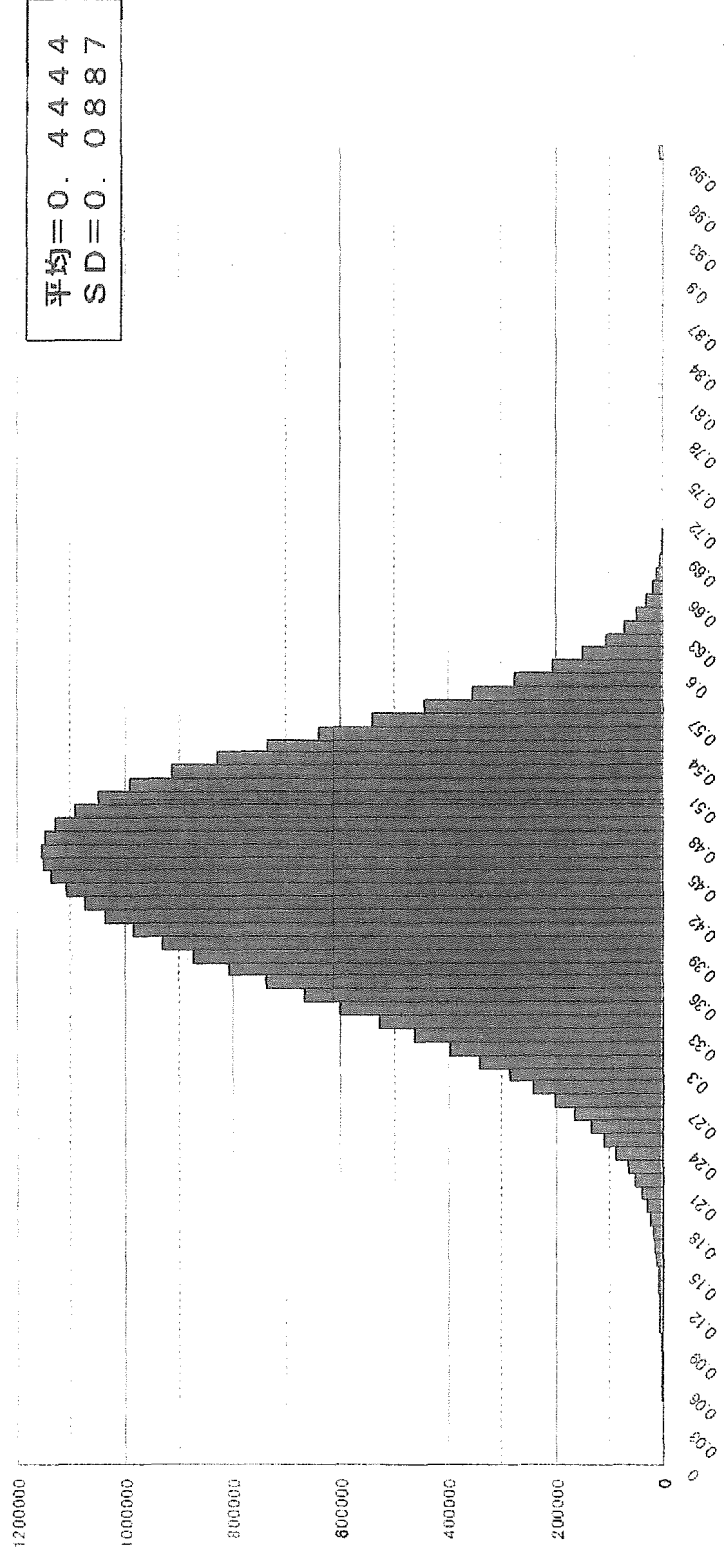


3. TMF計算

(1) TMF計算結果の分布

貴研究所、ご提示の遺伝子間の類似度計算方法 (TMF) を全組み合わせについて、計算しその分布を示す。



正規分布のような釣鐘型の形状となっている。TMFの特性上、このような形状が保障されるわけではないが、概ねこのような形状になると類推される。

3. TMF計算結果報告

(2) 性能考察

様々な改良を施した結果、以下のような時間となった

(1) TTG9, TTG10 5, 132Gene : 3時間30分
 レコード件数=5132 x 5132 ÷ 2 約1317万

(2) UTERO 6, 976Gene : 7時間30分
 レコード件数=6976 x 6976 ÷ 2 約2433万

(3) UTERO全件 12, 488Gene : 15時間30分
 レコード件数=12488 x 12488 ÷ 2 約7798万

(3) TTG9全件 22, 690Gene : 58時間30分
 レコード件数=22690 x 22690 ÷ 2 約2億5742万

ただし、Teradataは別PJと共同で使用しているため、クリアな状態ではない

4. 密度ベースと階層的クラスタリングの融合



(1) Zスコアを用いた階層的クラスタリング

Zスコアを用いた、全Geneによる階層的クラスタリングについて取り組んだ

(1) 前提条件

全Gene間の類似度を表すために、何らかの類似度（または非類似度）を定義できたと仮定する

(2) 前処理

- ・全Geneについて、全dose/時間を対象にZスコアを計算する
- ・上で計算した16個のZスコアを、16次元空間の座標とみなす
- ・全Geneについて、ユークリッド距離で類似度を計算する（※）

※2つのGene、A、Bを与えたときに、対称性があることから $A \geq B$ という制約をつける

・クラスタリング手法

閾値ごとに、DBSCANクラスタリング手法（改良版）を実施し、その結果を用いて、階層的な表現とする

・問題点

- ・閾値を決定するのが難しい
- ・ノイズの影響により、クラスタリングの結果に大きく影響を与える可能性がある

4. 密度ベースと階層的クラスタリングの融合

(2) T M F 計算結果を用いた階層的クラスタリング

T M F 計算結果を用いた、階層的クラスタリングを改良し、5 1 2 9 G e n e で実施した

(1) 前提条件

Zスコアと同じ

(2) 処理概要

Zスコアと同じ

(3) クラスタリング結果 (閾値=0.72 n=4)

1 4 4 コア G e n e : 4 6 クラスタ

※問題点

Zスコアと同様に、閾値とnの値を決定するのが難しい

4. 密度ベースと階層的クラスタリングの融合

(2) T MF 計算結果を用いた階層的クラスタリング (クラスタ結果)

T MF 計算結果を用いた、階層的クラスタリングを改良し、5129 Gene で実施した

