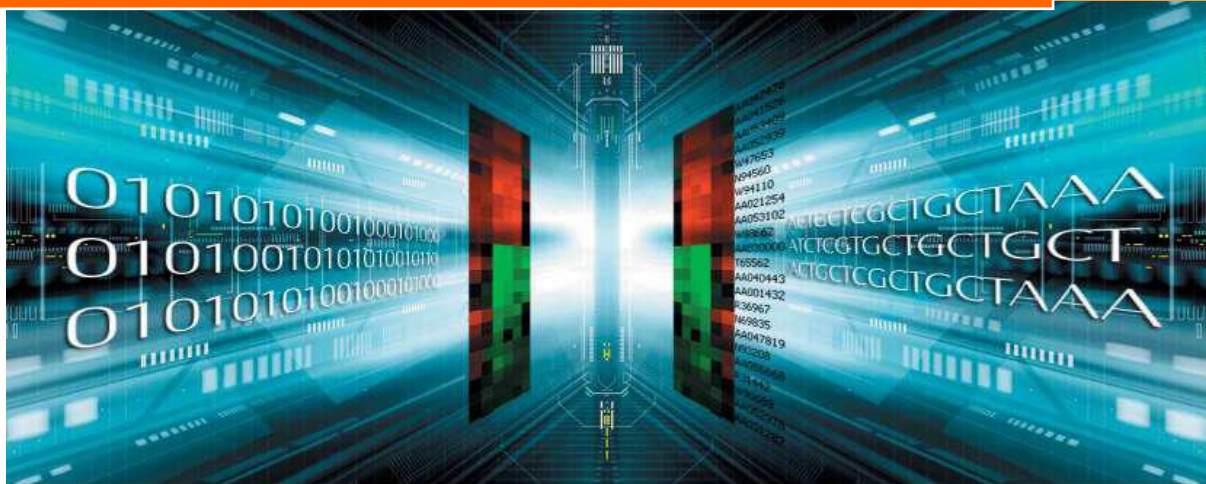


Version 3.0

RNA-Seq/Microarray DEG Analysis



(주)이바이오젠

서울특별시 영등포구 선유로13길 25
(문래동6가), 에이스하이테크시티2, 305호
Tel. 02-3141-0791

service@e-biogen.com

<http://www.e-biogen.com>

< 목 차 >

1. 엑셀기반 DEG 분석 (ExDEGA v.1.6.0)

2. Web 기반 Gene Set Enrichment 분석

2-1. DAVID tool을 이용한 Functional Annotation 분석

2-2. String-db tool을 이용한 gene set 분석

2-3. MSigDB기반 GSEA 분석

3. KEGG DB 기반 Pathway 분석

4. MeV Software 이용 Clustering Heatmap 작성

1. 엑셀기반 DEG 분석 (ExDEGA v.1.6.0)

㈜이바이오젠은 QuanSeq, mRNA-Seq, Total RNA-Seq 과 Micorarray data 를 엑셀 기반에서 DEG 를 쉽게 분석할 수 있도록 분석보고 시 ExDEGA (Excel based Differentially Expressed Gene Analysis) tool 을 함께 제공한다. ExDEGA 분석툴은 ㈜이바이오젠이 연구자들이 Microarray 및 RNA-Seq 데이터를 보다 쉽게 다루고 원하는 데이터를 쉽게 얻을 수 있도록 사용자 편의를 최대한 반영한 분석툴이고 엑셀 프로그램 안에서 다양한 분석을 직관적으로 수행할 수 있도록 개발되었다. ExDEGA 분석툴은 사용자들의 요구사항을 지속적으로 반영하여 데이터분석과 엑셀사용에 익숙하지 못한 연구자들도 쉽게 사용이 가능하도록 계속 업데이트 될 예정이다.

이바이오젠에서 제공하는 Microarray data 와 RNA-Seq data (엑셀 데이터)를 열기 전에 함께 제공한 ExDEGA(버전).zip 파일의 압축을 풀고 setup 을 실행하면 분석툴이 설치된다(그림 1-1).

설치가 완료되면 보고된 엑셀데이터를 열면 자동으로 ExDEGA 분석툴이 엑셀에 반영된 것을 확인할 수 있다. 참고로 ExDEGA 설치 전에 실행 중인 엑셀 파일이 있으면 종료시킨 후 다시 실행해야 ExDEGA 를 사용할 수 있다.

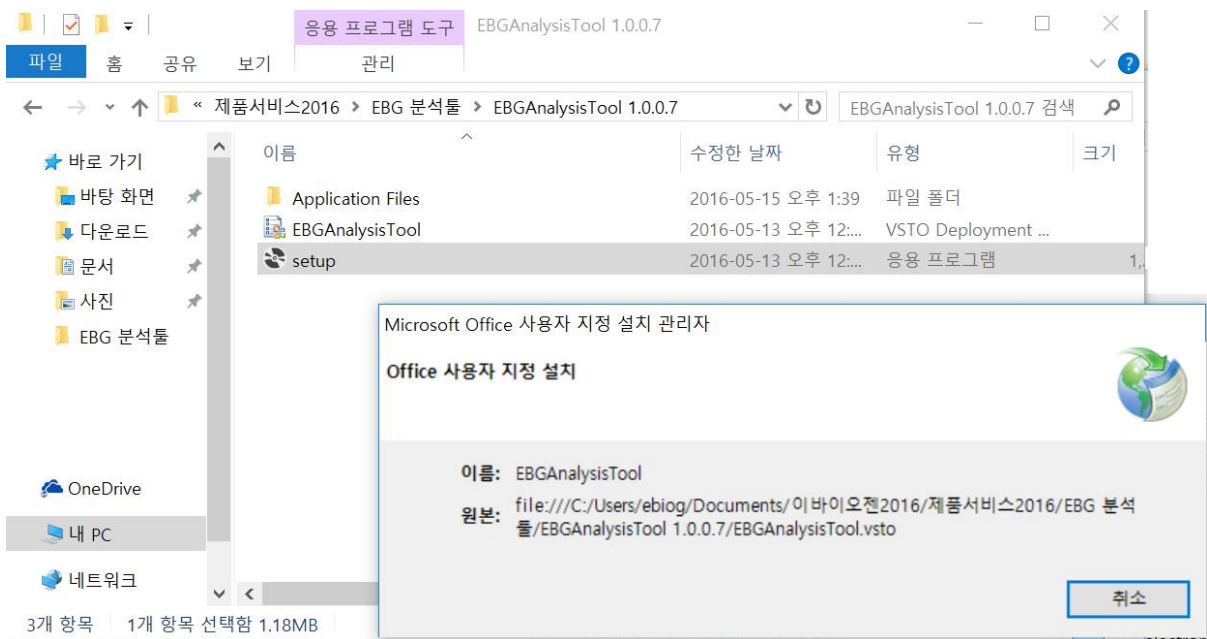


그림 1-1. ExDEGA set up

??? ExDEGA Report.xls 파일을 열면 왼쪽에 Gene Ontology (GO) 분석 창과 가운데에 mRNA expression data, 오른쪽에 DEG 분석 창이 나온다(그림 1-2).

GO 분석 창에서는 기본 설정된 GO 와 사용자가 원하는 대로 GO 를 구성하여 분석할 수 있고 DEG 분석과 함께 연동하여 데이터를 쉽게 얻을 수 있다. DEG 분석 창에서는 Fold change, Normalized RC, p-value 등을 선택하여 원하는 데이터를 쉽게 얻을 수 있고 GO graph 를 통해 전체적인 발현패턴을 확인할 수 있다. 뿐만 아니라, DEG 분석 창에서 Scatter Plot, Volcano Plot, Venn Diagram 을 직접 그릴 수 있고 필터링된 유전자들을 대상으로 Clustering heatmap 을 작성하기 위한 MeV 프로그램 input file 을 자동으로 만들 수 있고 Gene expression graph, Gene search 기능도 이용할 수 있어 연구자가 RNA-Seq data 를 쉽게 활용할 수 있다.

ID	Gene Symbol	Fold change				p-value				Average of Normalized RC (log2)					
		A/Cont	B/Cont	B/A	D/C	A/Cont	B/Cont	B/A	D/C	Control	A	B	C	D	Control
3	1.0610005238Rik	0.836	0.822	0.983	2.133	0.041	0.888	0.039	0.144	7.925	7.666	7.642	9.418	10.511	8.032
4	2.0610007919Rik	0.817	0.964	1.152	2.476	0.082	0.149	0.115	0.085	10.170	10.113	10.317	4.420	6.596	10.312
5	3.0610007919Rik	0.933	0.818	0.886	1.010	0.447	0.356	0.029	0.045	10.298	10.183	10.009	11.666	11.609	10.458
6	4.0610008107Rik	0.976	0.769	1.638	0.810	0.129	0.008	0.044	0.000	1.787	6.897	1.409	6.702	6.398	0.000
7	5.0610009814Rik	0.847	0.728	0.857	2.000	0.000	0.000	0.055	0.175	6.349	6.108	5.886	2.783	3.818	6.817
8	6.0610009822Rik	0.806	1.145	1.420	0.961	0.126	0.031	0.161	0.005	8.353	8.042	8.548	10.237	10.180	8.575
9	7.0610009900Rik	0.980	1.006	1.108	1.391	0.011	0.092	0.050	0.255	10.063	10.034	10.182	10.829	11.305	10.114
10	8.0610009118Rik	0.851	1.503	1.787	0.567	0.447	0.604	0.023	0.502	7.172	6.940	7.761	5.285	4.466	7.479
11	9.0610009200Rik	0.991	0.961	0.970	1.096	0.046	0.913	0.030	0.275	12.033	12.019	11.976	11.851	11.983	11.897
12	10.0610010808Rik	0.969	0.935	0.965	0.982	0.025	0.008	0.026	0.031	11.892	11.847	11.796	13.880	13.654	12.155
13	11.0610010905Rik	1.079	1.235	1.145	1.039	0.195	0.181	0.325	0.510	7.894	8.004	8.199	8.501	8.556	7.772
14	12.0610010918Rik	0.901	1.140	1.266	1.560	0.177	0.020	0.465	0.029	13.911	11.780	14.120	14.354	14.916	12.775
15	13.0610011068Rik	0.788	1.034	1.313	1.722	0.044	0.008	0.053	0.072	10.740	10.356	10.789	9.241	10.025	10.710
16	14.0610012689Rik	0.839	1.063	1.268	0.757	0.012	0.801	0.269	0.177	9.865	9.611	9.953	10.506	10.104	9.881
17	15.0610012838Rik	1.000	1.000	1.000	1.000	1.000	0.022	0.050	0.000	0.000	0.000	4.446	4.886	4.367	0.000
18	16.0610013029Rik	1.205	1.013	0.841	0.796	0.094	0.965	0.264	0.239	10.458	10.706	10.457	10.331	10.006	10.569
19	17.0610013108Rik	0.906	0.918	1.013	0.694	0.007	0.000	0.000	0.003	12.408	12.405	12.904	13.066	12.939	12.568
20	18.0610013116Rik	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000
21	19.0610013173Rik	0.867	1.140	1.315	0.788	0.028	0.478	0.190	0.099	10.036	9.830	10.225	9.969	9.621	9.974
22	20.0610013821Rik	0.811	0.786	0.920	1.164	0.235	0.331	0.680	0.300	6.312	6.010	5.890	5.454	5.673	6.557
23	21.0610013829Rik	1.119	1.000	0.947	1.132	0.354	0.885	0.737	0.466	8.124	8.287	8.206	4.537	8.716	8.302
24	22.0610013903Rik	0.822	0.730	0.863	0.626	0.155	0.161	0.470	0.097	12.194	11.911	11.699	9.216	8.093	12.132
25	23.0610014081Rik	0.847	1.141	1.764	0.846	0.173	0.841	0.047	0.148	3.857	3.228	4.047	4.940	1.895	3.878
26	24.0610014094Rik	1.633	1.121	0.688	0.511	0.097	0.892	0.275	0.077	8.192	8.900	8.357	10.137	8.025	8.475
27	25.0610014091Rik	0.795	0.847	1.066	0.908	0.022	0.213	0.021	0.434	10.076	9.744	9.836	5.313	4.324	10.168
28	26.0610014317Rik	1.049	1.169	1.114	0.658	0.780	0.090	0.015	0.000	2.774	2.444	2.600	4.005	3.977	2.972
29	27.1100011520Rik	1.351	1.945	2.550	1.000	0.455	0.000	0.284	1.000	4.569	5.103	6.454	0.000	0.000	4.830
30	28.1100011618Rik	1.005	1.766	1.756	0.816	0.032	0.011	0.154	0.000	8.064	8.072	8.884	9.878	9.585	8.243
31	29.1100011818Rik	0.842	1.944	2.311	0.962	0.198	0.014	0.304	0.014	8.103	7.855	9.063	8.872	8.816	7.920
32	30.1100012012Rik	1.009	0.969	0.960	1.209	0.040	0.047	0.061	0.189	7.980	7.993	7.934	8.611	8.885	8.080
33	31.1100014919Rik	0.774	1.115	1.442	1.314	0.045	0.910	0.021	0.117	4.904	5.121	4.851	9.065	9.458	9.451

그림 1-2. mRNA expression data format made in E-Biogen

1-1. Gene Category 사용 방법

mRNA expression data 는 수 만개의 유전자를 포함하기 때문에 유전자를 한 개씩 분석하기 보다 기능별로 그룹을 지어 분석을 하는 것이 용이하다. 이를 위해 많은 연구자들이 gene ontology (GO)를 활용한다. GO 는 비슷한 기능의 유전자들을 묶어 놓은 그룹이라고 생각하면 이해하기 쉽다.

Gene Category 창은 수많은 GO 중 임의로 15 개를 선택하여 관련 유전자를 필터링 할 수 있도록 만들어 놓은 것이다. 예를 들어, Aging 관련 유전자만 분석을 원할 경우, Gene Category 창에서 Aging 을 선택하면 해당 유전자 리스트만 필터링 된다(그림 1-3).

그리고 Gene Category 의 여러 항목들을 동시에 만족하는 유전자를 필터링할 수 있고 적어도 한 항목만이라도 포함하는 유전자를 보고자 하는 경우도 필터링이 가능하도록 "AND"와 "OR" 기능을 갖추고 있다.

Filter: 259		Fold change				p-value				Average of Normalized RC (log2)				
ID	Gene Symbol	A/Contr	B/Contr	B/A	D/C	A/Contr	B/Contr	B/A	D/C	Control	A	B	C	D
1775	1773 Abat	0.793	0.890	1.122	1.898	0.205	0.027	0.723	0.147	9.232	8.899	9.064	9.310	10.234
1989	1987 Adg	0.828	0.857	1.035	0.426	0.410	0.022	0.040	0.028	9.367	9.095	9.145	12.675	11.444
2094	2092 Adm	0.587	0.722	1.230	6.288	0.057	0.262	0.049	0.061	8.505	7.737	8.035	5.110	3.000
2108	2106 Adrala	1.097	0.846	0.771	0.873	0.286	0.023	0.040	0.011	11.544	11.678	11.304	3.859	3.702
2116	2114 Adrb3	0.935	0.989	1.058	0.466	0.819	0.944	0.042	0.040	9.378	9.281	9.363	11.878	9.298
2181	2179 Agt	1.169	0.891	0.762	0.238	0.497	0.044	0.324	0.111	5.895	6.120	5.728	2.672	0.498
2183	2181 Agtr1a	0.923	0.557	0.604	0.286	0.432	0.173	0.219	0.074	10.984	10.668	10.139	4.796	7.196
2308	2306 Akt1	0.988	1.040	1.053	0.987	0.802	0.000	0.044	0.048	14.471	14.452	14.527	12.348	12.330
2328	2326 Aldh3a1	0.632	0.552	0.874	1.724	0.110	0.009	0.045	0.107	12.316	11.653	11.460	8.772	6.827
2367	2365 Alox12	0.682	0.737	1.081	0.852	0.032	0.140	0.712	0.293	9.203	8.652	8.764	8.198	7.967
2402	2400 Amfr	1.000	0.977	0.977	0.463	0.000	0.088	0.001	0.037	13.080	13.080	13.046	11.417	10.324
2403	2401 Amh	0.681	2.794	4.101	1.050	0.523	0.219	0.146	0.017	3.548	2.995	5.031	5.470	5.541
2462	2467 Ankle1	1.760	2.715	1.543	0.840	0.217	0.006	0.032	0.035	3.529	4.345	4.970	8.828	8.576
2592	2590 Apaf1	1.239	1.083	0.875	0.704	0.187	0.040	0.440	0.008	10.751	11.060	10.867	13.604	11.885
2605	2603 Apex1	0.965	1.292	1.339	1.203	0.336	0.210	0.173	0.189	9.431	9.379	9.801	13.370	13.837
2634	2632 Apod	0.640	0.932	1.456	0.794	0.037	0.049	0.029	0.022	9.585	8.941	9.483	10.581	8.817
2698	2696 Arg1	1.384	0.774	0.560	3.447	0.234	0.478	0.205	0.500	5.052	5.520	4.683	0.000	1.785
2957	2955 Atg7	1.067	1.040	0.975	1.024	0.157	0.796	0.026	0.025	12.788	12.882	12.845	12.001	12.036
2965	2963 Atm	1.150	0.942	0.819	1.038	0.206	0.041	0.268	0.476	9.501	9.703	9.415	11.610	11.665
2967	2965 Atn1	1.124	0.866	0.771	1.320	0.168	0.037	0.324	0.179	13.817	13.985	13.609	14.134	14.535
2995	2993 Atp2b1	1.013	0.824	0.814	0.959	0.045	0.000	0.037	0.033	12.416	12.434	12.137	11.775	11.715
3049	3047 Atp8a2	1.283	1.022	0.796	1.182	0.194	0.856	0.044	0.042	5.129	5.488	5.160	5.078	5.319
3061	3059 Atr	0.978	0.898	0.918	0.942	0.047	0.133	0.181	0.319	9.364	9.332	9.208	10.770	10.684
3100	3098 Aurkb	1.071	1.368	1.277	1.131	0.026	0.167	0.058	0.152	7.273	7.373	7.725	13.114	13.292
3238	3236 Bak1	1.046	1.070	1.023	1.311	0.586	0.914	0.048	0.068	11.775	11.839	11.872	14.094	14.485
3267	3265 Bbc3	0.994	1.528	1.538	0.338	0.002	0.022	0.000	0.025	9.330	9.321	9.942	12.002	9.914
3401	3399 Bcl2	1.185	0.871	0.735	0.900	0.187	0.082	0.044	0.044	10.477	10.722	10.277	6.259	6.107
3436	3434 Becn1	0.991	1.083	1.093	0.744	0.773	0.026	0.046	0.169	10.938	10.925	11.053	10.723	10.297
3559	3557 Brca2	1.166	0.934	0.801	1.092	0.368	0.624	0.326	0.054	8.479	8.700	8.380	14.476	14.602
3759	15929 Pck1	0.391	1.248	3.194	3.468	0.286	0.042	0.046	0.214	8.897	8.542	10.217	0.602	2.396
3803	3801 Carvho	0.810	1.276	1.513	1.309	0.017	0.044	0.201	0.124	8.967	8.664	9.261	11.525	11.913

그림 1-3. Gene ontology (Aging) selection

'View All Data' 버튼을 누르면 필터를 해제하여 다시 전체 결과를 볼 수 있고 15개의 GO 중 관심 기능이 없다면 'Gene Category Settings' 버튼을 이용하여 Quick GO site 에서 다른 GO 를 추가할 수 있다(그림 1-4). '?' 버튼을 누르면 GO 추가하는 방법이 자세히 설명되어 있다.

