	106
	1028	4614	4617	45016	1013	1420920.at	111000P1.69a	12
	1028	4615	46054	4612	1575	141193.at	Taf1	106
	1028	4616	30282	46350	4873	1420948.at	Zfp6	106
発現値	1028	4617	4622	46250	477	1424189.at	Napst1	11
	508	25951	149914	4954	476	1423814.at	Wdr53	103
	508	25952	4846	4848	1188	1409403.at	-	708
	508	25953	49844	4942	1188	1413172.at	Mafk	208
	508	25954	14755	4947	1041	143400.at	14241130	211
	508	25955	19742	4941	1184	1441769.at	Ubp1	208
発現値	8028	30282	30483	1348	369	1422303.at	Taf1b2	278
	8028	30283	30282	1347	371	1430011.at	Sox10	208
	8028	30284	4142	1347	371	1422444.at	14241130	211
	8028	30285	30240	1341	371	1400928.at	Srsf11	578
	8028	30287	14493	1341	1184	1424189.at	-	208
Percellome Spike	158	30281	30	30031	444	1457472.at	-	24
	688	5965	30	30032	474	1457472.at	-	24
	248	1092	37	189242	7783	1457472.at	-	24
	578	508	71	492176	4621	1457472.at	-	24
	578	181	71	892821	38922	1457472.at	-	24

Langmuir式フィッティング検証 飽和遺伝子

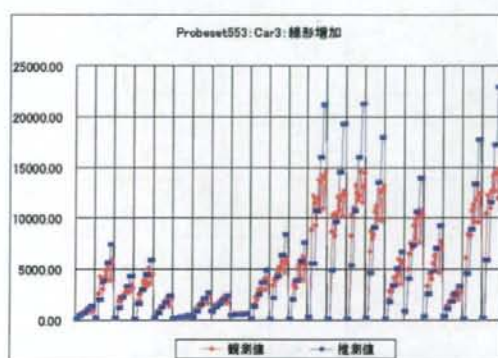
飽和遺伝子: 553: Car3: 定数



AIC=2568.4

Langmuir式フィッティング検証 飽和遺伝子

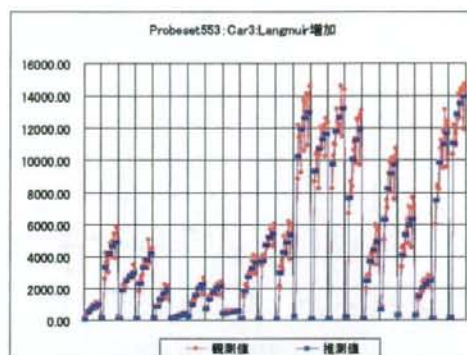
飽和遺伝子: 553: Car3: 線形(増加)



AIC=927.6

Langmuir式フィッティング検証 飽和遺伝子

飽和遺伝子: 553: Car3: Langmuir (増加)