id	clst_label	gene_core
1	1415682_at	1415682_at
		1415689_at, 1420524_at, 1433552_at
		1415702_at, 1421847_at, 1434646_s.at
		1415813_at, 1421893_at, 1438046_s.at
3	1415867_at	1450275_x.at
		1450735_at, 1437503_s.at
		1415946_at, 1423416_at, 1437850_x.at
		1415961_at, 1423912_at, 1433192_s.at
		1416252_at, 1423038_at, 1448665_at
		1416835_at, 1424311_at, 1445495_at
		1416742_at, 1424451_at, 1449644_at
2	1415699_s.at	1424488_at, 1448846_s.at
		1416764_at, 1425129_s.at, 1450740_at
		1417157_at, 1426111_x.at, 1451357_at
		1417390_at, 1426233_at, 1451362_at
		1417544_s.at, 1426783_at, 1451782_s.at
		1417744_s.at, 1427078_at, 1452092_s.at
		1417785_at, 1427078_at, 1452092_s.at
		1418112_at, 1427720_s.at, 1452147_at
		1418522_at, 1428373_at, 1455001_x.at
		1418794_at, 1429296_at, 1455855_x.at
		1419260_at, 1429534_at, 1456059_at
		1419445_s.at, 1429707_at, 1456312_x.at
		1419819_s.at, 1433451_at

id	clst_label	gene_core
7	1415986_at	1415986_at
8	1416215_at	1416215_at
9	1416372_at	1416372_at
		1416500_at
		1416667_at
		1417285_at
10	1416500_at	1424235_at
		1426162_s.at
		1433395_s.at
		1448697_s.at
11	1416595_at	1416595_at
12	1416671_at	1416671_at
13	1416789_at	1416789_at
14	1417508_at	1417508_at
15	1417934_at	1417934_at
16	1418631_at	1418631_at
17	1418897_at	1418897_at
		1418996_s.at
18	1418996_s.at	1428031_s.at
		1448823_at