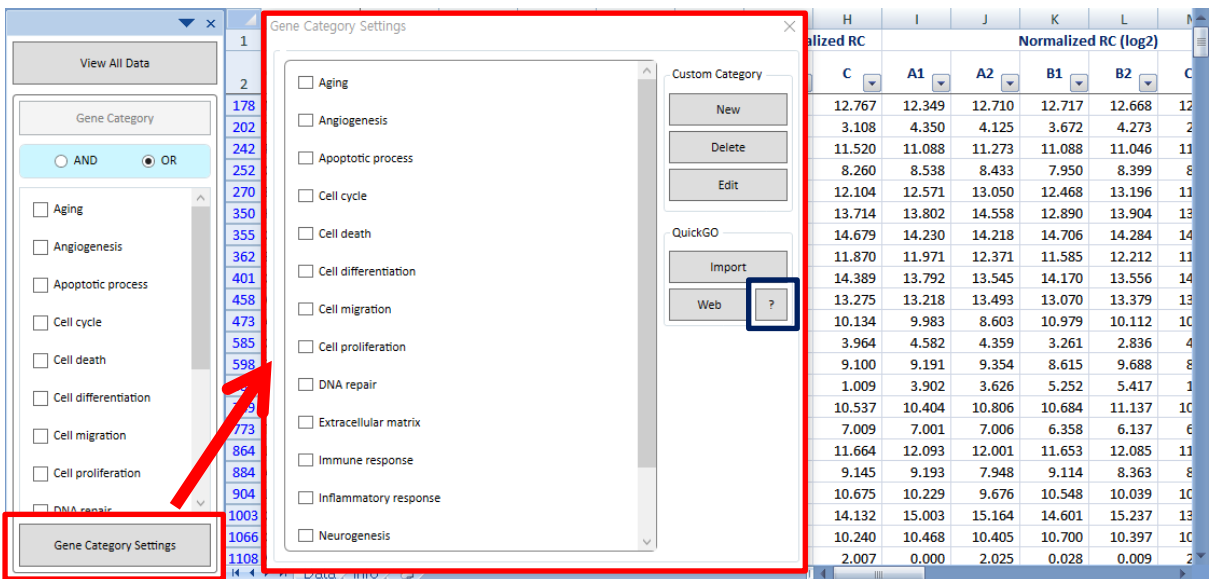


그림 1-4. Gene category settings

만약 원하는 유전자 그룹 목록이 있다면, 직접 입력하여 새로운 Gene Category 를 추가할 수도 있다. Gene Category Settings 버튼을 누른 후 New 를 선택하고 원하는 gene list 입력(or 복사-붙여넣기) 한 뒤, Gene category 이름 설정 후 저장하면 새로운 GO category 를 확인 할 수 있다(그림 1-5-a,b).

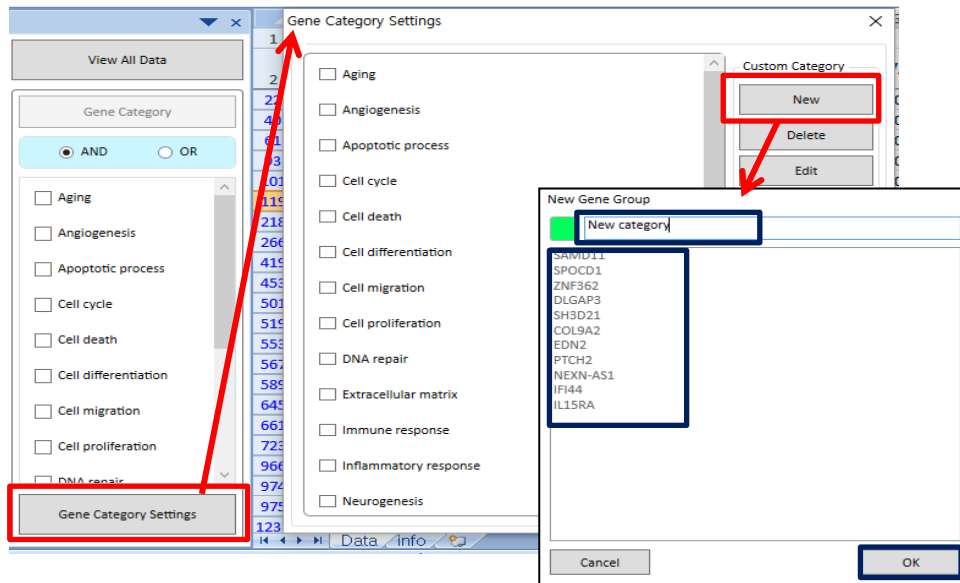


그림 1-5-a. Adding Genes to make a new gene category

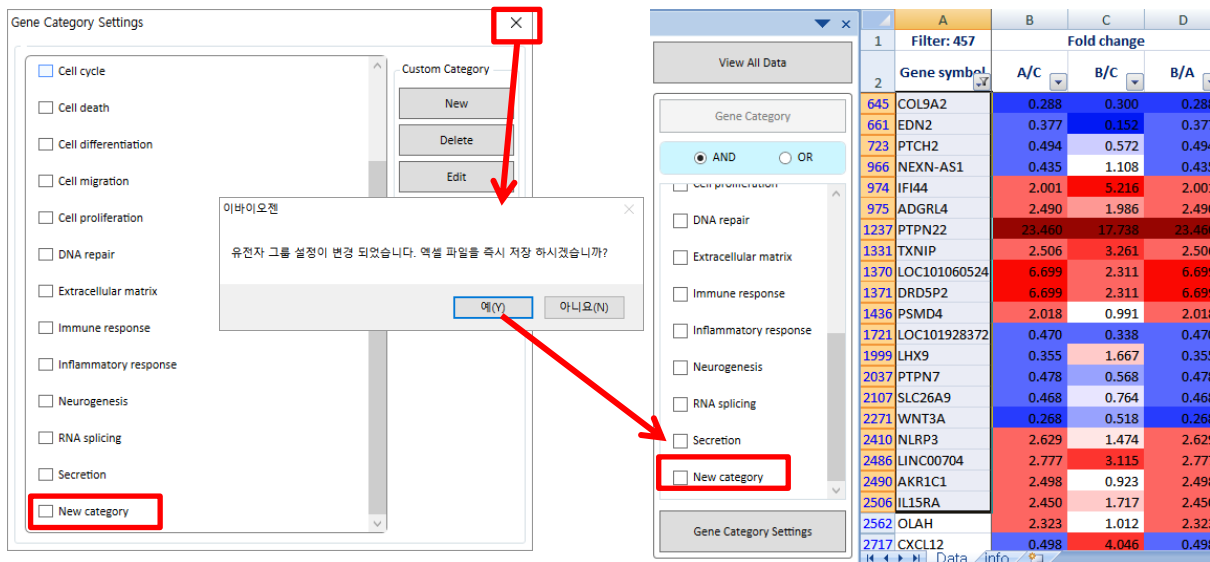


그림 1-5-b. Adding Genes to make a new gene category

1-2. Significant Gene Selection 사용 방법

오른편의 DEG Analysis 부분에서 "Significant Gene Selection" 창은 전체 결과 중 control 과 test 를 비교한 결과에서 유의하게 발현 차이가 나는 유전자를 필터링 할 수 있도록 만들어 놓은 것이다. 예를 들어, control 기준으로 A 에서 발현이 2 배 이상 증가 또는 감소하고, normalized RC(log)값이 4 이상이고, t-test 결과 p-value 값이 0.05 이하인 유전자(반복 실험한 데이터의 경우)를 선택하면 95 개의 유전자가 필터링 된다(그림 1-6).

그리고 여러 개의 비교그룹에서 동시에 Significant gene 을 선별하고자 할 경우와 적어도 한 비교그룹에서 Significant gene 을 선별하고자 할 경우에는 "AND"와 "OR" 기능을 사용하면 된다.

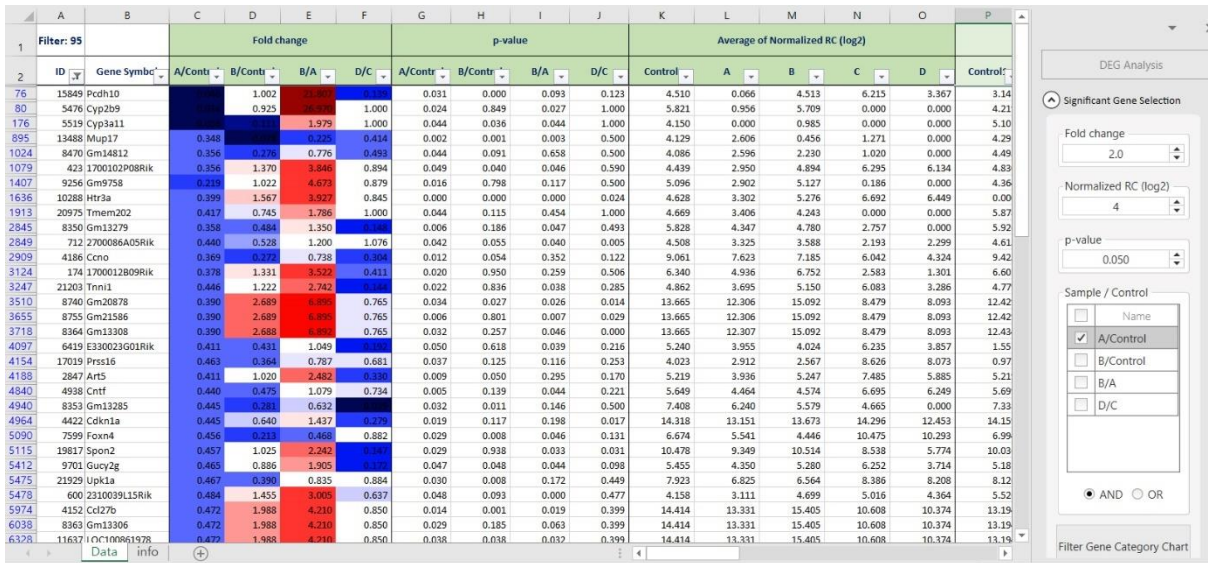


그림 1-6. Significant gene selection

Gene Category 와 Significant gene selection 은 연동 가능하다. 그림 1-7 에서 처럼 Gene Category 의 Cell differentiation 을 선택하면 10 개의 유전자가 필터링 된다(그림 1-7). 10 개의 유전자는 본 데이터에서 Cell differentiation 관련 유전자들 중 A/Control 비교그룹에서 유의하게 발현이 증가 또는 감소한 유전자를 의미한다.

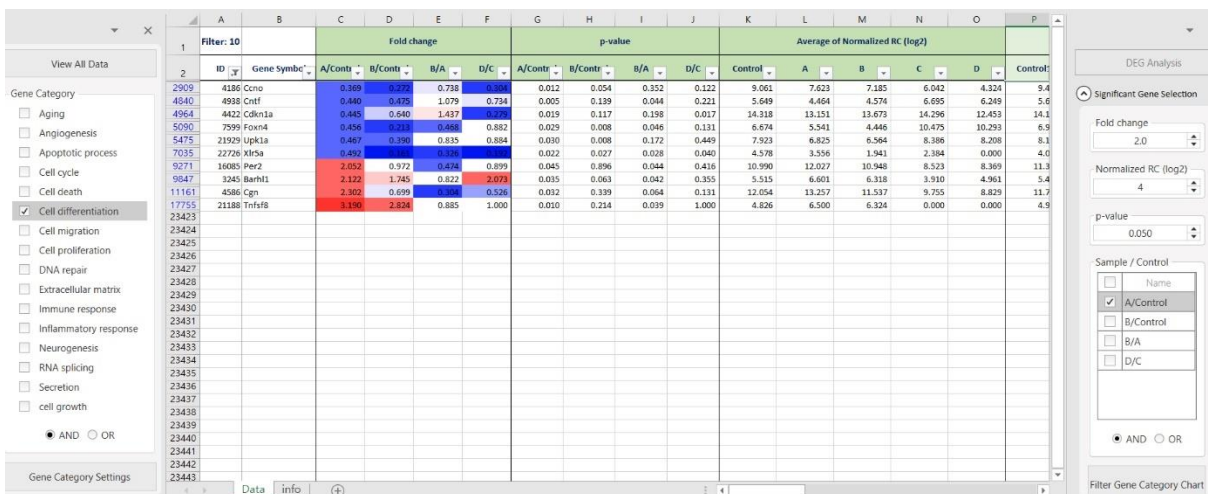


그림 1-7. Significant genes related to Cell cycle

실험 결과에 따라 발현 변화값 (fold change), p-value, normalized RC(log2) 기준을 조정할 수 있고 반복 실험인 경우만 p-value 를 선택할 수 있다.

“View Gene Category Chart” 버튼을 누르면 각 GO 관련 유전자 중 발현이 유의하게 차이 나는 유전자의 %와 수가 그래프로 그려진다. 본 분석을 통해 어떤 GO 의 유전자들이 상대적으로 많은 발현 변화가 있었는지를 확인할 수 있다. 전체 데이터 상태에서 Significant Gene Selection 의 비교 그룹을 선택하고 “View Gene Category Chart”를 클릭하면 증가/감소한 유전자 들 대상으로 GO Chart 가 생성된다. 그래프의 각 영역을 클릭하면 해당 유전자들이 필터링 된다. 예를 들어 왼쪽의

Pie chart 의 특정영역을 클릭하면 해당 GO 의 증가/감소된 유전자가 함께 필터링 되고 오른쪽의 증가/감소된 bar chart 에서 bar 상단의 숫자는 해당 유전자 수이고 bar 를 클릭하면 해당 유전자가 필터링 된다(그림 1-8).

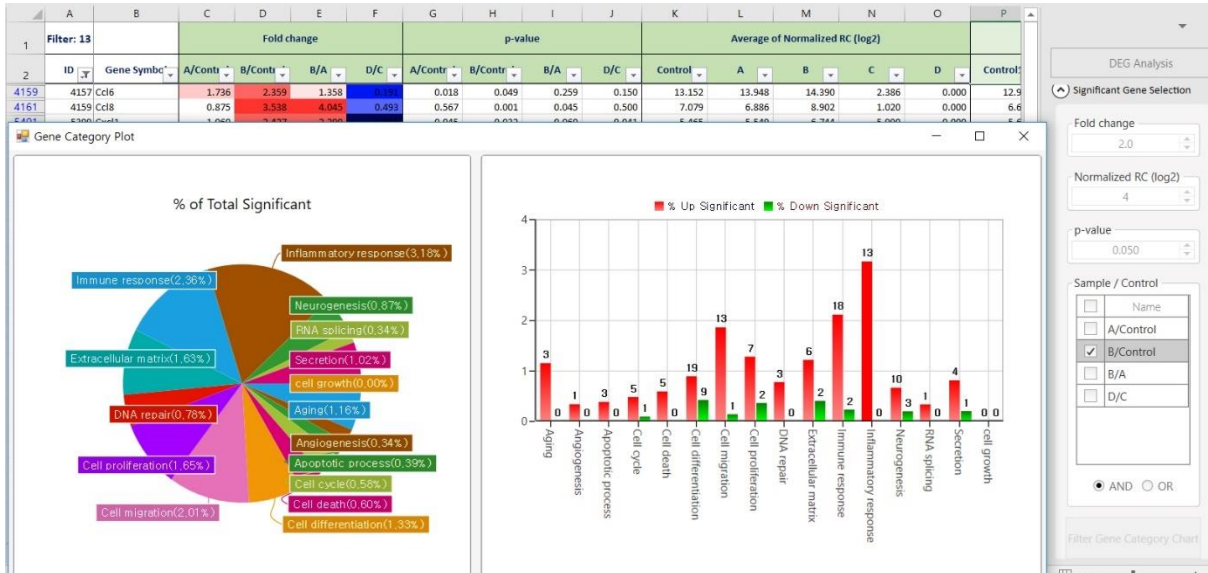


그림 1-8. View Gene Category Chart

1-3. Analysis Graph 사용 방법

DEG Analysis 부분에서 "Analysis Graph" 창을 펼치면 아래 그림 1-9 와 같이 Scatter Plot, Volcano Plot, Venn Diagram 을 엑셀에서 쉽게 그릴 수 있다.

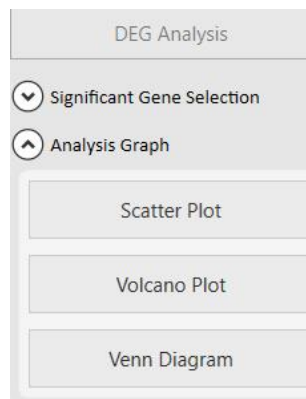


그림 1-9. Analysis Graph Tool

첫번째 Scatter Plot 은 오른쪽에 샘플 비교그룹을 선택하고 Fold threshold line 을 선택하고 "Graph View"를 클릭하면 왼쪽에 선택한 비교그룹을 대상으로 Scatter Plot 이 자동 생성된다. Plot 에서 특정 spot 을 클릭하면 해당 유전자가 표시되고 마우스 오른쪽을 클릭하여 표시를 지울 수도 있다. 그리고 여러 개의 유전자를 동시에 표시하고 싶다면 "Gene Select(ID Input)" 창에 해당 유전자 ID 를 복사하여 입력하고 "Add"를 클릭하면 Gene Symbol 이 자동 생성된다(그림 1-10).

