

Users Guide for ExP to Perform EPIG

Extracting Patterns and Identifying co-expressed Genes



Version 1.2 (April 15th, 2011)

Analysis of Gene Expression Data using EPIG (Extracting Patterns and Identifying co-expressed Genes)

Launching ExP

Double click on the EPIG.bat file (dated 10/30/2009). The Expression Predictor dialog box will launch. See below.



Data Format

The gene expression data must be relative change in a fixed/standard tab-delimited text format:

- First row: unique array name
- Second row: names for replicate groups
- First col: probe ID or gene accession

See below:

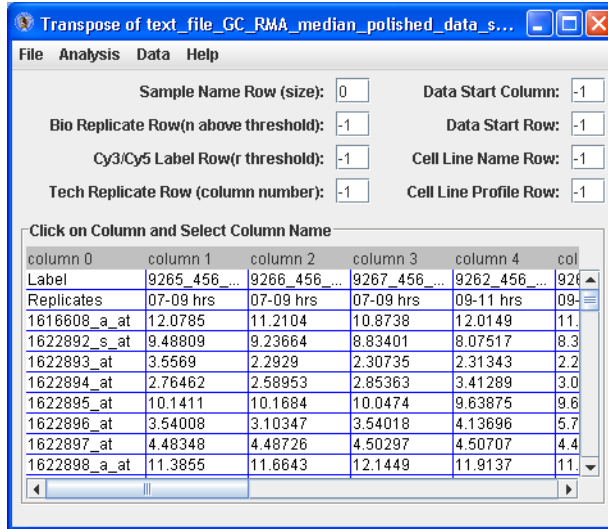
Probe	betweenG_1_1	betweenG_1_2	betweenG_2_1	betweenG_2_2	betweenG_3_1	betweenG_3_2
replicate	betweenG_1	betweenG_1	betweenG_2	betweenG_2	betweenG_3	betweenG_3
A_42_P534203	4.18	3.14	1.91	2.39	1.24	1.02
A_43_P22195	2.47	2.56	2.54	1.84	1.17	0.61

Convert Intensity Data to Relative Change

If the data is single channel intensity data (Affymetrix or Agilent for example), the data must be converted to relative change using intensity measurements from a baseline

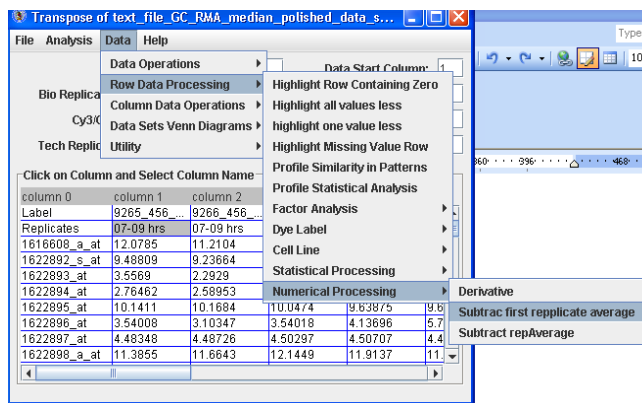
experiment, control sample or sham group. Otherwise, skip to the **Load Relative Change Data** section.

Click the Load Compiled Expression button. The data set will appear in a table. See below.



Highlight by clicking the first cell under “column 1”, right click and choose the “Select Data Start Column” option. The “Data Start Column” box value should switch from -1 to 1. Next, highlight by clicking the very first cell (upper most left) with a data value, right click and choose the “Select Data Start Row” option. The “Data Start Row” box value should switch from -1 to 2. Finally, highlight by clicking a cell in the row with the name of one of the replicate samples (i.e. 07-09 hrs in the figure above), right click and choose the “Replicate Row” option. The “Bio Replicate Row” box value should switch from -1 to 1.