id	clst_label	gene_core
19	1420460_s.at	1420460_s.at
20	1420525_s.at	1420525_s.at
21	1422487_at	1422487_at
		1452651_at
22	1422576_at	1422576_at
23	1423394_at	1423394_at
24	1423739_x.at	1423739_x.at
		1424041_s.at
25	1424008_s.at	1424008_s.at
		1424406_at
26	1424039_at	1424039_at
		1428138_s.at
		1452059_at
27	1424117_at	1424117_at
28	1424694_at	1424694_at
29	1424708_at	1424708_at
30	1426414_s.at	1426414_s.at
31	1426679_at	1426679_at
32	1427060_at	1427060_at
33	1427896_at	1427896_at

id	clst_label	gene_core
34	1428218_s.at	1428218_s.at
35	1431423_s.at	1431423_s.at
36	1431431_s.at	1431431_s.at
37	1434251_at	1434251_at
38	1435995_at	1435995_at
39	1438159_x.at	1438159_x.at
40	1448206_at	1448206_at
41	1448279_at	1448279_at
42	1448621_s.at	1448621_s.at
		1450431_s.at
43	1450720_at	1450720_at
		1455152_at
44	1452683_at	1452683_at
45	1452917_at	1452917_at
46	1454955_at	1454955_at

4. 密度ベースと階層的クラスタリングの融合

(2) T MF 計算結果を用いた階層的クラスタリング (クラスタ情報)

id	clst_label	min_gene	n	tmf_val
1	1415682_at	1415682_at	1	1.0000
2	1415699_a_at	1426233_at	62	0.5708
3	1415867_at	1416284_at	7	0.5917
4	1415870_at	1416153_at	9	0.6164
5	1415876_a_at	1415876_a_at	1	1.0001
6	1415979_x_at	1438477_a_at	5	0.7111
7	1415986_at	1415986_at	1	1.0000
8	1416215_at	1416215_at	1	1.0000
9	1416372_at	1416372_at	1	0.9999
10	1416500_at	1448697_s_at	7	0.6705
11	1416595_at	1416595_at	1	0.9999
12	1416671_a_at	1416671_a_at	1	0.9999
13	1416789_at	1417468_at	2	0.7280
14	1417508_at	1417508_at	1	1.0001
15	1417934_at	1417934_at	1	1.0000

id	clst_label	min_gene	n	tmf_val
16	1418631_at	1419177_at	2	0.7238
17	1418897_at	1418897_at	1	1.0001
18	1418996_a_at	1448823_at	3	0.6623
19	1420460_a_at	1420460_a_at	1	1.0000
20	1420525_a_at	1420525_a_at	1	1.0000
21	1422487_at	1452691_at	2	0.7474
22	1422576_at	1422576_at	1	1.0000
23	1423394_at	1423394_at	1	1.0001
24	1423739_x_at	1449710_s_at	3	0.7140
25	1424008_a_at	1424406_at	2	0.7278
26	1424039_at	1428138_s_at	3	0.6615
27	1424117_at	1424117_at	1	1.0001
28	1424694_at	1424694_at	1	1.0000
29	1424708_at	1424708_at	1	1.0000
30	1426414_a_at	1426414_a_at	1	1.0000

id	clst_label	min_gene	n	tmf_val
31	1426679_at	1426679_at	1	1.0001
32	1427060_at	1427060_at	1	0.9999
33	1427896_at	1451700_a_at	2	0.7241
34	1428218_a_at	1428218_a_at	1	1.0001
35	1431423_a_at	1431423_a_at	1	1.0000
36	1431431_a_at	1431431_a_at	1	1.0000
37	1434251_at	1434251_at	1	1.0000
38	1435995_at	1435995_at	1	1.0000
39	1438159_x_at	1438159_x_at	1	1.0001
40	1448206_at	1448206_at	1	1.0000
41	1448279_at	1448279_at	1	1.0001
42	1448621_a_at	1450431_a_at	2	0.7364
43	1450720_at	1455152_at	2	0.7400
44	1452683_at	1452683_at	1	0.9999
45	1452917_at	1452917_at	1	1.0000
46	1454955_at	1454955_at	1	1.0000