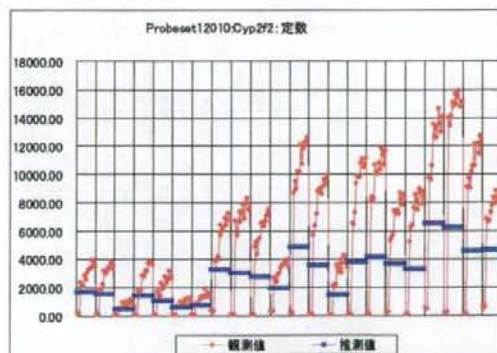


最もAICの小さい最適モデル

AIC=714.3

Langmuir式フィッティング検証 飽和遺伝子

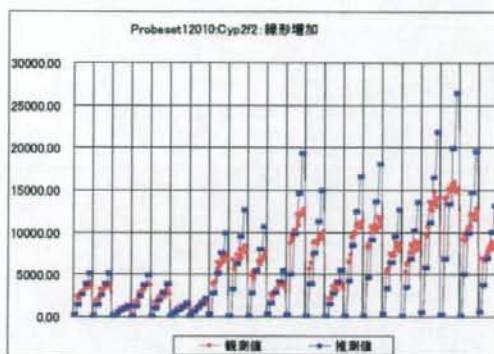
飽和遺伝子: 12010: Cyp2f2: 定数



AIC=2768.1

Langmuir式フィッティング検証 飽和遺伝子

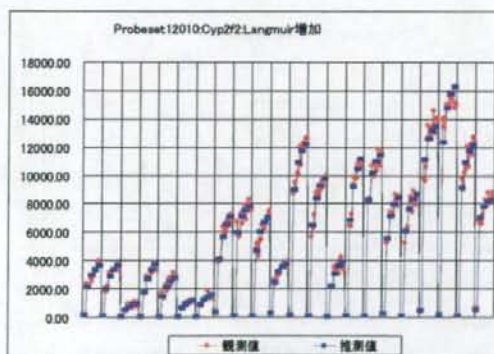
飽和遺伝子:12010:Cyp2f2:線形(増加)



AIC=930.2

Langmuir式フィッティング検証 飽和遺伝子

飽和遺伝子:12010:Cyp2f2:Langmuir(増加)

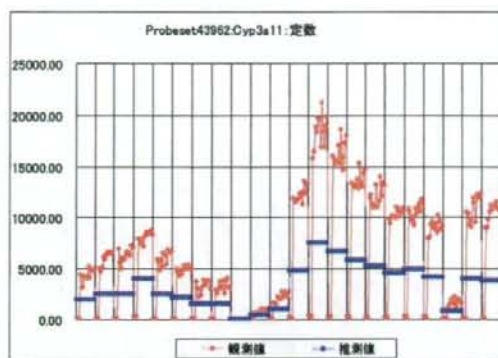


最もAICの小さい最適モデル

AIC=590.1

Langmuir式フィッティング検証 飽和遺伝子

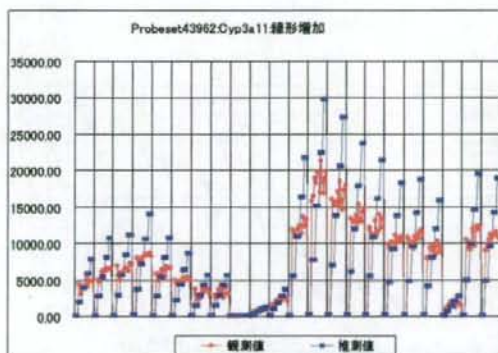
飽和遺伝子: 43962: Cyp3a11: 定数



AIC=2995.2

Langmuir式フィッティング検証 飽和遺伝子

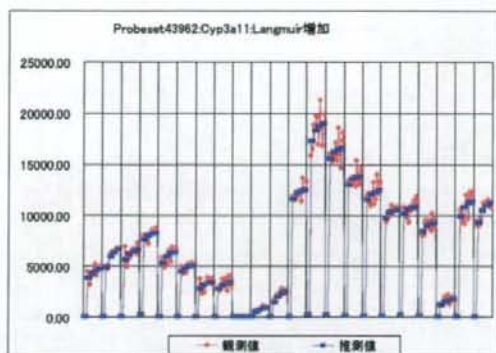
飽和遺伝子: 43962: Cyp3a11: 線形(増加)



AIC=1076.8

Langmuir式フィッティング検証 飽和遺伝子

飽和遺伝子: 43962:Cyp3a11:Langmuir(増加)

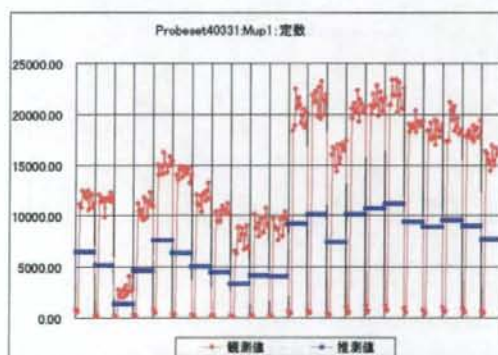


最もAICの小さい最適モデル

AIC=622.9

Langmuir式フィッティング検証 飽和遺伝子

飽和遺伝子: 40331:Mup1:定数



AIC=2865.0