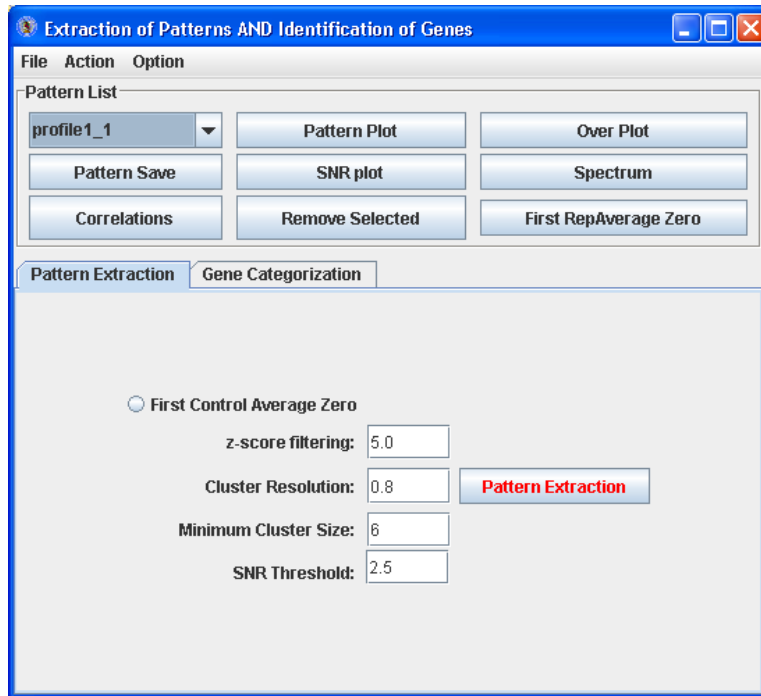
From the menu option select Data→Row Data Processing→Numerical Processing→Subtract First Replicate Average.



From the Menu option select File -> Save Data to save the converted data.

Load Relative Change Data

With the data already saved and it is relative change measurements, it can be loaded directly into ExP for EPIG analysis by clicking the EPIG button option on the ExP splash screen. Or, if the converted relative change data is in memory as above, from the Menu option select Analysis -> EPIG. The EPIG analysis dialog box will launch. See below.



Cell lines, multiple controls/references/shams adjustment or alignment

You have a case as below.

T1 ctl1 trt11 trt12 trt13...

T2 ctl2 trt21 trt22 trt23...

T3 ctl3 trt31 trt32 trt33...

At each time point, there is a time matched control. Your input file can be like this –

sampleid	S1	S2	S3	S4	S5	S6	...	S7	S8	S9	S10	S11	S12	...
time	T1	T1	T1	T1	T1	T1	...	T2	T2	T2	T2	T2	T2	...
trt	Ctl1	Ctl1	Ctl1	Trt11	Trt11	Trt11	...	Ctl2	Ctl2	Ctl2	Trt21	Trt21	Trt21	...
Gene1	Data start						...							
Gene1														

You can select (right mouse button click) time row as cell line name row, and trt row as cell line profile row. Also you select trt row as replicate row.

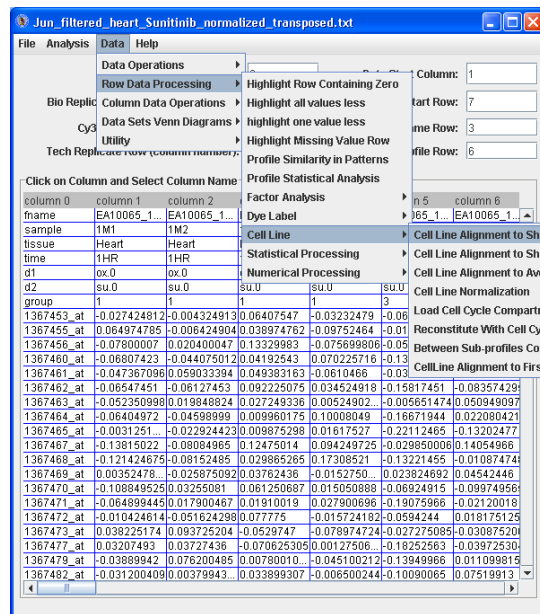
(1)

Then you go to menu Data>Row data processing>cell line>cell line alignment to sham zero. – this will use the ctl as a baseline (average to be zero and all others adjusted accordingly)within its own time point. Then you can save the data if you want and run EPIG.

(2)

Another approach is you go to menu Data>Row data processing>Numerical processing>subtract first replicate average – in this case the first control group used as a baseline, then run EPIG.

These 2 different approaches will give you different results. Both are interesting. Depending on your focus, you may use one or both results.