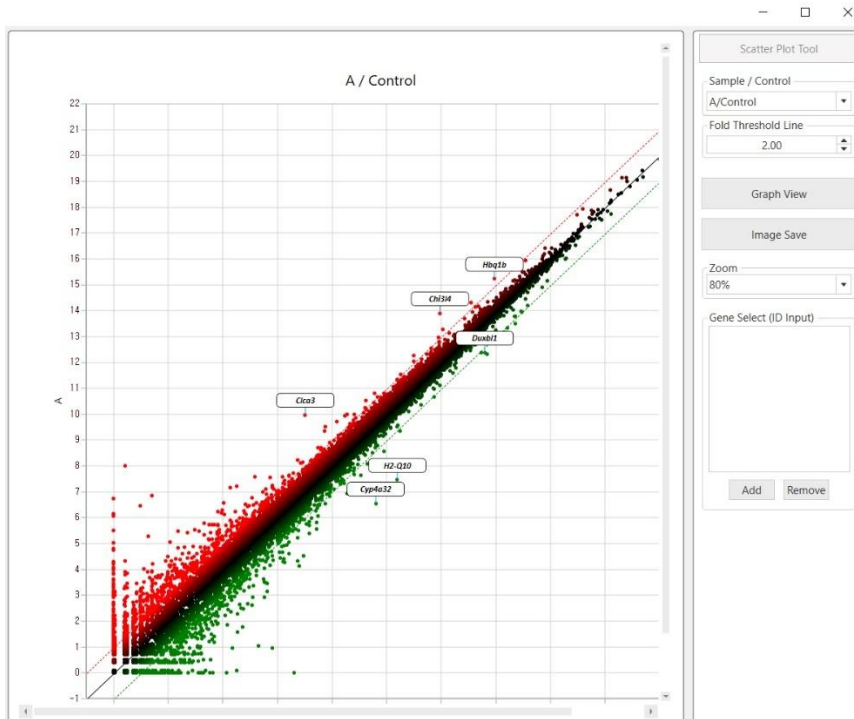


그림 1-10. Analysis Graph Tool – Scatter Plot

두번째 Volcano Plot 은 Scatter Plot 의 기능과 거의 동일한데 오른쪽에 샘플 비교그룹을 선택하고 Fold threshold line 과 p-value 를 선택하고 “Graph View”를 클릭하면 왼쪽에 선택한 비교그룹을 대상으로 Scatter Plot 이 자동 생성된다. Plot 에서 특정 spot 을 클릭하면 해당 유전자가 표시되고 마우스 오른쪽을 클릭하여 표시를 지울 수도 있다. 그리고 여러 개의 유전자를 동시에 표시하고 싶다면 “Gene Select(ID Input)” 창에 해당 유전자 ID 를 복사하여 입력하고 “Add”를 클릭하면 Gene Symbol 이 자동 생성된다(그림 1-10).

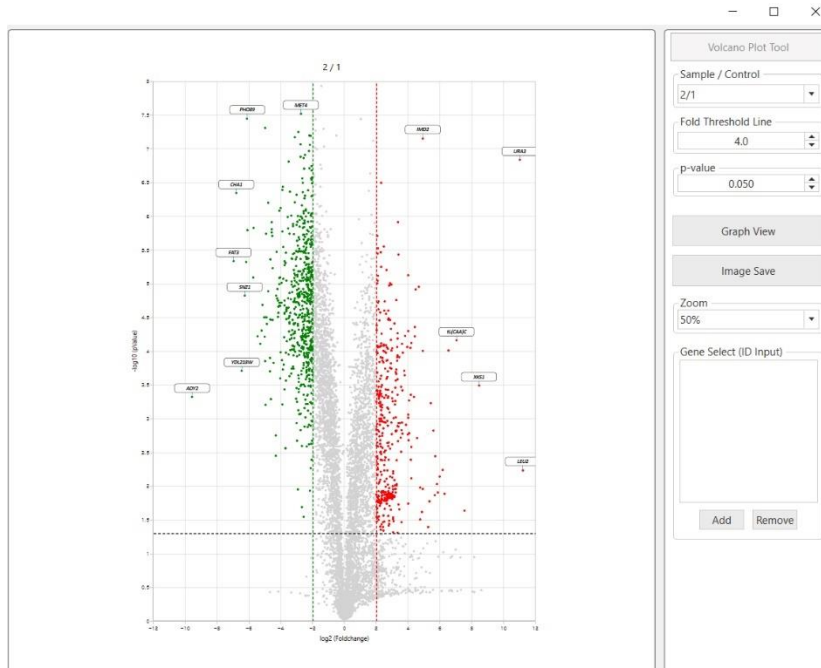


그림 1-11. Analysis Graph Tool – Volcano Plot

세번째 Venn Diagram 을 통해 2 개, 3 개 또는 4 개 까지의 비교그룹을 대상으로 Venn Diagram 을 작성할 수 있다. Venn Diagram 을 그릴 샘플 비교그룹과 Fold Change, p-value(반복실험시)을 선택 후, Diagram View 를 클릭하면 결과를 확인할 수 있으며 그룹은 최대 4 그룹까지 선택 가능하다. 아래의 그림은 A/C 와 B/C, B/A 결과 중, 2fc 이상 up, down 된 list 를 가지고 Venn Diagram 을 작성한 결과이다(그림 1-12).

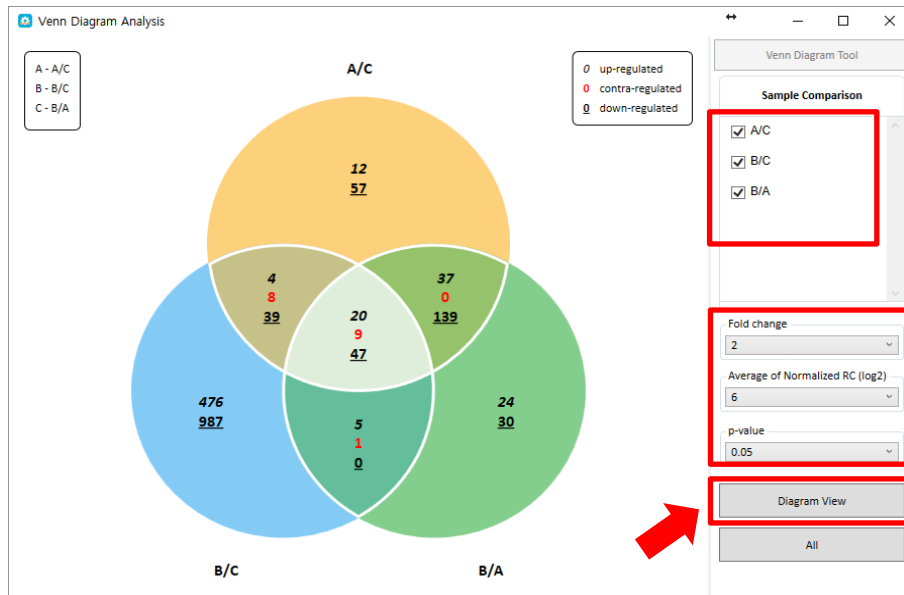


그림 1-12. Analysis Graph Tool – Venn Diagram

Venn Diagram 결과에서 표시되는 형식은 다음과 같다(그림 1-13).

1. **기울어진 숫자** : 2fold 이상 up-regulated 된 gene 수
2. **빨간색 숫자** : regulation 이 대조되는 gene 수
3. **밑줄 친 숫자** : 2fold 이상 down-regulated 된 gene 수



그림 1-13. For example of up ,down, contra-regulated in Venn Diagram

Venn Diagram 이미지를 오른쪽 클릭하면 Venn Diagram 각 영역에 어떤 유전자들이 있는지 확인할 수 있다. 예를 들어, A/C 에서만 2fold up 이 되는 유전자를 보고 싶으면, Venn Diagram 에서 A/C 에서만 해당되는 영역을 찾아 마우스 오른쪽 클릭 하면 2fold up 된 유전자 list 4 개가 엑셀 sheet 에 filter 된다(그림 1-14).

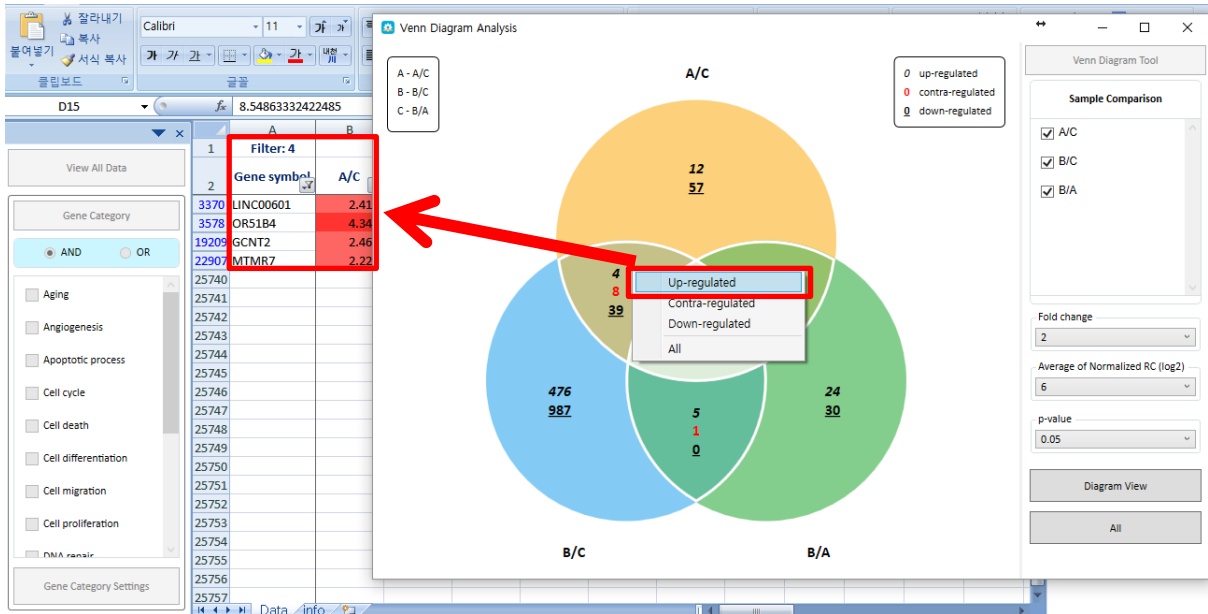


그림 1-14. Filtering 2fold up-regulated gene list in Venn Diagram

ExDEGA 에서 제공되는 모든 이미지는 오른쪽마우스를 눌러 'Save image' 버튼을 통해 저장이 가능하다(그림 1-15).

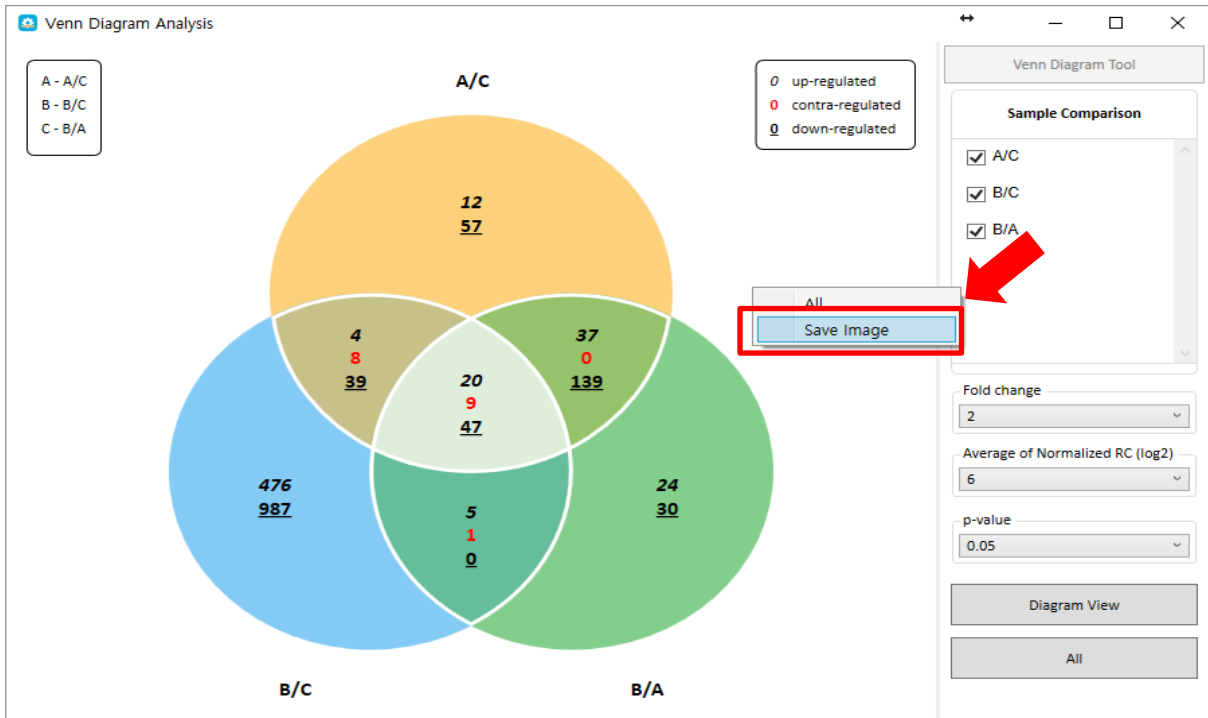


그림 1-15. Save image

1-4. Clustering Heatmap Support 사용 방법

ExDEGA 의 DEG Analysis 에서는 Significant Gene Selection 또는 Venn Diagram 등을 통해 Data Mining 을 수행한 후 정리된 유전자 리스트를 대상으로 Clustering Heatmap 을 쉽게 작성할 수 있도록 지원한다.

당사에서 추천하는 Clustering Heatmap 프로그램은 MeV 인데 ExDEGA 에서 MeV 용 Input file 을 자동 생성해 주고 MeV 에서 해당 파일을 불러오면 된다. 이후의 Clustering 방법 및 이미지 가공 및 저장 방법은 본 매뉴얼 “4. MeV Software 이용 Clustering Heatmap 작성” 부분을 참고하면 된다.

그림 1-16 에서 필터링된 유전자 리스트를 대상으로 Clustering Heatmap 을 작성하려면 크게 2 종류의 데이터를 이용할 수 있는데 첫번째는 Fold change 값을 이용할 시 Type 부분에 Fold change 를 체크하고 Export Data Select 에서 Heatmap 에 표현할 비교그룹을 체크하여 “Data Export”를 클릭한 후 “???.txt”로 저장하면 된다. 두번째는 발현값(Raw Data(RC))으로 표현하고자 할 때 Raw Data 를 체크하고 샘플이 3 개 이상이면 z-score 를 체크하고 샘플이 2 개면 median 을 체크하고 Export Data Select 에서 Heatmap 에 표현할 비교그룹을 체크하여 “Data Export”를 클릭한 후 “???.txt”로 저장하면 된다

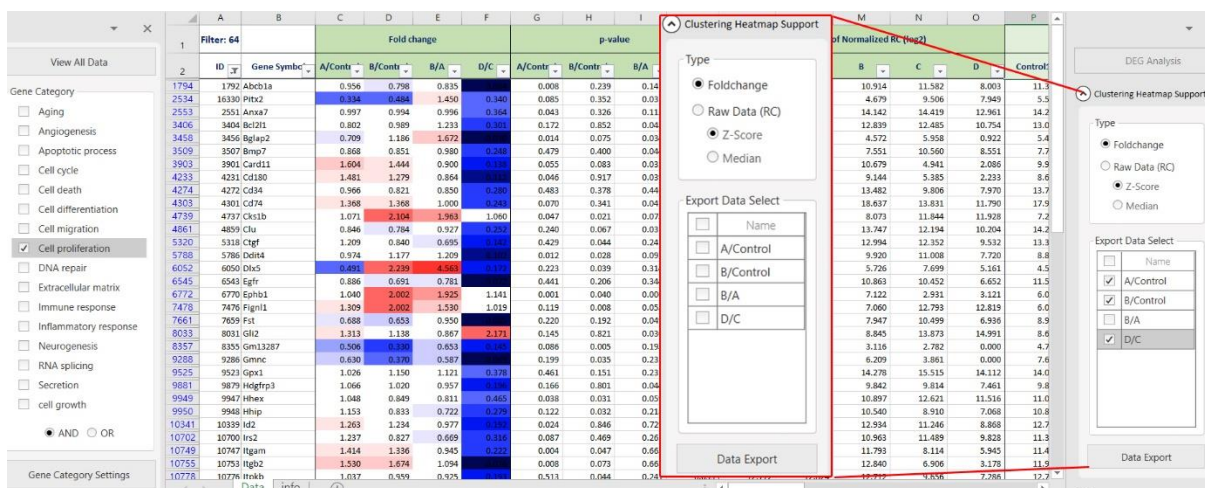


그림 1-16. Clustering Heatmap Support

1-5. Selected Gene Plot & Gene Search 사용 방법

ExDEGA 의 기능 중에 선별한 유전자 또는 연구자가 관심있는 유전자들을 대상으로 발현패턴을 그래프로 표현하고자 할 때는 “Selected Gene Plot” 기능을 사용하면 된다.

선별한 유전자의 gene symbol 을 복사하여 Selected Gene Plot 창에 붙여 넣고 “Expression Plot View”를 누르면 normalized RC(log2) 값, fold change 값으로 line graph 가 그려진다(그림 1-17).

그리고 특정 keyword 관련 유전자를 검색하고 싶을 때는 gene search 창을 이용하면 된다. 예를 들어 ‘insulin’을 검색하면 엑셀 Data Sheet 에 ‘insulin’ keyword 을 포함하는 모든 유전자가 검색되어 필터링 된다(그림 1-18).