4. 密度ベースと階層的クラスタリングの融合

(3) d を動かしたときの階層的クラスタリングの違いについて

TMFクラスタリングに関して、 $d=3$ と $d=4$ について、簡単に比較してみました
 下記に $d=3$ と $d=4$ におけるクラスタ数と出現 Gene 数についてまとめています

	dens	0.700	0.705	0.710	0.715	0.720	0.725	0.730	0.735	0.740	0.745
N of Clst	3	220	212	197	171	154	128	104	83	62	45
N of Gene		3575	2723	2015	1453	1021	667	438	286	179	120
N of Clst	4	91	84	55	69	51	52	44	35	28	14
N of Gene		2348	1723	1186	814	518	330	200	126	78	48

	dens	0.750	0.755	0.760	0.765	0.770	0.775	0.780	0.785	0.790	0.795
N of Clst	3	30	18	9	7	4	3	2	2	2	2
N of Gene		72	52	37	32	25	23	22	22	22	21
N of Clst	4	7	3	2	2	2	2	2	2	2	2
N of Gene		33	28	25	22	22	22	20	20	20	20

	dens	0.800	0.805	0.810	0.815	0.820	0.825	0.830	0.835	0.840	0.845
N of Clst	3	2	2	2	2	1	1	1	1	1	1
N of Gene		20	20	18	17	14	14	14	14	14	13
N of Clst	4	2	2	2	1	1	1	1	1	1	1
N of Gene		19	17	17	14	14	14	14	14	13	12

	dens	0.850	0.855	0.860	0.865	0.870
N of Clst	3	1	1	1	1	1
N of Gene		13	12	12	12	10
N of Clst	4	1	1	1	1	1
N of Gene		12	12	11	10	9

5. 関数近似手法

(1) 前提条件

前提条件①

時間とドーズを次のように置き換える

時間	0h	2h	4h	8h	24h
T	0	1	2	3	4

投与量	0	0.1	1.0	10.0
S	0	1	2	3

前提条件② (近似関数の構成方法)

- (1) $t = 0$ においてドーズに依存せず、一定
- (2) 今回のデータの最も自由度の高い次数で表現する

$$\begin{aligned}
 f(t, s) = & \beta_{00} + \beta_{10}t + \beta_{20}t^2 + \beta_{30}t^3 + \beta_{40}t^4 \\
 & + \beta_{11}ts + \beta_{21}t^2s + \beta_{31}t^3s + \beta_{41}t^4s \\
 & + \beta_{12}ts^2 + \beta_{22}t^2s^2 + \beta_{32}t^3s^2 + \beta_{42}t^4s^2 \\
 & + \beta_{13}ts^3 + \beta_{23}t^2s^3 + \beta_{33}t^3s^3 + \beta_{43}t^4s^3
 \end{aligned}$$

$\beta_{01}, \beta_{02}, \beta_{03}$ は、条件(1)を満たすため、恒等的に0

5. 関数近似手法

(1) 前提条件

AICによる係数推定

遠伝子によっては、ノイズの部分が強く、高次方程式による近似が適切でない場合もある。これらの方程式の選択をAICを用いて行う。

$$AIC = -2 \cdot \text{最大対数尤度} + 2 \cdot \text{パラメータ数}$$

$$\text{最大対数尤度} = -\frac{n}{2} \log 2\pi - n \log \sigma - \frac{n}{2}$$

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

時間、投与量、各々、何次までと区切って、係数を求める。それぞれの方程式に対するAICを計算し、比較する。AICは、小さい値のものが良好な結果であると認める。絶対値ではなく、一連の比較対象の相対値として比較を行う必要がある。

5. 関数近似手法

(2) 関数近似手法を用いた分析

関数近似手法を行った後、相関関数を用いた時間軸の分析を実施しました

- (1) 関数近似実施
- (2) クラスタリング実施
関数近似結果を表現値に戻し、ZSCORE化を実施
単純な3パターンを排除（フラット、時間で1次増加、時間で1次減少）
- (3) クラスタ毎に中心変動を求め
- (4) クラスタ同士の相互相関係数関数を求める