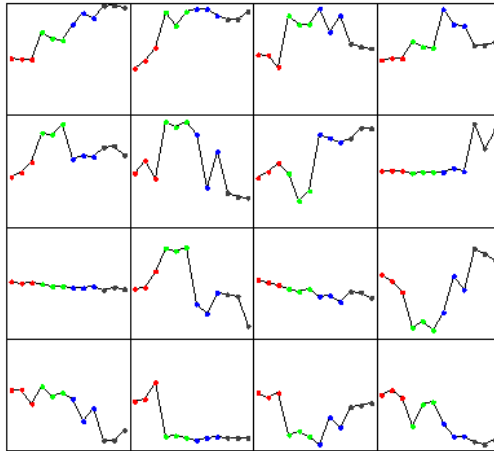


EPIG Analysis: Pattern Extraction

The Pattern Extract tab contains an option and fields for parameters that control the extraction of the significant patterns from the gene expression data:

- First control Average Zero – click this radio button if the first array data set is a “sham” control or baseline which all the other arrays are compared to
- z-score filtering - filters those probes/genes profiles which have one or more data points with z-score larger than the set value
- Cluster resolution - correlation threshold to filter those patterns generated with the r-value larger than that to another pattern.
- Minimum Cluster size – the number probe/gene profiles to be merged for generating the patterns
- SNR Threshold – the threshold for the signal-to-noise for the profiles to be filtered out and thus not used for generating the patterns

Click the “Pattern Extraction” button to generate the patterns. A plot of the extracted patterns will be displayed (see below).

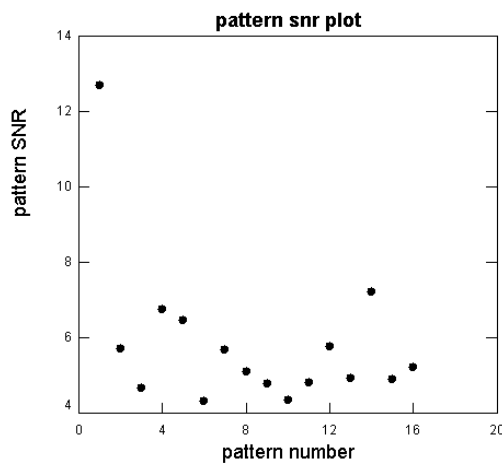


The patterns are arranged in order from the upper-left corner going right. Leave the default settings or go to the menu at the top, Action -> Reset to start over and modify the parameters according to get more or less patterns to be generated. However, from empirical analysis, the default parameters work well. Go to Action -> Rename Pattern Name to name the patterns numerically in the Pattern List drop down.

Some patterns will look very similar (i.e. 6 and 10, 9 and 11, 13 and 16). They will likely have a high correlation value. Click the “Correlation” button to see the pairwise correlation values in the Message Board dialog box. See below.

Pattern	SNR	Mag	SNR_pValue	ANOVA	ANOVA_pValue			
0.79288954	1.0	0.6950458	0.7307569	0.67089677	0.15500996	0.41037086	0.24798886	-0.7914782
0.3403344	0.6950458	1.0	0.7653934	0.5245207	0.69786763	-0.033109333	-0.3243777	-0.34570765
0.588281	0.7307569	0.7653934	1.0	0.116767496	0.18393864	0.5294783	-0.076751314	-0.58972293
0.32274702	0.67089677	0.5245207	0.116767496	1.0	0.48075536	-0.2858558	0.0446259	-0.36152065
-0.374996	0.15500996	0.69786763	0.18393864	0.48075536	1.0	-0.663471	-0.6837736	0.27259564
0.74349666	0.41037086	-0.033109333	0.5294783	-0.2858558	-0.663471	1.0	0.5735619	-0.6161712
0.65567696	0.24798886	-0.3243777	-0.076751314	0.0446259	-0.6837736	0.5735619	1.0	-0.6639985
-0.89916277	-0.7914782	-0.34570765	-0.58972293	-0.36152065	0.27259564	-0.6161712	-0.6639985	1.0
-0.5194785	-0.10762728	0.22067894	-0.38123578	0.5885514	0.74226606	-0.84852946	-0.5132226	0.47542062
-0.7984953	-0.82875437	-0.59602004	-0.8034659	-0.28331652	0.037248634	-0.5908825	-0.26019767	0.7889088
0.42425033	-0.17937848	-0.61513513	-0.16279803	-0.50017375	-0.95260346	0.6689673	0.7401924	-0.31240857
-0.8104283	-0.3974954	0.24086808	-0.16486362	-0.03869681	0.7881939	-0.76719177	-0.8393417	0.7424558
-0.8178918	-0.8866073	-0.74107856	-0.6686368	-0.630424	-0.17765951	-0.31700006	-0.3004827	0.78730756
-0.25138468	-0.7046596	-0.9653431	-0.68762827	-0.6462589	-0.7756052	0.14925231	0.36086887	0.32431388
-0.971262	-0.7455766	-0.30232435	-0.5790937	-0.2387012	0.41646773	-0.7996987	-0.8120091	0.79442596

Click the SNR plot button to see a plot of the SRN for each pattern. See below.