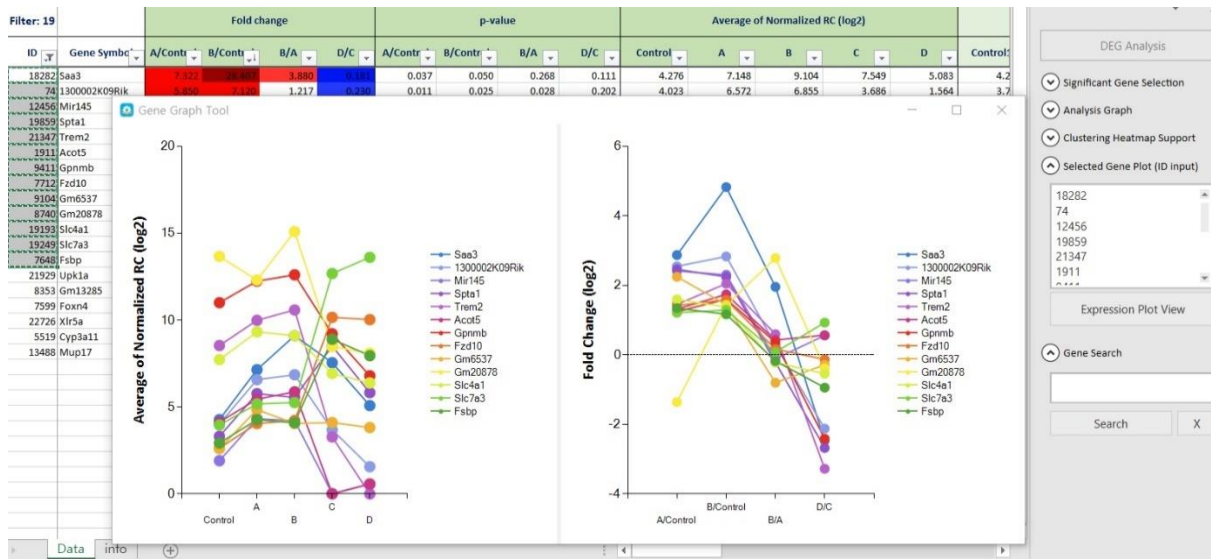


그림 1-17. Gene graph

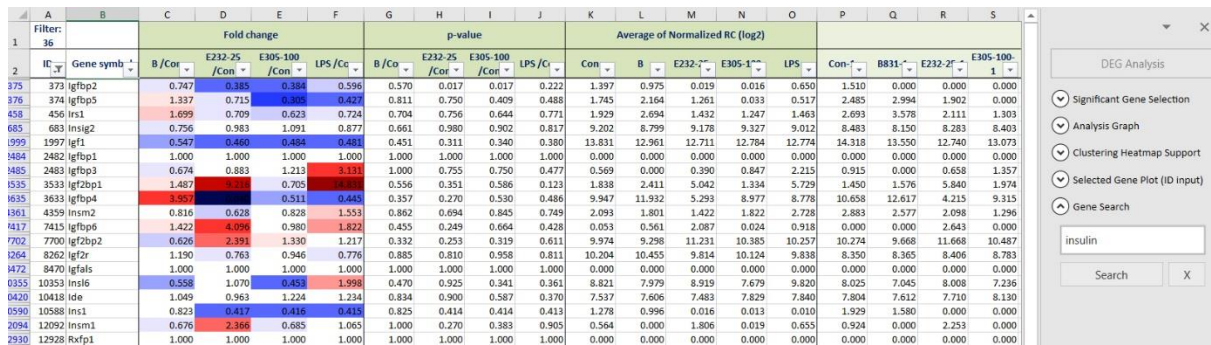


그림 1-18. Genes related to insulin

2. Web 기반 Gene Set Enrichment 분석

2-1. DAVID tool을 이용한 Functional Annotation 분석

DAVID 는 다양한 데이터 베이스를 기반으로 유전자의 상관관계를 통계적으로 분석하여 유전자의 주요 기능을 예측하는 analysis tool 이다. 분석과정은 그림 2-1 과 같다.

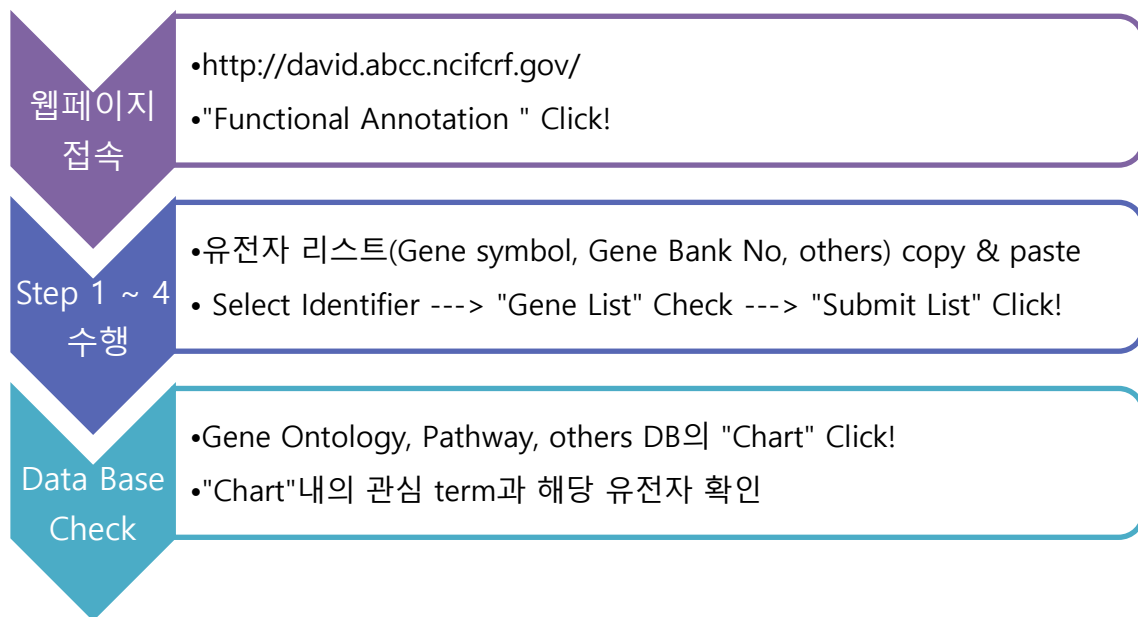


그림 2-1. DAVID tool analysis process

DAVID 에서는 3 천 개 이상의 유전자는 분석할 수 없으므로 3 천 개 이하로 유전자를 선별해야 한다. mRNASeq 결과에서 significant gene 을 선별하여 DAVID 분석을 한다. DAVID 홈페이지 (<http://david.abcc.ncifcrf.gov/>)에 접속하여 "Functional Annotation"을 클릭한다(그림 2-2).

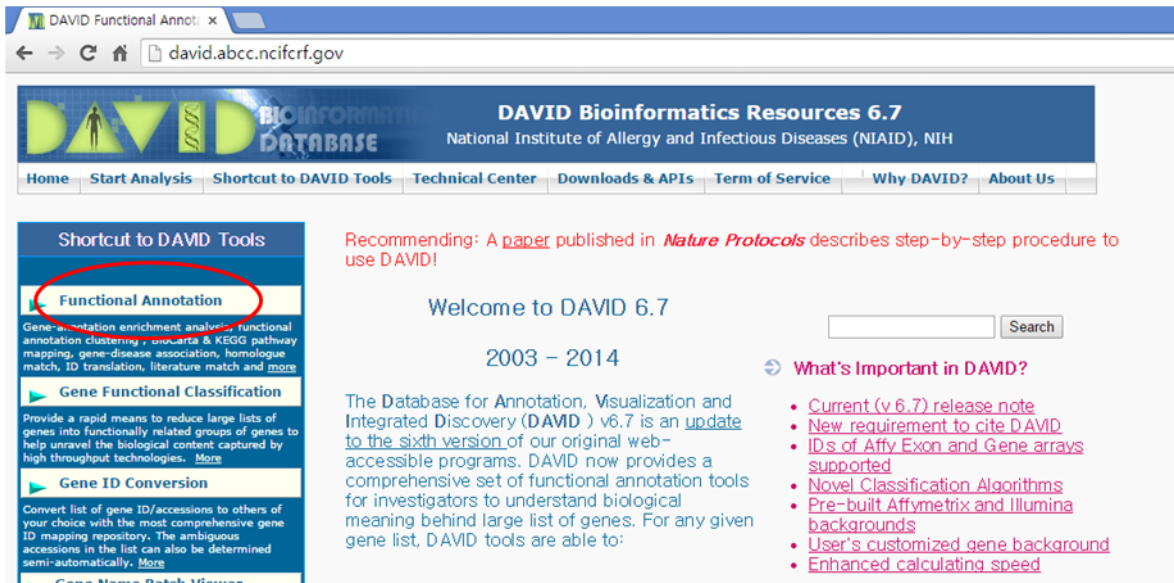


그림 2-2. DAVID tool webpage

“Upload” 탭에서 Step 1 에서 Step 4 까지 수행한다(그림 2-3). Step 1 에서선별한 유전자의 Gene Symbol 을 복사하고 “A: Paste a list” 창에 붙여 넣는다. Step 2 에서“OFFICIAL_GENE_SYMBOL”를 선택한다. 만약 step 1 에서 Gene Bank No.를 넣었다면 “GENEBANK_ACCESSION” 을 선택한다. Step 3 에서 “Gene List”를 체크하고 Step 4 에서 “Submit List”를 누른다. Gene Symbol 을 넣은 경우,“multiple species have been detected in your gene list”라는 창이 뜨면“확인”을 누른다.

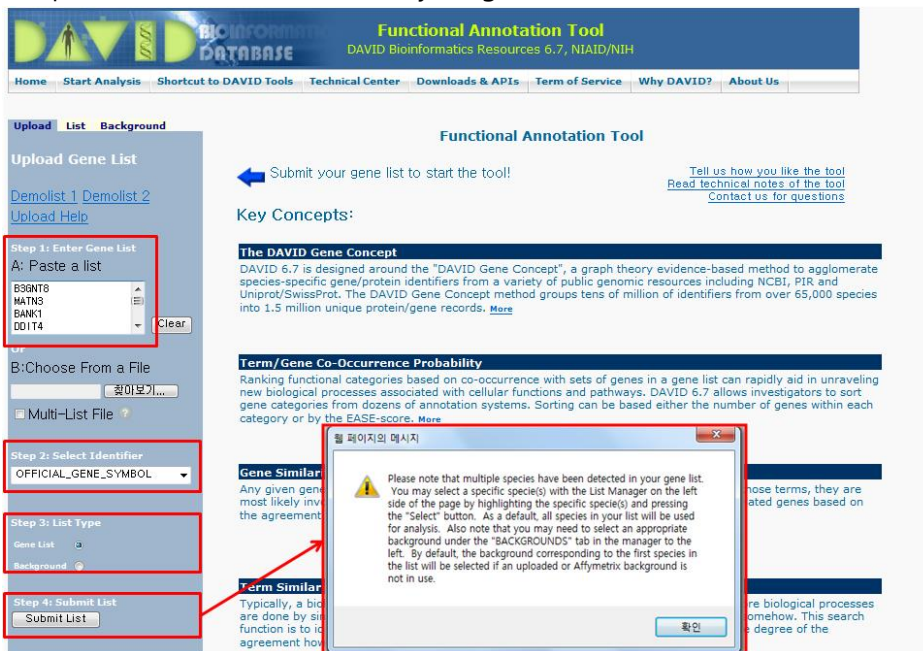


그림 2-3. DAVID tool : Step 1 ~ Step 4

실험한 종을 선택하고 “Select Species”를 누르면 해당 종의 유전자를 대상으로 다시 분석된다. 예시에서는 160개의 유전자 리스트를 넣었지만 데이터베이스에서 기능이 밝혀진 94개이기에 최종 94 개 유전자를 대상으로 Functional Annotation 분석이 완료되었다(그림 2-4).

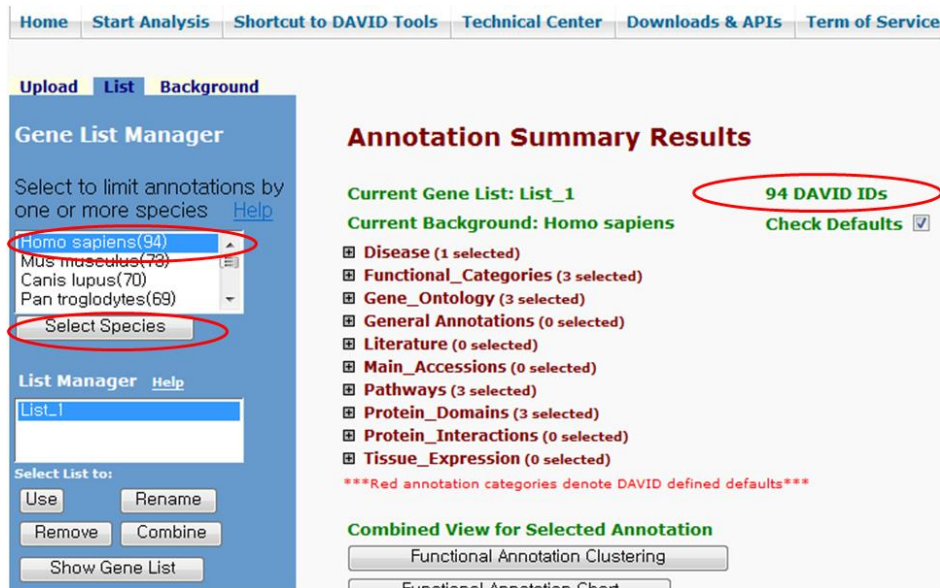


그림 2-4. DAVID tool : Select Species

분석 결과를 확인하기 위해 예로 Gene Ontology 중 Biological Process 를 확인한다.“Gene_Ontology”의 “+” 표시를 클릭하여 결과 창을 열고 “GOTERM_BP_FAT”의 “Chart”를 누르면 94 개 유전자들이 관여하는 Biological Process 에 속하는 GO 를 확인할 수 있다(그림 2-5). 관심 GO 를 클릭하면 QuickGO 데이터베이스로 연결되어 각 GO 의 정보를 확인할 수 있다.GO 의 Gene 막대를 클릭하면 해당 GO 관련 유전자들을 확인할 수 있다.

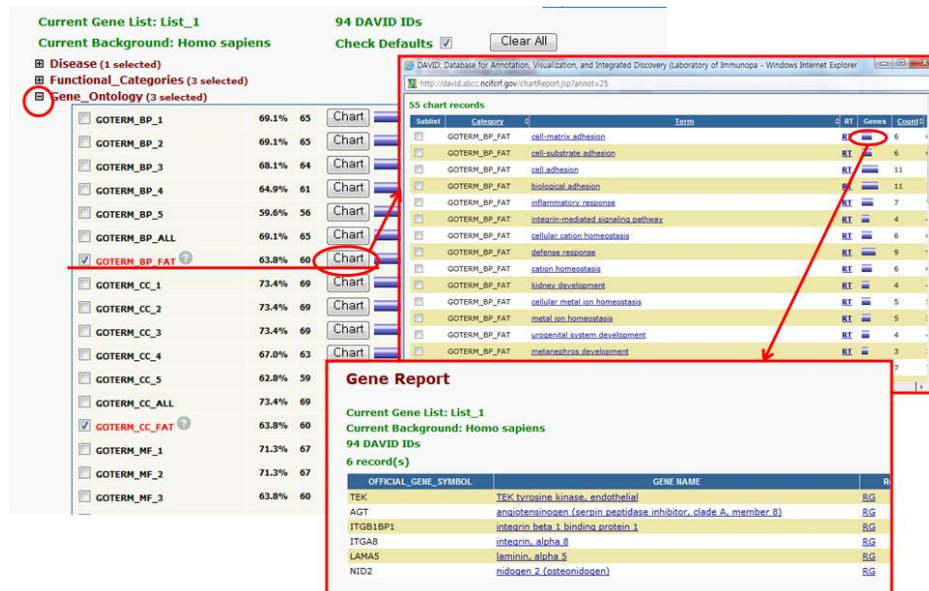


그림 2-5. DAVID tool : exploring Gene Ontology analysis result

이와 같은 방법으로 Pathway 결과를 확인해 보면 KEGG_PATHWAY database 에서 주요 Pathway 가 나온다(그림 2-6).각 pathway 를 누르면 pathway 그림을 확인할 수 있다. pathway 그림에서 별 표시가 되어 있는 유전자가 input 유전자(160 개) 중 해당 pathway 에 관여하는 유전자이다. 유전자를 클릭하면 유전자 정보를 자세히 알 수 있다.

Annotation Summary **Functional Annotation Chart** [Help and Manual](#)

Current Gene List: List_1
Current Background: Homo sapiens
156 DAVID IDs

Options

2 chart records

Sublist	Category	Term	RT	Genes	Count	%	P-Value	Benjamini
	KEGG_PATHWAY	Pathways in cancer	KE		7	4.5	3.7E-2	9.1E-1
	KEGG_PATHWAY	Regulation of actin cytoskeleton	KE		5	3.2	8.1E-2	9.3E-1

Annotations: Disease (1 selected), Functional_Categories (3 selected), Gene_Ontology (3 selected), General_Annotations (0 selected), Literature (0 selected), Main_Accessions (0 selected), Pathways (3 selected), BBID, BiDCARTA, EC_NUMBER (23.7% 37), KEGG_PATHWAY (25.6% 40), PANTHER_PATHWAY (12.8% 20), REACTOME_PATHWAY (21.2% 33)

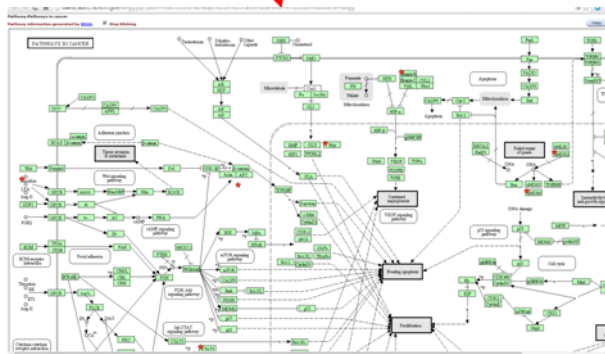


그림 2-6. DAVID tool : exploring Pathway analysis result

DAVID 분석은 input 한 유전자들이 유의하게 관련되는 GO, pathway 등을 분석하는 tool 이다. 즉, input 한 유전자에서 많은 유전자들이 관련되는 GO, pathway 만 결과로 나오기 때문에 input 유전자 중 적은 수가 관련되는 GO, pathway 는 결과에 나오지 않는다. 또한 input 유전자의 수가 적으면 분석 결과가 없을 수도 있다. DAVID 에서는 유전자 2 개 이상, EASE score 0.1 이하를 default 로 분석하여 이 기준에 적합한 결과를 보여준다. option 에서 이 기준을 조정할 수 있다. David 분석 결과의 각 항목은 DAVID 홈페이지의 Help and Tool Manual 에 자세히 설명되어 있다(그림 2-7).