関数近似を行うことにより、誤差を排除することはできた。しかし、その結果を用いてクラスタリングを行うと、再現性の低い遠伝子ほど先にクラスタとしてまとまるという結果となった。

適切ではない結果 になった。このアプローチは、クラスタリングの前処理として用いるべきではない と考えられる。

5. 関数近似手法

(3) 関数近似手法の改善

単純な多項式による関数近似だけでなく、生物学的な制約条件を組み込んだ多項式を作成し近似を行う。近似結果として改善されていると考える。

- ① 投与量=0の場合には、時間に依存せず発現量は一定である
- ② 投与量=0の場合には、0hと24hは同じ発現量になる
- ③ 全投与量において、時間が2h以下では、同じ発現量になる
- ④ 全投与量において、時間が4h以下では、同じ発現量になる
- ⑤ 全投与量において、時間が8h以下では、同じ発現量になる

生物学的制約は、手法自身は有効な手段と考えられる。しかし、新規の制約を組み込む際に数式への展開が必要のため、研究者による自由な拡張が制限されるため、分析手法として活用しにくいものと考えられる。

国立医薬品食品衛生研究所

MilleFeuille システム 運用説明書

2004年4月1日 STEP 1版

NTTコムウェア株式会社

目次

- 1 はじめに
 - 2 プロジェクト情報の作成
 - 2-1 ジョブの新規作成
 - 2-2 プロジェクト情報の登録
 - 3 実験データの登録
 - 4 EIFファイルの作成
 - 5 Surface情報の作成
 - 5-1 ジョブの新規作成
 - 5-2 EIFファイルの登録
 - 5-3 一次QCの実行
 - 5-4 実験データのロード、Surfaceデータ作成
 - 6 TMF計算
 - 6-1 ジョブの新規作成
 - 6-2 TMF計算の実行
 - 7 CIFファイルの作成
 - 8 クラスタリング情報の作成
 - 8-1 ジョブの作成
 - 8-2 CIFファイルの登録
 - 8-3 クラスタリングの実行
 - 9 その他
 - 9-1 バックアップとリストアについて
 - 9-2 故障発生時の連絡先
-
- 参考1 全ジョブネット一覧
 - 参考2 システムフロー
 - 参考3 ディレクトリ構成
 - 参考4 ログの確認方法