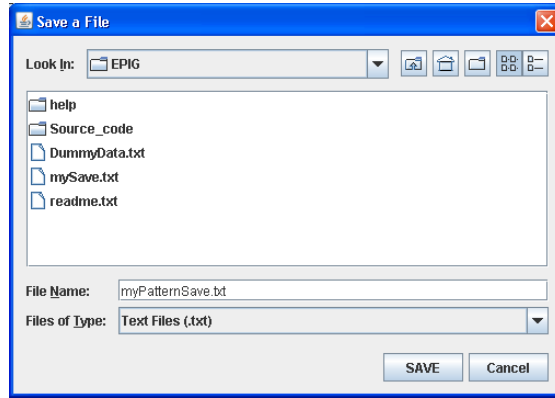


The “Message Board” dialog box displays the SNR values, magnitude of change, SNR p-value, ANOVA test statistic and ANOVA p-value for each pattern. See below.

Pattern	SNR	Mag	SNR_pValue	ANOVA	ANOVA_pValue
1	12.679711	4.143225	0.000000	0.080079	0.000001
2	5.714427	1.590073	0.000077	0.058070	0.000130
3	4.652817	1.532069	0.000340	0.081318	0.000429
4	6.740316	1.218798	0.000023	0.024522	0.000089
5	6.465257	0.893779	0.000031	0.014333	0.000117
6	4.318239	2.080281	0.000572	0.174057	0.001955
7	5.674582	0.875256	0.000082	0.017843	0.000120
8	5.098265	0.841472	0.000177	0.020431	0.000218
9	4.773723	0.068562	0.000284	0.000155	0.000736
10	4.337179	0.901183	0.000555	0.032380	0.000925
11	4.804509	0.166823	0.000271	0.000904	0.000835
12	5.753014	1.726544	0.000074	0.067550	0.000294
13	4.912333	0.889690	0.000232	0.024602	0.000423
14	7.220184	0.874827	0.000014	0.011011	0.000015
15	4.891305	1.293093	0.000239	0.052417	0.000309
16	5.207103	1.642817	0.000152	0.074653	0.000330

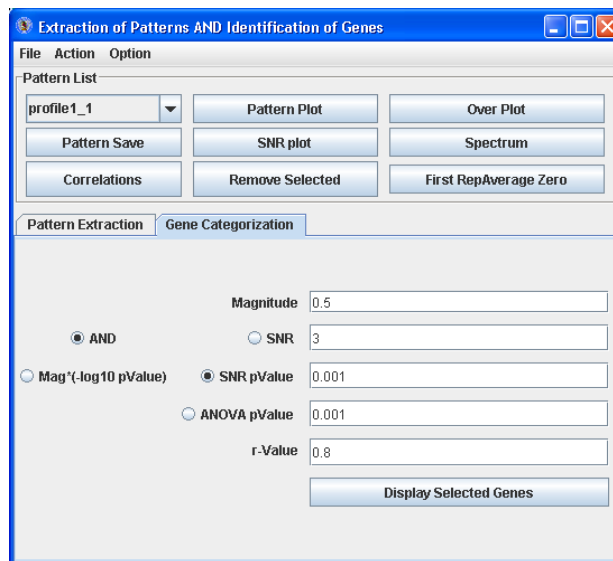
The pattern (of the two similar ones) with the lower SNR (of patterns 13 and 16, pattern 13 has a lower SNR) may be deleted by selecting the pattern from the drop down menu and click the “Remove Selected” button.

Once you have the number of desired patterns generated, click “Pattern Save” and enter in a name for a file to be generated which will have the patterns and associated parameter values stored in it. See below.



EPIG Analysis: Gene Categorization

Click the Gene Categorization tab (see below).



This tab contains parameters that control the binning of probes/genes to the patterns that were generated:

- Magnitude (S) – the amount of variation in gene expression within profiles.

$$S = \begin{cases} \max\{\bar{g}_i\}, & \text{if } \min\{\bar{g}_i\} > 0 \\ -\min\{\bar{g}_i\}, & \text{elseif } \max\{\bar{g}_i\} < 0 \\ \max\{\bar{g}_i\} - \min\{\bar{g}_i\} & \text{otherwise} \end{cases} .$$

where \bar{g}_i is the intra-group average.