Annotation Summary Results [Help and Tool Manual](#)

Current Gene List: List_1
Current Background: Homo sapiens
94 DAVID IDs
Check Defaults Clear All

Annotations: Disease (1 selected), Functional_Categories (3 selected), Gene_Ontology (3 selected), General_Annotations (0 selected), Literature (0 selected)

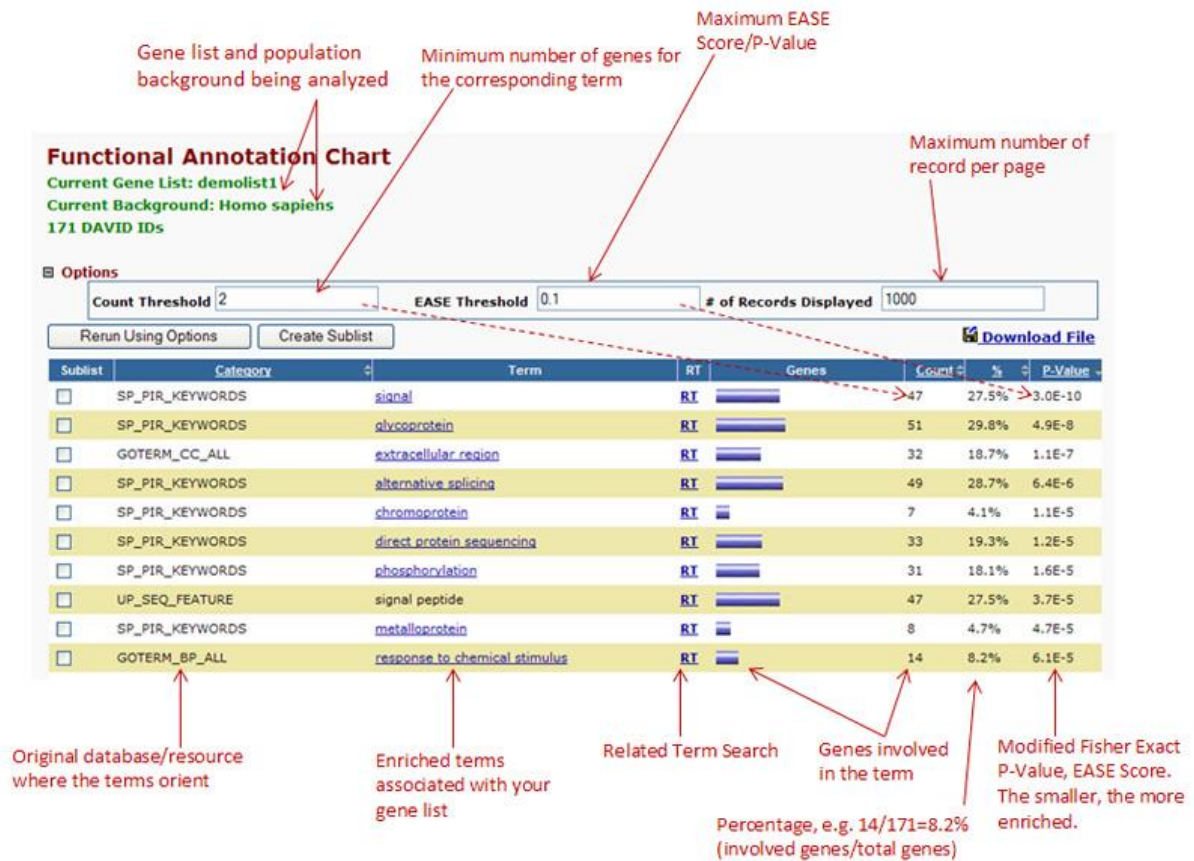


그림 2-7. DAVID Help and Tool Manual

2-2. String-db tool을 이용한 gene set분석

String-db tool 은 Protein-Protein Interaction 데이터 베이스를 기반으로 유전자의 상관관계를 통계적으로 분석하여 유전자의 주요 기능을 예측하고 Network 을 build 해 주는 분석툴이다. 분석과정은 그림 2-2-1 과 같다.

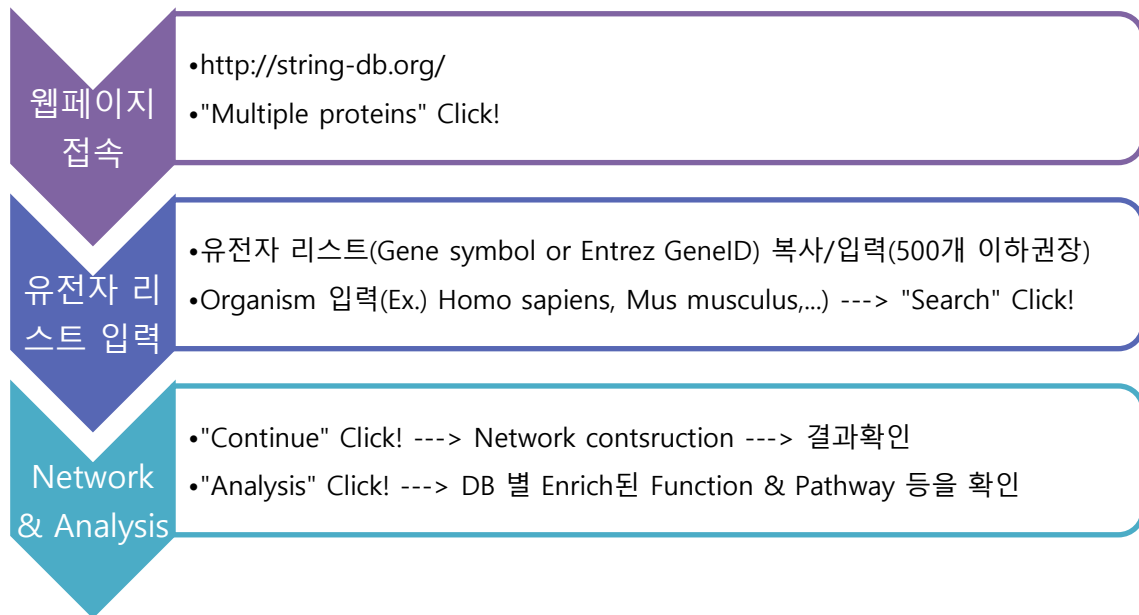


그림 2-2-1. String-db tool analysis process

String-db 에서는 500 개 이하의 유전자를 input 하는 것을 권장하고 있고 여러 public ID 중 EntrezGeneID 사용이 좀더 편리하다. mRNA-Seq 결과에서 significant gene 을 선별하고 String-db 홈페이지 (<http://string-db.org/>)에 접속하여 "Multiple proteins"을 클릭하고 "List of names" 입력창에 유전자 리스트를 복사한다.그리고 "Organism" 입력창에 해당 species 학명을 입력하고 "Search"를 클릭한다(그림 2-2-2).

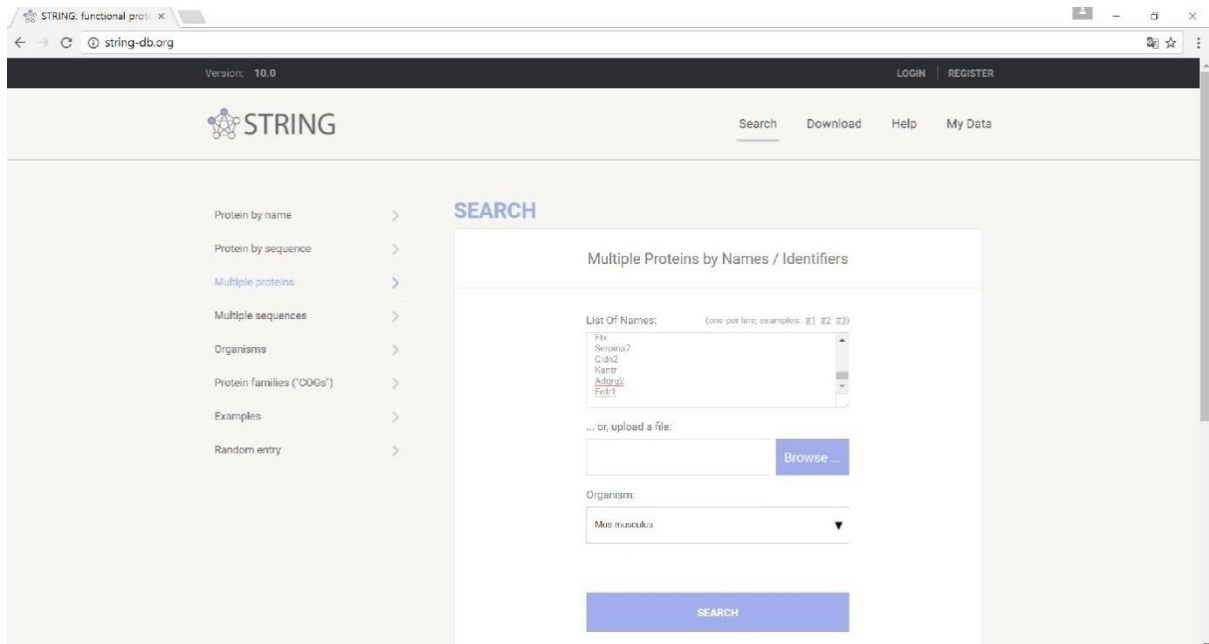


그림 2-2-2. Multiple proteins search

“Search” 결과 중간에 아래 그림과 같은 유전자 확인 단계가 있고 별 이상이 없으면 “continue”를 클릭하여 계속 진행한다(그림 2-2-3).

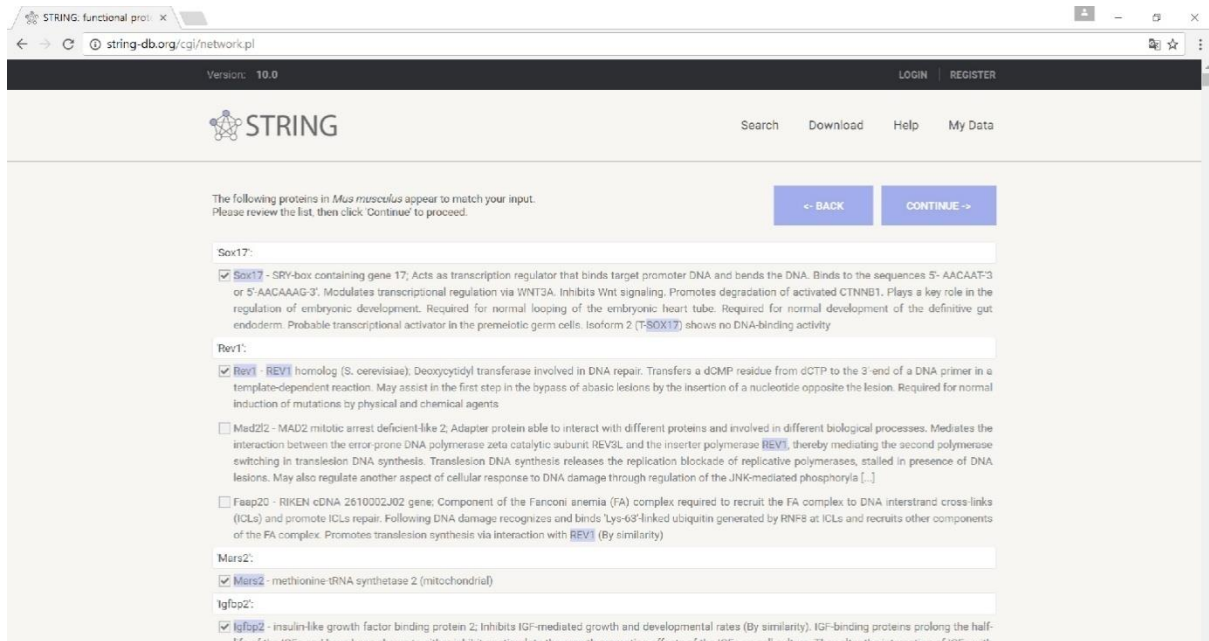


그림 2-2-3. Gene confirmation step

분석이 완료되면 그림 2-2-4와 같이 String DB 기반 Network 결과를 확인할 수 있고 “Analysis” 탭을 클릭하면 “Functional enrichments in your network” 결과를 확인할 수 있다(그림 2-2-5). 각 Functional DB 결과의 오른쪽 하단에 “more”를 클릭하면 FDR<0.05 이하에 해당하는 항목을 모두 볼 수 있다.

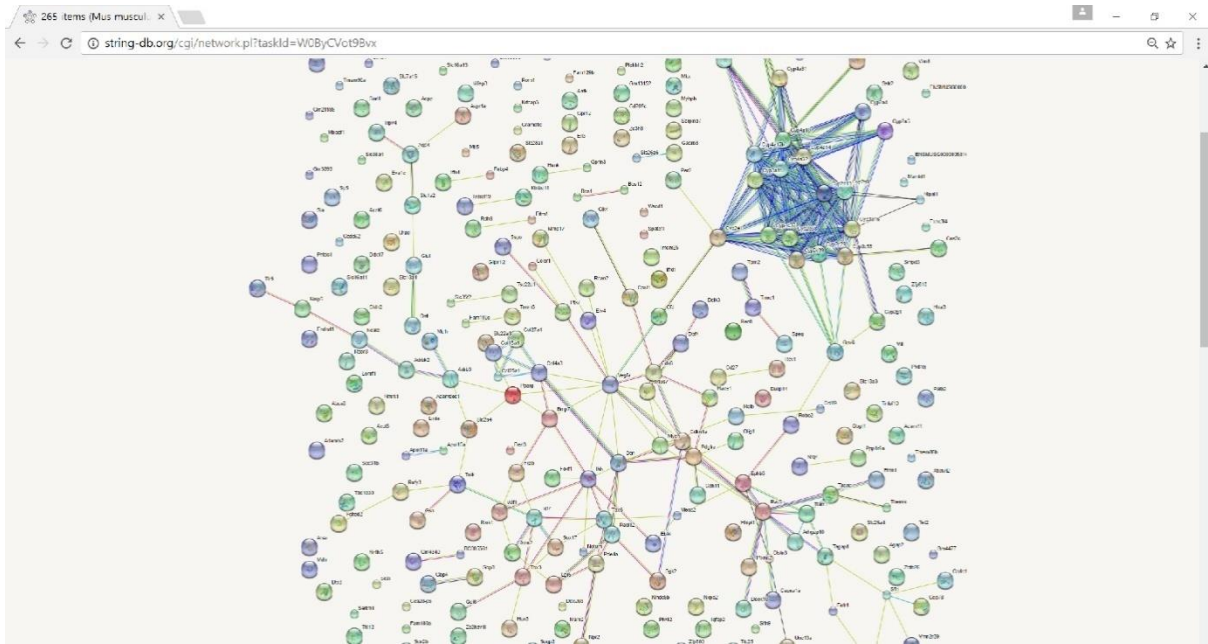


그림 2-2-4. String network result

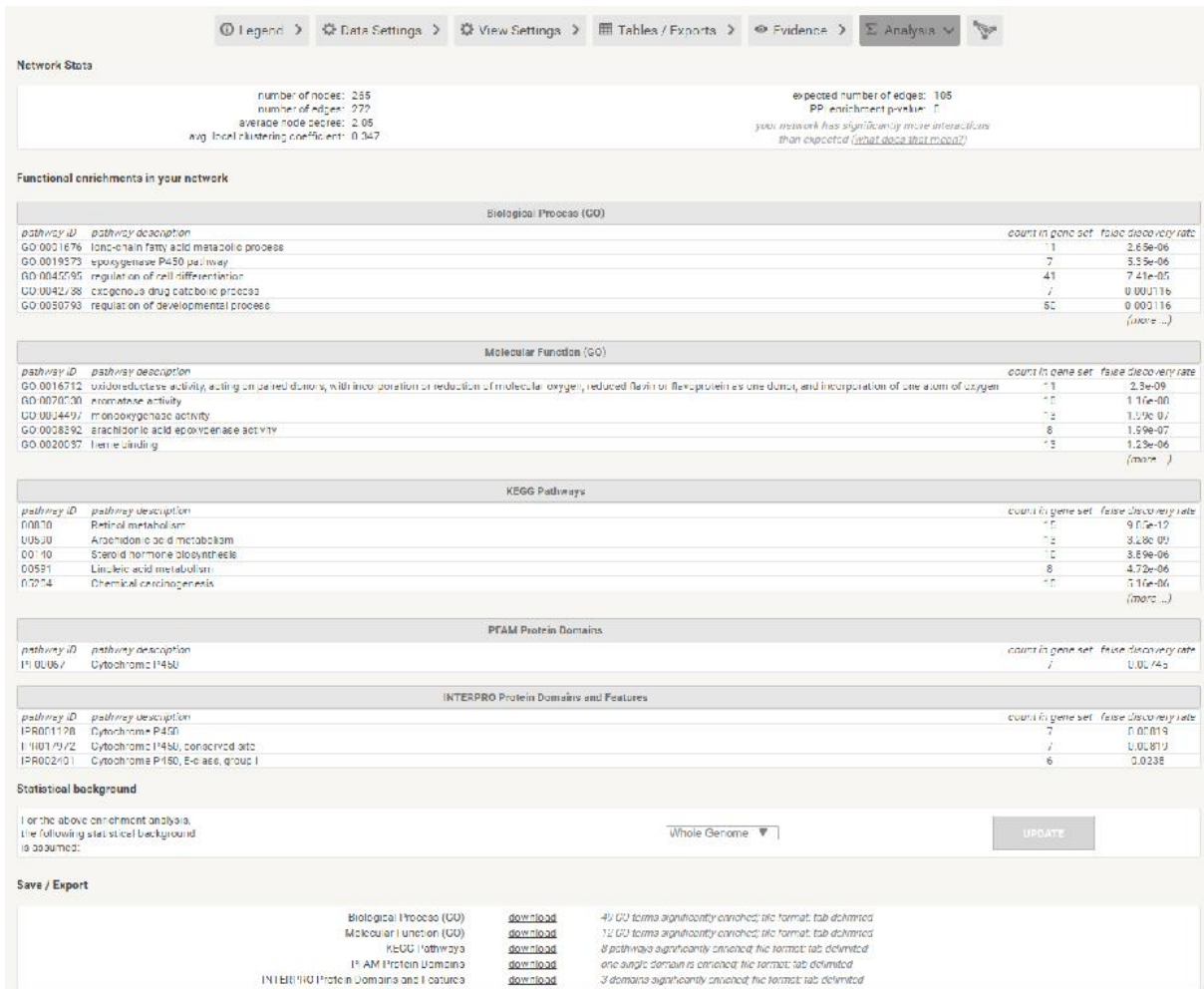


그림 2-2-5. Functional enrichments result

관심 있거나 중요한 Function을 클릭하면 Network상에서 해당 유전자들이 붉은색으로 표시되고 (그림 2-2-6) 관심 있는 유전자를 클릭하면 해당 유전자의 자세한 정보를 추가로 얻을 수 있다(그림 2-2-7).