1 はじめに

NTT COMWARE

本書は、国立医薬品食品衛生研究所毒性部における
Millefeuille システム STEP 1 版
についての運用説明書です。

プロジェクトの登録からクラスタリングの実行までを、
実際の実験フローに従って記述しています。

2004年3月31日
NTTコムウェア株式会社

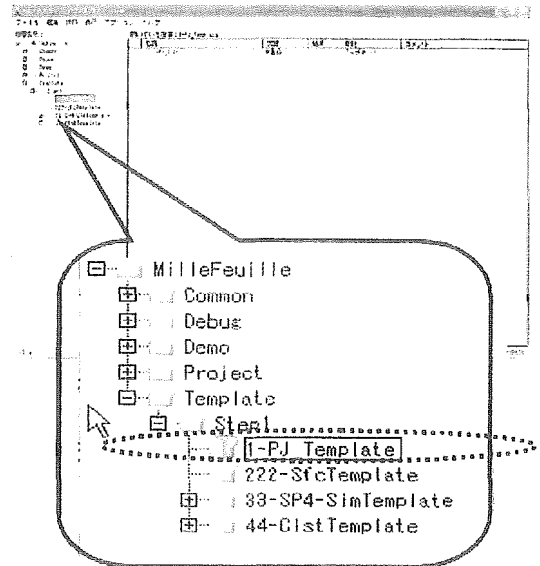
? プロジェクト情報の作成

プロジェクト情報をJP1上で登録し、初期環境とEISファイルの作成を行います。

2-1 JP1におけるジョブの新規作成

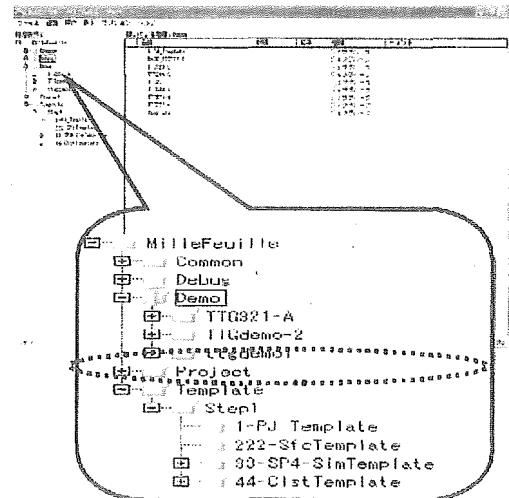
① プロジェクトテンプレートのコピー

JP1の左側の、
/MilleFeuille/Template/Step1/1-PJ_Template
を右クリックし、コピーを選択する。



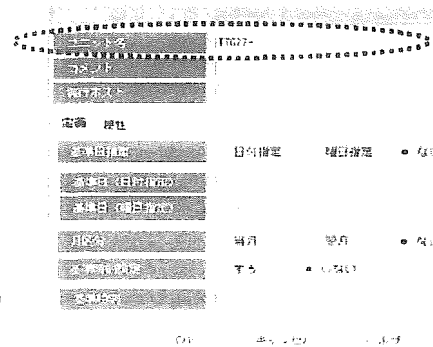
② プロジェクトテンプレートの貼り付け

JP1の左側より、
/MilleFeuille/Projectを右クリックし、
貼り付けを選択する。



③ プロジェクト名の入力

/MilleFeuille/Demoの下に
(1-PJ_Template)がコピーされるので、
それを右クリックし、プロパティを
選択する。すると、右のような画面が
表示されるので、ユニット名に②で
作成したジョブグループと同じPJ名
(例：TTG27-L)を入力後、OKボタンを押す。