- SNR – the signal to noise ratio for the profile

$$SNR = S/N$$

where S is denoted above,

$$N = \sqrt{\frac{\sum_i^m [(n_i - 1) s_i^2]}{\frac{m}{\sum_i (n_i - 1)}} \sum_i^m \frac{1}{n_i}}$$

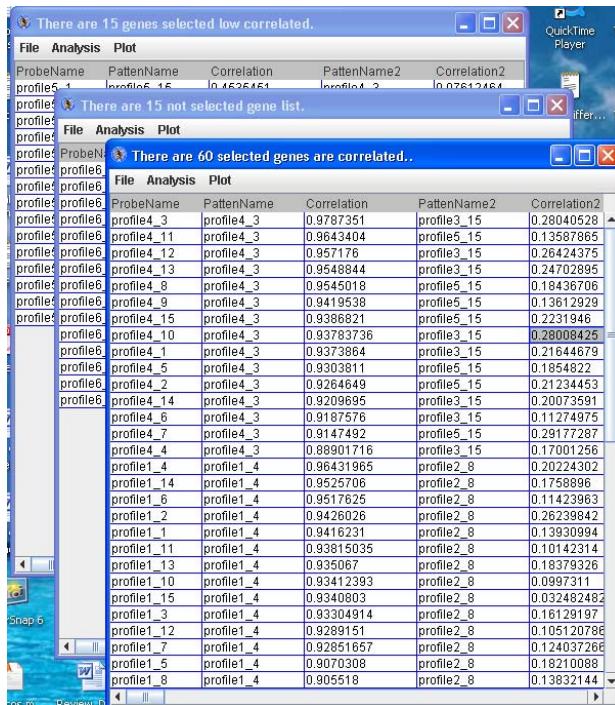
and the sample variance is

$$s_i^2 = \frac{\sum_j^{n_i} (g_{ij} - \bar{g}_i)^2}{n_i - 1}.$$

- SNR p-value – the significance level for the SNR
- ANOVA p-value – the significance level for an ANOVA of the profile
- r-value – the correlation value for the profile to be similar to the pattern

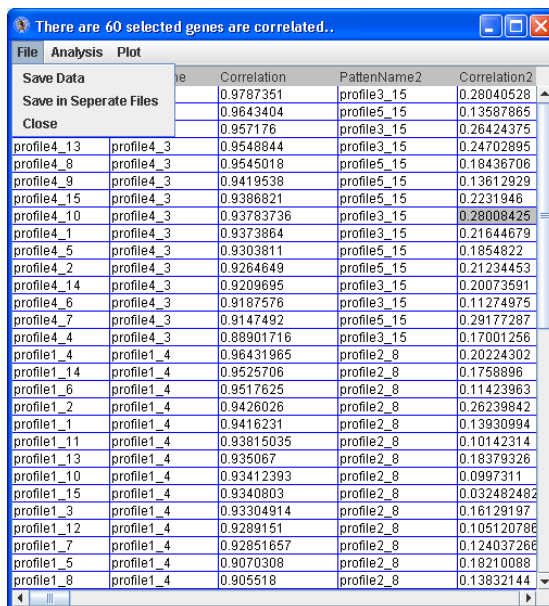
From empirical analysis, the default parameter settings work well. To adjust the settings, enter a value for Magnitude, click “And”, click the radio button and enter a value for either SNR, SNR p-value or ANOVA p-value, and then enter an r-value. Click Mag*(log10pvalue) if you want to use the John Zheng modification for categorization of the profiles to the patterns that takes into account both the magnitude of change and the p-value.

Click “Display Selected Genes” to launch up to three dialog boxes as tables with the probe/gene profiles categorized to the patterns, the probe/gene profiles not categorized to particular patterns and the probe/gene profiles categorized to patterns but have low correlation. See below.



Click “Spectrum” to launch a dialog box as a table with the probe/gene profiles categorized to the particular pattern in the drop down list.

For any of the tables, from the menu go to File -> Save Data or Save in Separate Files to save the probe/gene profiles categorized to patterns. See below.



If saving to separate files, the file names will be generated automatically (such as selectedGenesInPattern_2_261.txt) and saved in the working directory for EPIG. This is

file name contains the pattern number (#2) and the number of probes\genes (n=261) categorized to it.

Reference:

Chou JW, Zhou T, Kaufmann WK, Paules RS, Bushel PR. Extracting gene expression patterns and identifying co-expressed genes from microarray data reveals biologically responsive processes. *BMC Bioinformatics*. 2007 Nov 2;8:427