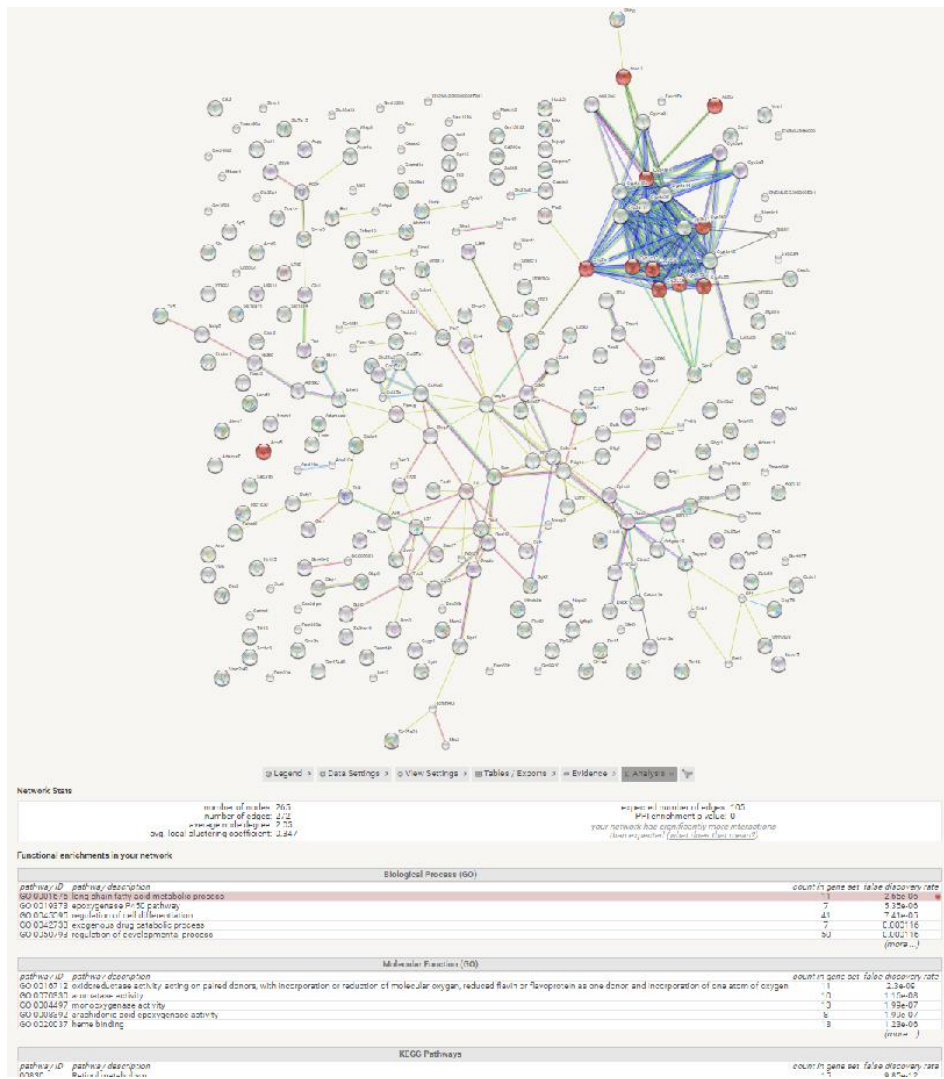


그림 2-2-6. Function selection on your network

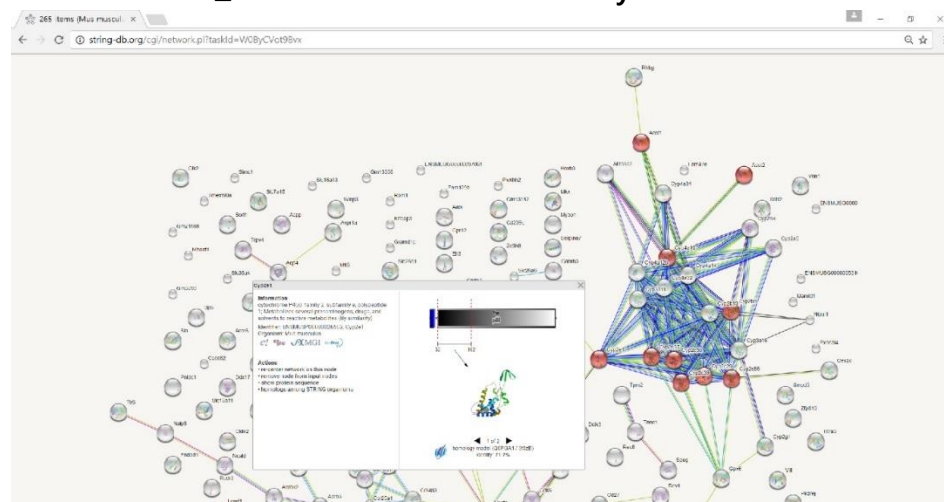


그림 2-2-7. Gene selection on your network

“Legend” 탭에서는 Node, Edge, Input 유전자의 설명을 자세히 볼 수 있고(그림 2-2-8) “Tables/Exports” 탭에서는 Network와 유전자 정보를 파일로 저장할 수 있다.(그림 2-2-9)

Legend | Data Settings | View Settings | Tables / Exports | Evidence | Analysis

Nodes:

- Network nodes represent proteins:** splice isoforms or post-translational modifications are collapsed, i.e. each node represents all the proteins produced by a single, protein-coding gene locus.
- Node Size:**
 - small nodes: protein of unknown 3D structure
 - large nodes: some 3D structure is known or predicted
- Node Color:**
 - colored nodes: query proteins and first shell of interactors
 - white nodes: second shell of interactors

Edges:

- Edges represent protein-protein associations:** associations are meant to be specific and meaningful, i.e. proteins jointly contribute to a shared function; this does not necessarily mean they are physically binding each other.
- Known Interactions:**
 - from curated databases
 - experimentally determined
- Predicted Interactions:**
 - gene neighbourhood
 - gene fusions
 - gene co-occurrence
- Others:**
 - text mining
 - on expression
 - protein homology

Your Input:

- Pparg:** peroxisome proliferator activator receptor gamma; Receptor that binds peroxisome proliferators such as hypolipidemic drugs and fatty acids. Once activated by a ligand, the receptor binds to a promoter element in the gene for acyl-CoA oxidase and activates its transcription. It therefore controls the peroxisomal beta oxidation pathway of fatty acids. Key regulator of adipocyte differentiation and glucose homeostasis. ARF6 acts as a key regulator of the tissue-specific adipocyte P2 (aP2) enhancer. Acts as a critical regulator of gut homeostasis by suppressing NF-kappa-B-mediated prolifia [..] (805 aa)
- Pdgfra:** platelet derived growth factor receptor, alpha polypeptide, Tyrosine-protein kinase that acts as a cell-surface receptor for PDGFA, PDGFB and PDGFC and plays an essential role in the regulation of embryonic development, cell proliferation, survival and chemotaxis. Depending on the context, promotes or inhibits cell proliferation and cell migration. Plays an important role in the differentiation of some marrow-derived mesenchymal stem cells. Required for normal skeleton development and cephalic closure during embryonic development. Required for normal development of the mucosa lining in [..] (1089 aa)
- Irf1:** interferon-related developmental regulator 1; Could play a role in regulating gene activity in the proliferative and/or differentiative pathways induced by NF- κ B. May be an oncogene factor that attenuates or amplifies the initial ligand induced signal (449 aa)
- Tmen25:** transmembrane protein 25 (365 aa)
- Rcc8:** REC8 homolog (yeast); Required during meiosis for separation of sister chromatids and homologous chromosomes. Proteolytic cleavage of REC8 on chromosome arms by separin during anaphase I allows for homologous chromosome separation in meiosis I and cleavage of REC8 on centromeres during anaphase II allows for sister chromatid separation in meiosis II (591 aa)
- Tiam1:** T cell lymphoma invasion and metastasis 1; Modulates the activity of WIPK-like proteinase and connects extracellular signals to cytoskeletal activities. Acts as a GTP-dissociation stimulator protein that stimulates the GDP-GTP exchange activity of RHO-like GTPases and activates them. Activates RAC1, RHOA2, and to a lesser extent RHOA (fly similarity). Affects invasiveness of T-lymphoma cells (1991 aa)
- Cyp2c29:** cytochrome P450, family 2, subfamily c, polypeptide 29; Metabolizes arachidonic acid to produce 14,15-epoxyoctadecatrienoic acid (EET) (490 aa)
- Cyp4x:** cytochrome P450, family 4, subfamily x, polypeptide 4 (271 aa)
- Cyp2b13:** cytochrome P450, family 2, subfamily b, polypeptide 13 (431 aa)
- Cyp2c5:** cytochrome P450, family 2, subfamily c, polypeptide 5 (431 aa)

그림 2-2-7. Legend of your network

Legend | Data Settings | View Settings | Tables / Exports | Evidence | Analysis

Export your current network:

- ... as a bitmap image: [download](#) file format is PNG; portable network graphic
- ... as a high-resolution bitmap: [download](#) some PNG format, but resolution at 400 dpi
- ... as a vector graphic: [download](#) SVG: scalable vector graphic - can be opened and edited in Illustrator, CorelDraw, Dia, etc
- ... as a simple tabular text output: [download](#) TSV: tab separated values - can be opened in Excel
- ... as an XML summary: [download](#) structured XML interaction data, according to the PSM-ML data standard
- ... network coordinates: [download](#) a flat-file format describing the coordinates and colors of nodes in the network
- ... protein sequences: [download](#) MFA: multi-fasta format - containing the amino acid sequences in the network
- ... protein annotations: [download](#) a tab-delimited file describing the names, domains and annotated functions of the network proteins

Browse Interactions In tabular form:

*node1	node2	node1 accession	node2 accession	node1 annotation	node2 annotation	score
Aasn1	Cyp4a10	FNSMJSP00000126448	FNSMJSP00000061126	acyl-CoA thioesterase 1; /acyl-CoA/...	cytochrome P450, family 4, subfam.	0.561
Aasn1	Cyp4a14	FNSMJSP00000126448	FNSMJSP0000030487	acyl-CoA thioesterase 1; /acyl-CoA/...	cytochrome P450, family 4, subfam.	0.475
Aasn1	Rhbg	FNSMJSP00000126448	FNSMJSP00000130767	acyl-CoA thioesterase 1; /acyl-CoA/...	Rhesus blood group-associated B...	0.473
Aasn2	Cyp4a10	FNSMJSP00000126448	FNSMJSP000001176	acyl-CoA thioesterase 2; /acyl-CoA/...	cytochrome P450, family 4, subfam.	0.433
Adrb3	Adrb3	FNSMJSP00000030162	FNSMJSP0000070445	adrenergic receptor, beta 3; Beta-a...	adrenergic receptor kinase, beta 2...	0.640
Adrb3	Me1r	FNSMJSP00000030162	FNSMJSP0000065929	adrenergic receptor, beta 3; Beta-a...	melanocortin 1 receptor; Receptor...	0.900
Adrb3	Pparg	FNSMJSP00000030162	FNSMJSP0000004050	adrenergic receptor, beta 3; Beta-a...	peroxisome proliferator activated r...	0.791
Adrb3	Slc2a4	FNSMJSP00000030162	FNSMJSP0000016710	adrenergic receptor, beta 3; Beta-a...	solute carrier family 2 (facilitated g...	0.493
Adrbk2	Adrb3	FNSMJSP00000070445	FNSMJSP000000162	adrenergic receptor kinase, beta 2...	adrenergic receptor, beta 3; Beta-a...	0.640
Adrbk2	Ncald	FNSMJSP00000070445	FNSMJSP0000007611	adrenergic receptor kinase, beta 2...	neurocalin delta; May be involved...	0.660
Alfh3a2	Cyp4a10	FNSMJSP00000073764	FNSMJSP00000061126	aldehyde dehydrogenase family 3...	cytochrome P450, family 4, subfam.	0.943
Alfh3a2	Cyp4a17b	FNSMJSP00000073764	FNSMJSP0000002487	aldehyde dehydrogenase family 3...	cytochrome P450, family 4, subfam.	0.916
Alfh3a2	Cyp4a14	FNSMJSP00000073764	FNSMJSP0000030487	aldehyde dehydrogenase family 3...	cytochrome P450, family 4, subfam.	0.947
Alfh3a2	Cyp4a32	FNSMJSP00000073764	FNSMJSP0000021369	aldehyde dehydrogenase family 3...	cytochrome P450, family 4, subfam.	0.904
Apol10a	Apo11f	FNSMJSP00000060650	FNSMJSP00000132565	apolipoprotein L 10A	apolipoprotein L 11a	0.900
Apol11a	Apo11f	FNSMJSP00000060650	FNSMJSP00000060650	apolipoprotein L 11a	apolipoprotein L 10A	0.900
Aqp4	Axpr1a	FNSMJSP00000078088	FNSMJSP00000070323	aquaporin 4; Forms a water-specifi...	arginine vasopressin receptor 1A; ...	0.430
Aqp4	Slc1a2	FNSMJSP00000078088	FNSMJSP00000079100	aquaporin 4; Forms a water-specifi...	solute carrier family 1 (glut high al...	0.630
Aqp4	Tpwy4	FNSMJSP00000078088	FNSMJSP00000071359	aquaporin 4; Forms a water-specifi...	transient receptor potential cation ...	0.910
Arhgap10	Rac3	FNSMJSP00000078658	FNSMJSP00000018156	Rho GTPase activating protein 10; ...	RAS-related G3 botulinum substra...	0.917

Server load: low (39%) | [Permalink](#) | [Share](#)

© STRING CONSORTIUM 2017 | [ABOUT](#) | [INFO](#) | [ACCESS](#) | [CREDITS](#)

그림 2-2-8. Tables/Exports of your network

2-3. MSigDB기반 GSEA 분석

GSEA 분석은 MSigDB 기반으로 유전자의 상관관계를 통계적으로 분석하여 입력한 유전자 세트의 주요 기능을 예측하고 각 유전자가 어떤 기능들에 포함되는지 overlap 분석을 제공해 준다. 분석과정은 그림 2-3-1 과 같다.

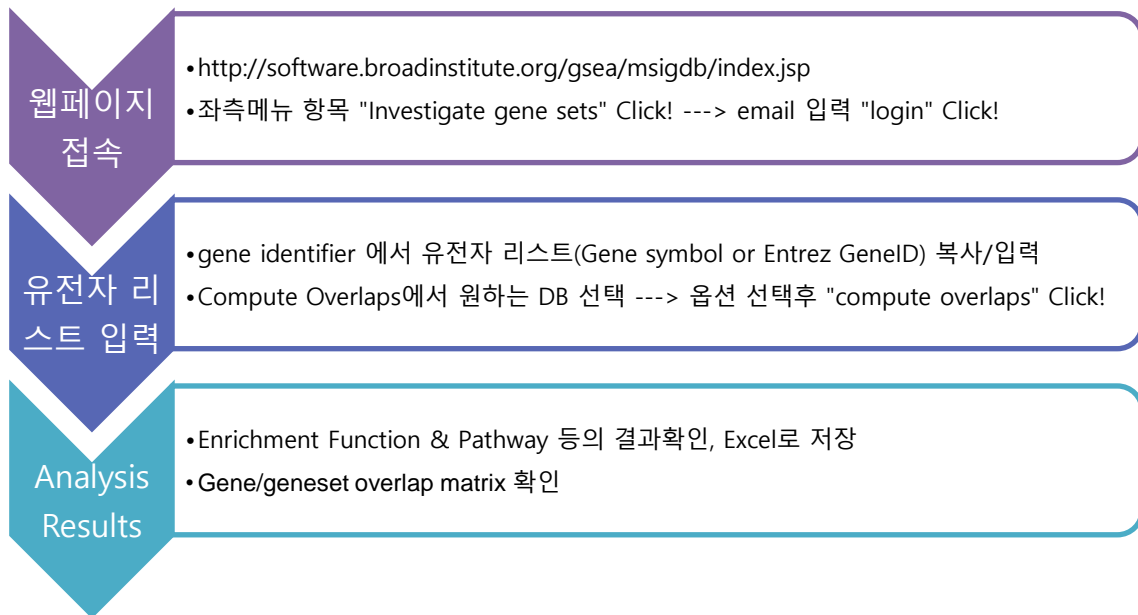


그림 2-3-1. Web based GSEA tool analysis process

MSigDB 에 접속하여 "Investigate gene sets"을 클릭하고 등록된 이메일을 입력하여 로그인을 수행한다.(그림 2-3-2).만약 등록이 필요할 시 "Click here"을 클릭하여 등록을 진행하면 된다.(그림 2-3-3).