2 プロジェクト情報の作成

2-2 プロジェクト情報の登録

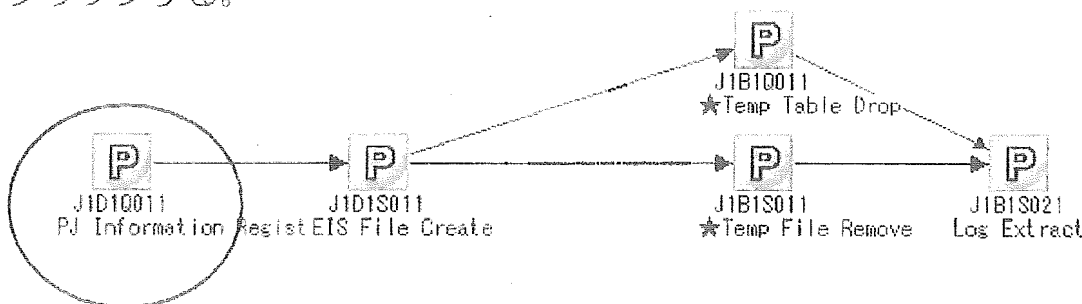
① ジョブネットの選択

JP 1の画面左側で、2-1で作成したジョブグループを（例：TTG27-L）を選択し、画面右側に表示されるジョブネット（PJ_Init）をダブルクリックする。

名	機	線	樹
PJ_Init	登録		海外

② ユニットの選択

①後、新たにウィンドウが立ち上がる（ジョブネットエディタ）ので、そのウィンドウに表示されているアイコン（JID1Q011）をダブルクリックする。



③ プロジェクト情報の入力

②後、右図のような画面が立ち上がる。

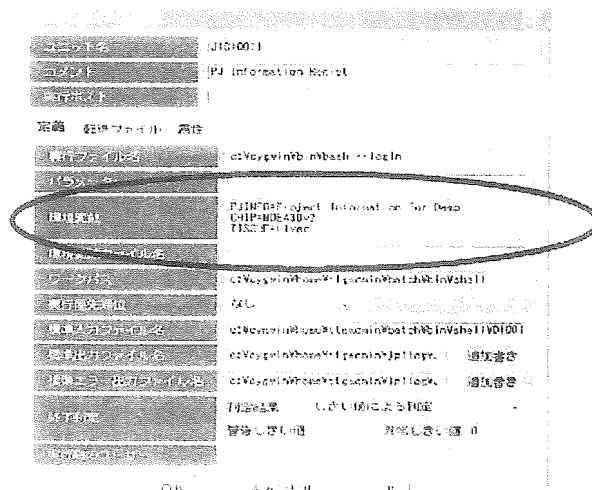
この環境変数エリアに、

- ・ PJINFO（プロジェクト情報）
- ・ CHIP（チップ名）
- ・ TISSUE（臓器情報）

を入力し、OKボタンを押す。

各入力情報は“=”の後に、記入して頂きますが、下記制約があります。

- ・ 各文字列の中で、改行は入れない
- ・ 半角英数字（大文字小文字区別あり）
- ・ PJINFO（255文字以内）
- ・ CHIP（255文字以内）
- ・ TISSUE（255文字以内）



2 プロジェクト情報の作成

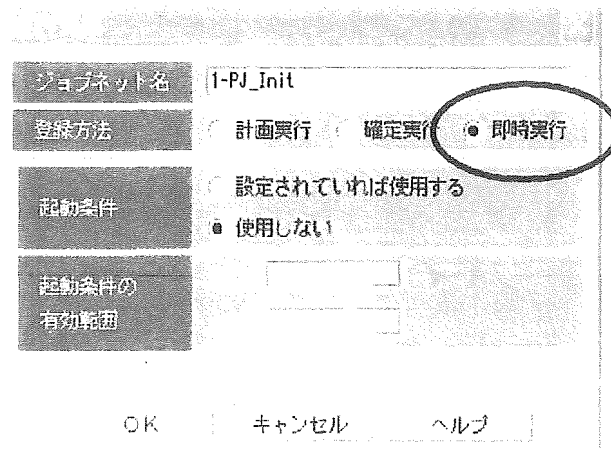
NTT COMWARE

④ 実行登録

③でプロジェクト情報を登録したプロジェクト（ジョブグループ）にあるPJ_Init（ジョブネット）を右クリックし、実行登録を選択する。

⑤ 即時実行

実行方法から即時実行を選択し、OKボタンを押す。



⑥ 実行結果の確認

ジョブが正常終了すると、

の下に、入力したプロジェクト名（例：TTG27-L）と同じ名前のフォルダが作成され、

の下には、EISファイルが作成される。

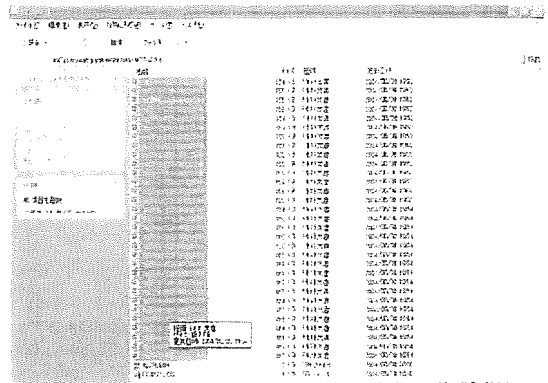
EISファイル名：プロジェクト名.EIS
（例：TTG27-L.EIS）

TTG27-L.EIS 1 KB EISファイル

実験データ (Raw Data) を下記手順に従って、MilleFeuilleシステムに登録して下さい。

① 実験データのコピー

の下
に2で入力したプロジェクト名
でフォルダが作成されているので、
この下に対象となる実験データを
コピーして下さい。



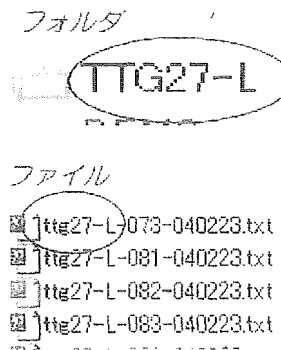
② ファイルの確認

コピーした実験データのファイル名の
先頭部分が、作成したフォルダの
プロジェクト名と等しいことを確認
して下さい。

【確認例】

フォルダ名：TTG27-L

ファイル名：ttg27-L-083-040223.txt



4 EIFファイルの作成

EIFファイルの作成方法を簡単に解説します

各クライアントより、EIC.exe を実行してください



EIC.exe

- ① 各項目に必要な事項を登録してください（画面（1））

画面（1）

#	Fix	Experiment Name	Time	Dose	Serial#	Date	Comment
1	<input type="checkbox"/>	TTG27-L-011-040220	02	000	011	04/02/20	
2	<input type="checkbox"/>	TTG27-L-012-040223	02	000	012	04/02/23	
3	<input type="checkbox"/>	TTG27-L-013-040223	02	000	013	04/02/23	
4	<input type="checkbox"/>	TTG27-L-021-040223	02	1	021	04/02/23	
5	<input type="checkbox"/>	TTG27-L-022-040223	02	1	022	04/02/23	

- ② 次に、ID,Comment など必要情報を登録してください（画面（2））

画面（2）

Surface ID	Surface Name	Normalization	P	Q	R
001	123-TTG27-L-001	Absolutization (Std)	1.1	1.5	

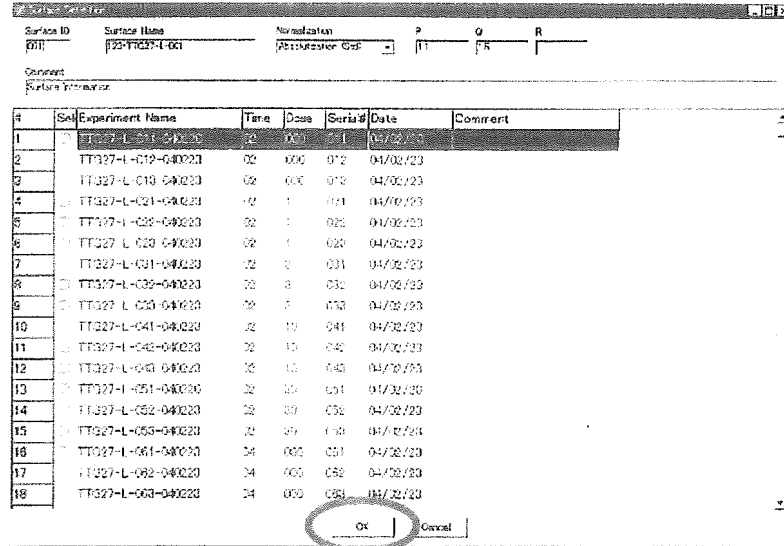
Comment
Surface Information

#	Sel	Experiment Name	Time	Dose	Serial#	Date	Comment
1	<input type="checkbox"/>	TTG27-L-011-040220	02	000	011	04/02/20	
2	<input type="checkbox"/>	TTG27-L-012-040223	02	000	012	04/02/23	
3	<input type="checkbox"/>	TTG27-L-013-040223	02	000	013	04/02/23	
4	<input type="checkbox"/>	TTG27-L-021-040223	02	1	021	04/02/23	
5	<input type="checkbox"/>	TTG27-L-022-040223	02	1	022	04/02/23	

4 EIFファイルの作成

③ 最後に、OKボタンを押せばEIFファイルが作成されます（画面（3））

画面（3）



（参考）EIFファイルサンプル

0123	TTG27-L	←←←←	PJ-No	PJ-Name		
	1,2,3-Triazole	←←←←	PJ-Information			
	MOE430v2	←←←←	Chip-Name			
	Liver	←←←←	Tissue-Name			
	001	←←←←	Surface-No			
	Test EIF	←←←←	Surface-Information			
1	1.1	1.5	←←←←	Absolutized-Mode	Param1	Param2
1	TTG27-L-011-040220			2	0	04/02/20
2	TTG27-L-012-040223			2	0	04/02/23
3	TTG27-L-013-040223			2	0	04/02/23
...						
58	TTG27-L-161-040224			24	30	04/02/23
59	TTG27-L-162-040224			24	30	04/02/23
60	TTG27-L-163-040224			24	30	04/02/23