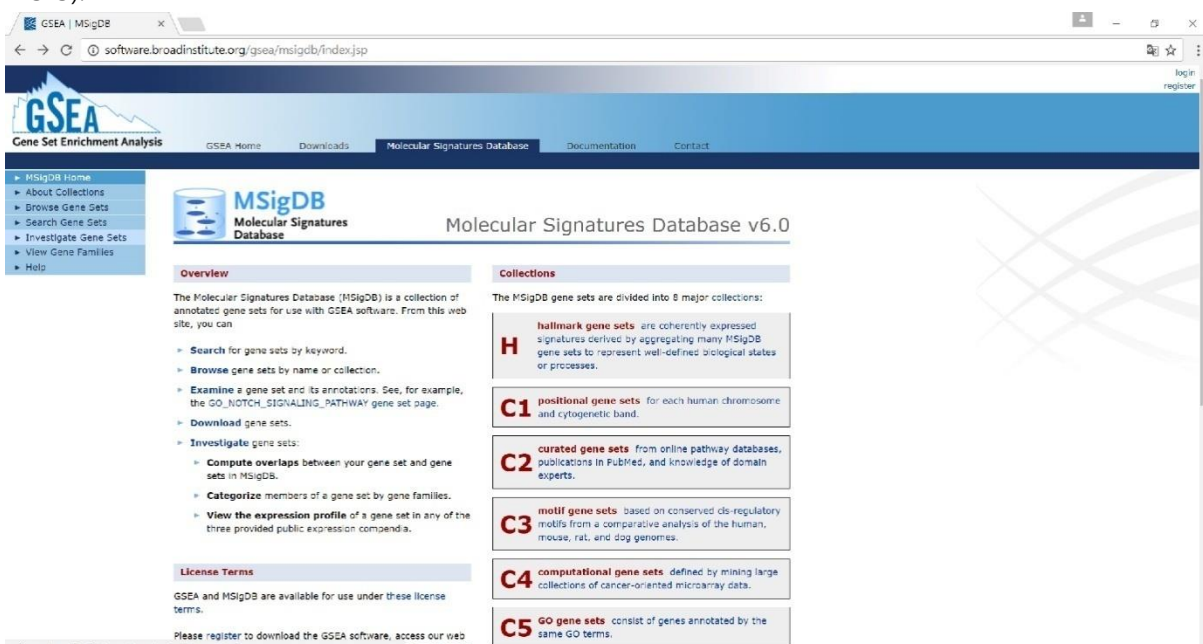


그림 2-3-2. GSEAmain page

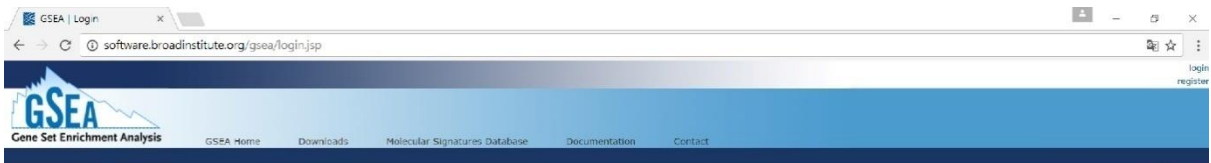


그림 2-3-3. GSEALogin page

"Gene Identifiers"입력창에 유전자 리스트(Gene Symbol, EntrezGeneID 또는 public ID)를 입력하고 "Compute Overlaps"에 원하는 DB 를 클릭한 후 맨 아래 "compute overlaps" 버튼을 클릭한다.(그림 2-3-4).DB 선택시 DB 명 앞의 파란색 글자를 누르면 해당 DB 정보를 확인할 수 있다.

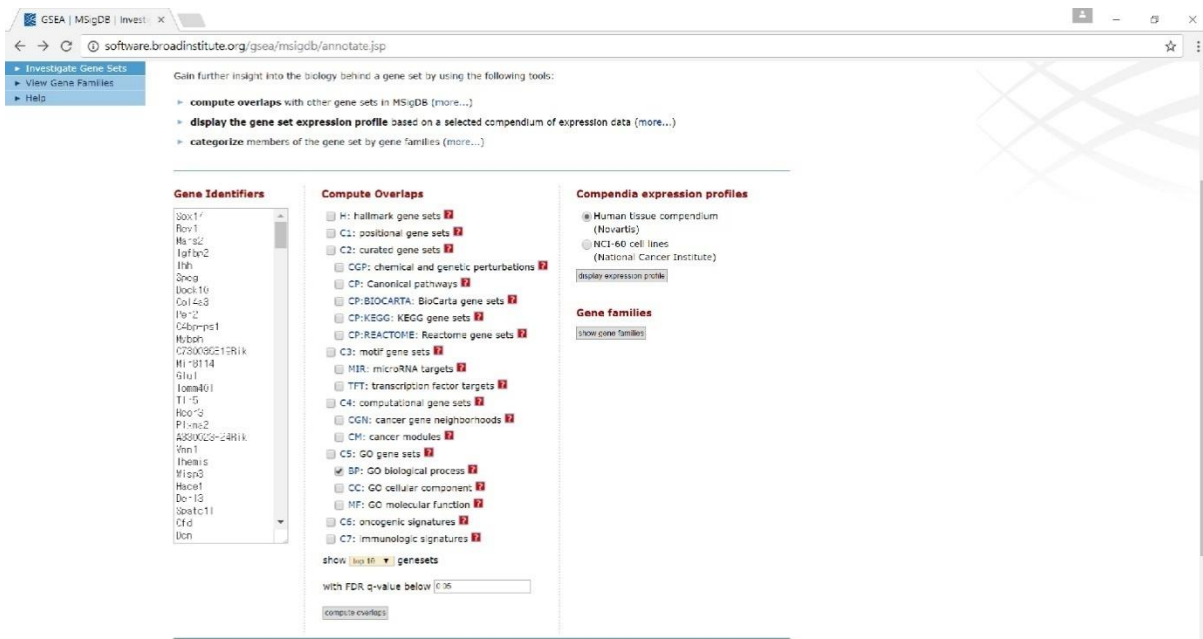


그림 2-3-4. GSEAAalysis

분석이 완료되면 그림 2-3-5 와 그림 2-3-6 과 같이 통계적으로 유의한 Gene Set List 와 Gene/Gene-set Overlap Matrix 결과를 확인할 수 있다.

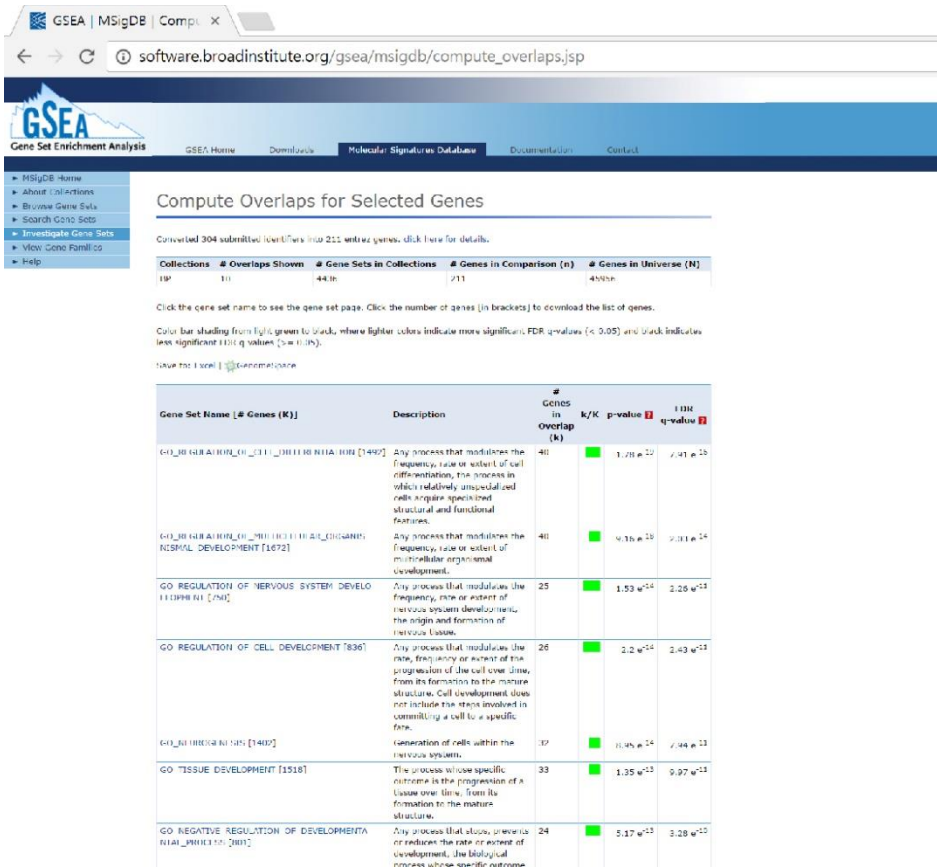


그림 2-3-5. GSEAAalysisResult(Gene Set)

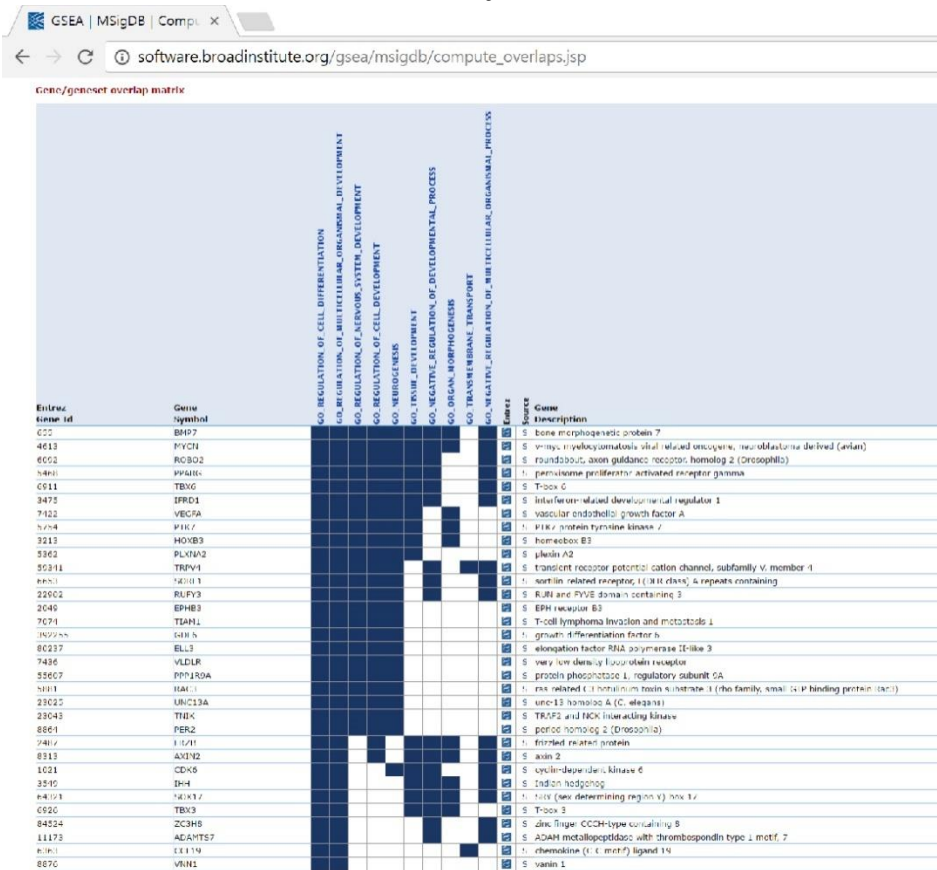


그림 2-3-6. GSEAAalysisResult(Gene/Gene-set Overlap Matrix)

3. KEGG DB 기반 Pathway 분석

mRNA-Seq 분석 결과에서 up/down-regulated 유전자들이 어떤 Pathway에 속하는지 확인하고자 한다면 KEGG에서 제공하는 KEGG Mapper를 이용하면 된다. 사용방법은 그림 3-1과 같은 순서로 진행된다.

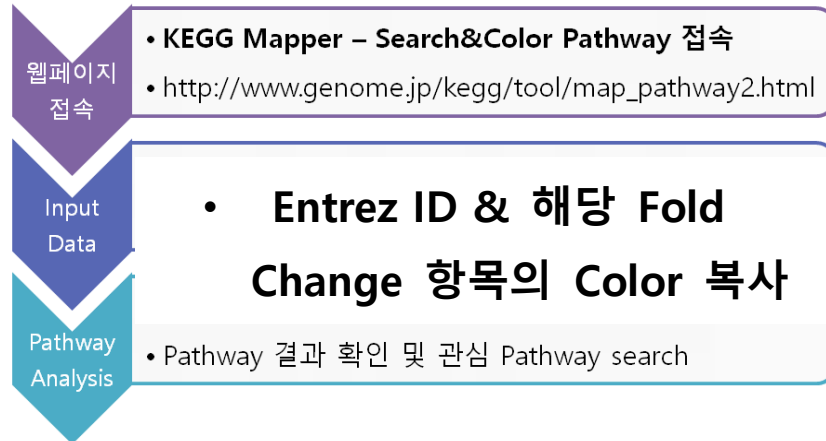


그림 3-1. KEGG Mapper tool analysis process

그림 3-2는 mRNA-Seq report에서 2fold, normalized RC(log2)>6을 기준으로 선별한 유전자를 KEGG 분석하는 과정이다.

*KEGG input 값은 excel 파일의 Annotation 항목 앞에 제작되어 있다.

오른쪽 필터에서 Fold change와 Normalized RC (반복실험의 경우 p-value) 값을 지정하고, 확인하고자 하는 Fold change 조합을 선택하여 필터를 적용 한다.

필터를 적용하여 선별 된 유전자의 KEGG input [Entrez ID, FC Color(#숫자,black)] cell을 복사하여, KEGG 분석에 사용할 것이다.

Gene symbol	B/A	C/A	Normalized RC (log2)			Raw data (RC)			KEGG mapper input		
			A	B	C	A	B	C	Entrez_ID	B/A	C/A
724 EIF2B3	2.021	1.579	9.681	8.706	10.560	936	405	1416	8891	#FF6347,black	#FFA07A,black
751 UQCRRH	3.086	1.929	11.872	10.406	13.375	4277	1318	9966	7388	#FF6347,black	#FFA07A,black
947 TYW3	2.051	1.946	10.192	9.584	11.211	1334	745	2224	127253	#FF6347,black	#FFA07A,black
966 NEXN-AS1	0.435	1.108	8.972	8.607	7.677	572	378	191	374987	#87CEEB,black	#87CEEB,black
1001 LOC646626	0.498	0.586	8.060	7.484	6.778	303	173	102	346626	#87CEEB,black	#B0E0E6,black
1023 GBP3	2.319	2.148	7.301	6.399	8.231	179	81	281	2635	#FF6347,black	#FF6347,black
1184 LAMTOR5	2.047	1.839	11.397	10.069	12.313	3076	1043	4775	10542	#FF6347,black	#FFA07A,black
1271 CD101	0.454	2.828	4.446	6.170	4.004	24	69	14	8398	#87CEEB,black	#FF6347,black
1331 TXNIP	2.506	3.261	6.363	5.797	7.349	93	53	152	10628	#FF6347,black	#FF6347,black
1403 MRPS21	2.122	1.345	9.250	8.035	10.133	694	254	1053	34460	#FF6347,black	#FFE4B3,black
1409 ADAMTSL4	0.481	0.591	10.394	10.734	8.768	1535	1655	408	34507	#87CEEB,black	#B0E0E6,black
1428 SEMA6C	0.440	0.343	9.007	9.265	7.518	586	597	171	10500	#87CEEB,black	#87CEEB,black
1436 PSMD4	2.018	0.591	12.673	12.047	13.476	7455	4111	10689	3710	#FF6347,black	#FF6347,black
1441 SELENBP1	0.168	2.248	6.289	5.086	3.102	88	32	7	8991	#00BFFF,black	#FF6347,black
1465 HRNR	0.391	0.935	10.512	10.985	9.284	1666	1969	584	88697	#87CEEB,black	#87CEEB,black
1550 SHE	0.386	1.612	6.004	5.289	4.553	72	37	21	126669	#87CEEB,black	#FFA07A,black
1570 EFN3	0.481	0.398	10.079	9.817	8.449	1233	876	327	1944	#87CEEB,black	#87CEEB,black
1575 TRIM46	0.430	0.437	8.955	9.185	7.081	565	565	126	80128	#87CEEB,black	#87CEEB,black
1632 CRABP2	0.379	1.222	11.041	10.557	9.249	2403	1463	570	1382	#87CEEB,black	#87CEEB,black

그림 3-2. KEGG Mapper tool analysis process

그림 3-3과 같이 KEGG Mapper 웹페이지(http://www.genome.jp/kegg/tool/map_pathway2.html)에 접속하고 Search & Color pathway 링크에 들어가면 아래와 같은 화면이 보여진다. 분석하고자 하는 유전자의 species를 선택하고, 'primary ID'는 KEGG identifiers로 선택한 뒤 'Enter objects one per line followed bgcolor, fgcolor' 창에 엑셀에서 준비해 놓은 Entrez ID, Color 항목을 복사-붙여넣기를 한다. 마지막으로 "Include aliases"와 "Use uncolored diagram" 항목에 체크를 한 후 Exec 버튼을 누른다.

KEGG Mapper - Search&Color Pathway

Search agains: Enter: map, ko, ec, rn, hsadd, or

Primary ID: Outside IDs for organism-specific pathways only

Enter objects one per line followed by bgcolor, fgcolor:

```
633 #87CEEB,black
105373383 #87CEEB,black
1852 #87CEEB,black
10134 #87CEEB,black
2157 #87CEEB,black
283981 #87CEEB,black
1438 #87CEEB,black
9189 #87CEEB,black
114758 #87CEEB,black
```

Alternatively, enter the file name containing the data:

Include aliases

Use uncolored diagrams

Display objects not found in the search

Search pathways containing all the objects (AND search)

Find organism - Chrome

Find three- or four-letter KEGG organism code

(human) Homo sapiens (human) [hsa]
 (human body louse) Pediculus humanus corporis (human body louse) [ph]

그림 3-3. KEGG Mapper tool analysis process

분석결과, 입력한 유전자들이 관여하는 pathway list가 나온다(그림 3-4). pathway 이름 옆에 있는 괄호 안 숫자는 입력한 유전자 중 각 pathway에 관여하는 유전자의 수이다. 괄호 안 숫자를 클릭하면 해당 유전자 목록을 볼 수 있다. pathway 이름을 클릭하면 해당 pathway chart가 열리고 입력한 유전자의 발현 up/down (red/green)이 색으로 표시되어 있다. Pathway 이미지는 "다른 이름으로 저장"이 가능하고 "html"으로 저장하면 이미지에 링크된 항목을 그대로 유지해서 저장이 가능하다.

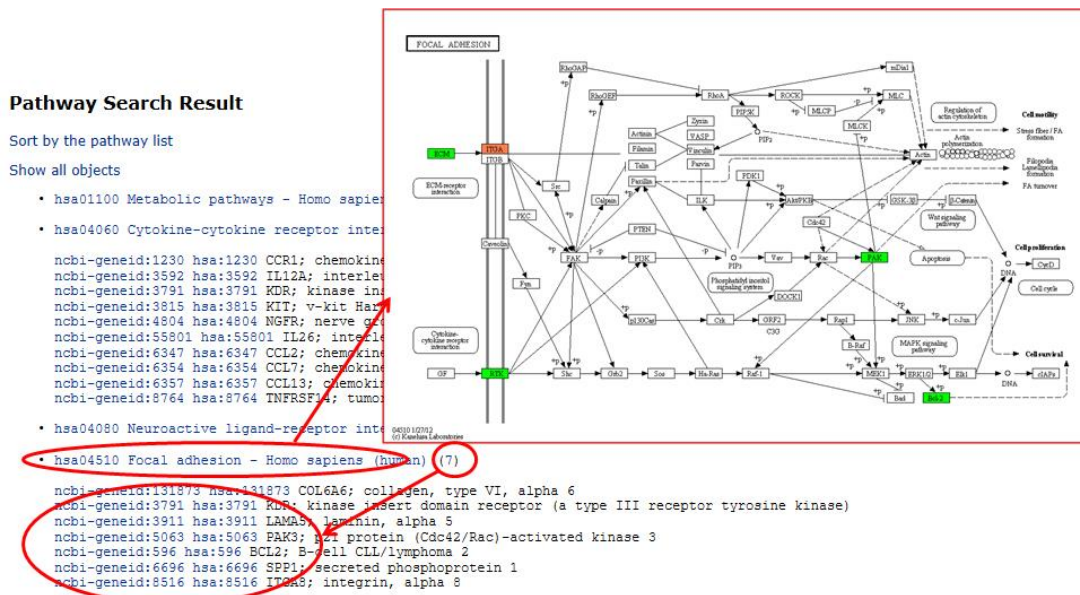


그림 3-4. KEGG Mapper tool analysis result

4. MeV Software 이용 Clustering Heatmap 작성

MeV 소프트웨어는 미국의 Dana-Farber Cancer Institute에서 개발한 Microarray, mRNA-Seq 전용 분석 프로그램으로 연구자들에게 무료로 공급하고 있다. 주로 clustering 분석과 통계분석(K-means clustering, Hierarchical clustering, t-test, Significance Analysis of mRNA-Seqs, Gene Set Enrichment Analysis, EASE)을 할 수 있는 프로그램이다. 아래 웹페이지에 접속하면 최신의 업데이트된 프로그램과 매뉴얼을 다운받을 수 있다.

<http://www.tm4.org> >> 오른쪽 Browse 항목내 "TM4 MeV Stand-Alone Client" 클릭
 프로그램을 다운받아 압축을 풀고, MeV 또는 TMEV를 클릭해서 프로그램을 실행시킨다(그림4-1).MEV프로그램을 실행시키면 세 개의 창이 나타난다(그림4-2). 분석창은 프로그램창의 메뉴에서 file->New multiple array viewer를 통해 여러개를 생성할 수 있고 데이터 분석은 분석창을 통해 진행한다.

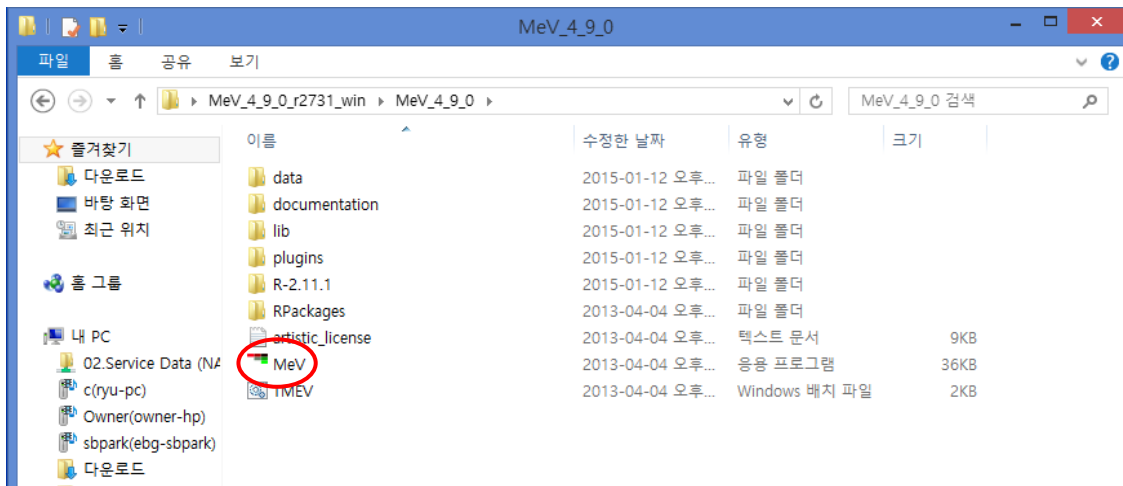


그림 4-1. MeV program folder and files

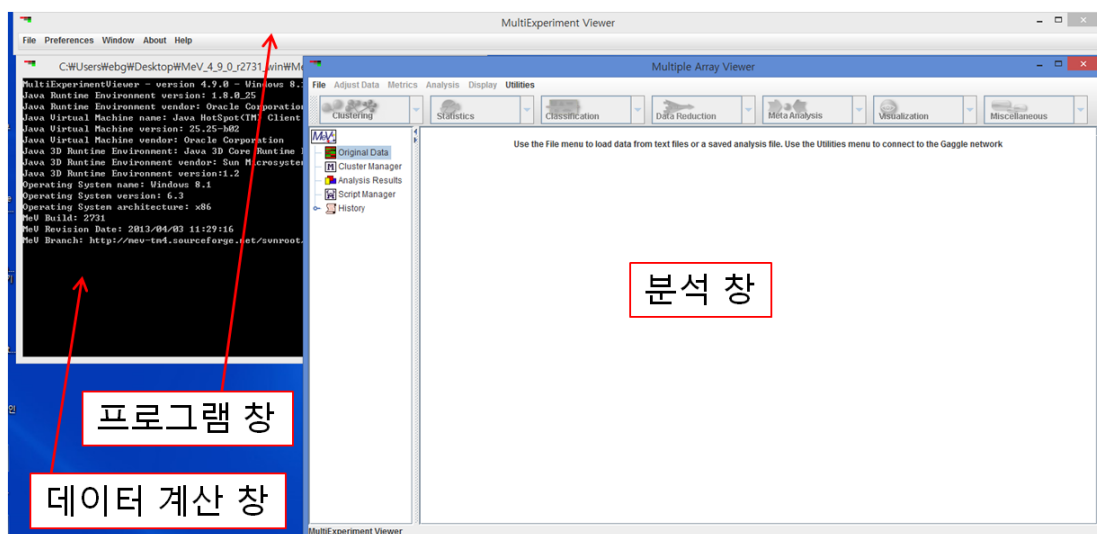


그림 4-2. MeV program windows

본 자료에서는 MeV 프로그램을 이용하여 Clustering 분석 방법을 설명한다. 우선 MeV 프로그램에 input할 데이터를 엑셀에서 파일 양식에 맞춰 저장해야 한다. 엑셀에 clustering 하고자 하는 유전자 이름과 fold change 또는 발현값(intensity)를 정리한다(그림 4-3). 그리고 '텍스트 (탭으로 분리)'파일 형식으로 저장해야 MeV에 upload 할 수 있다. MeV에서는 2만 개 이상의 유전자는 clustering 분석을 할 수 없으므로 2만 개 이하로 유전자를 선별해야 한다.

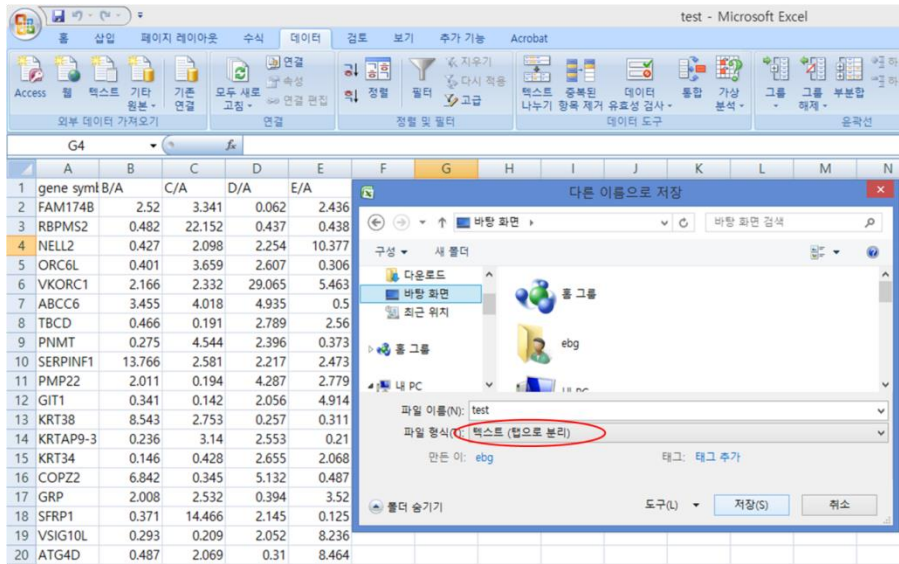


그림 4-3. Data format example

input 데이터 저장이 완료되면 MeV 프로그램의 분석창에서 file -> load data를 실행한다(그림 4-4). Browse를 클릭하여 input 데이터를 선택한다. 데이터가 fold change인 경우 "Two-color Array"로 체크하고 데이터가 intensity인 경우는 "Single-color Array"에 체크한다. 마우스로 데이터가 시작되는 부위를 클릭한 후 load를 누른다.

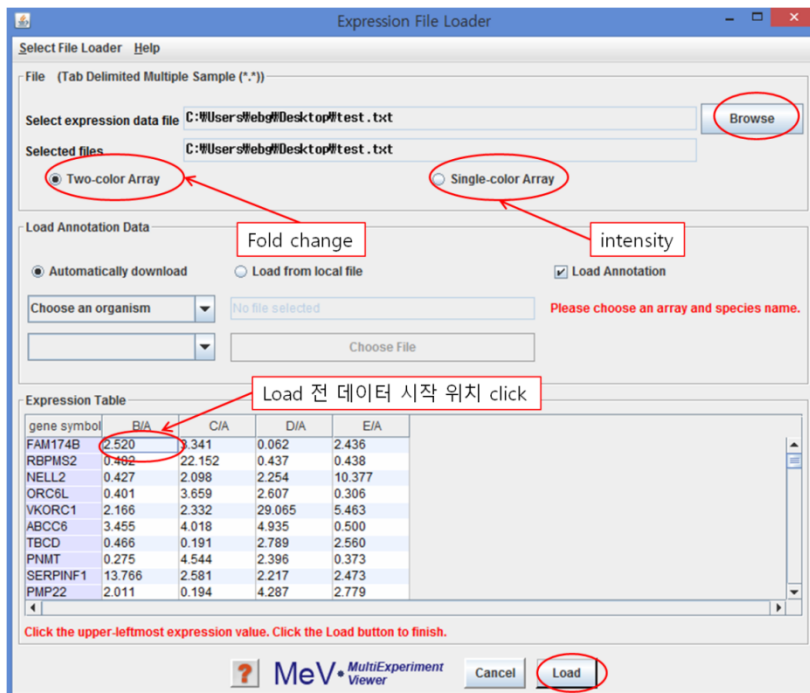


그림 4-4. Data uploading method

데이터가 열리면 Adjust Data -> Log Transformation -> Log2 Transform을 선택하여 fold change는 $\log_2(\text{fold change})$ 로, intensity는 $\log_2(\text{intensity})$ 로 바꿔준다(그림 4-5). 왼쪽 메뉴의 Original Data -> Expression image를 보면 \log_2 값으로 바뀌어 색이 변한 것을 확인할 수 있다.

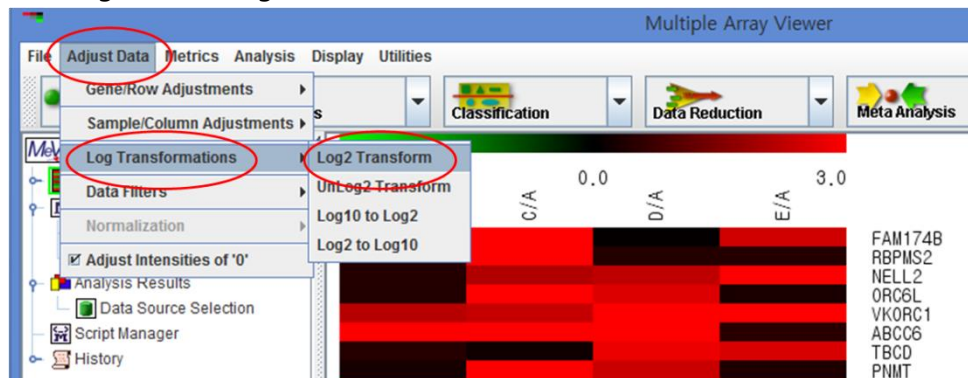


그림 4-5. Log2 transformation

Analysis-> Clustering-> HCL을 선택하여 Clustering 분석을 시작한다(그림 4-6).

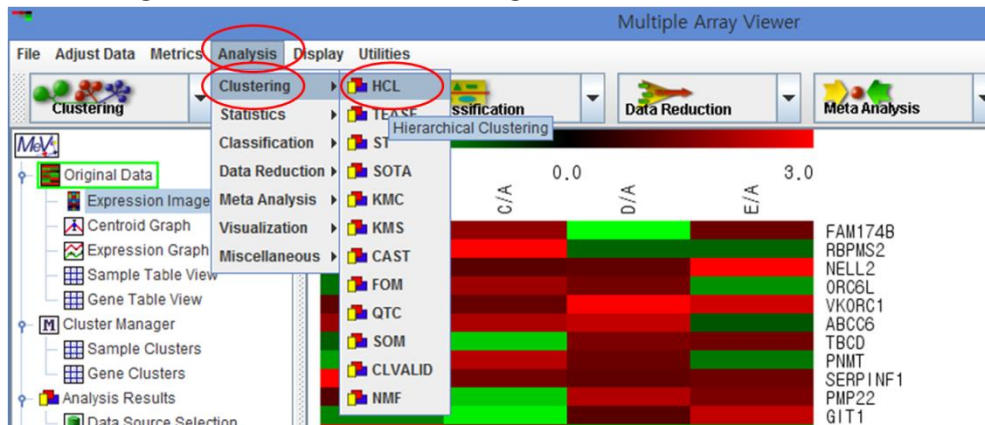


그림 4-6. Hierarchical Clustering Selection

Clustering 분석 시 다양한 옵션을 선택할 수 있다(그림 4-7). Gene tree를 선택하면 fold change 또는 intensity가 유사한 유전자끼리 clustering한 결과가 나온다. Sample tree를 선택하면 발현이 유사한 샘플끼리 clustering한 결과가 나온다.당사에서 clustering 분석을 할 때 Distance Metric는 Euclidean Distance로 Linkage Method Selection은 Average linkage clustering으로 설정한다. 다른 옵션을 선택해도 된다. 옵션을 선택하고 OK를 누른다.

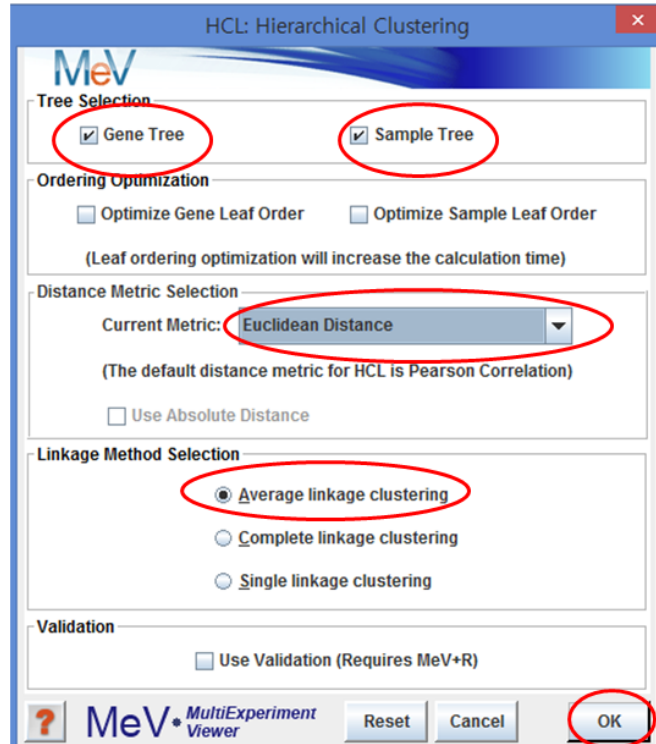


그림 4-7. Hierarchical Clustering Method

clustering이 완료되면 왼쪽 메뉴에 Analysis Results에 HCL 결과가 생긴다. HCL -> HCL tree를 클릭하면 clustering 결과가 화면에 나온다(그림 4-8). 위의 tree는 sample clustering 결과이고 왼쪽 tree는 gene clustering 결과이다. 각 tree에는 distance scale bar가 있어서 tree의 길이를 가늠할 수 있다. tree의 길이는 distance이며, distance가 짧을수록 유전자 간 또는 샘플 간의 발현이 비슷한 것, 길수록 발현이 다른 것이다.

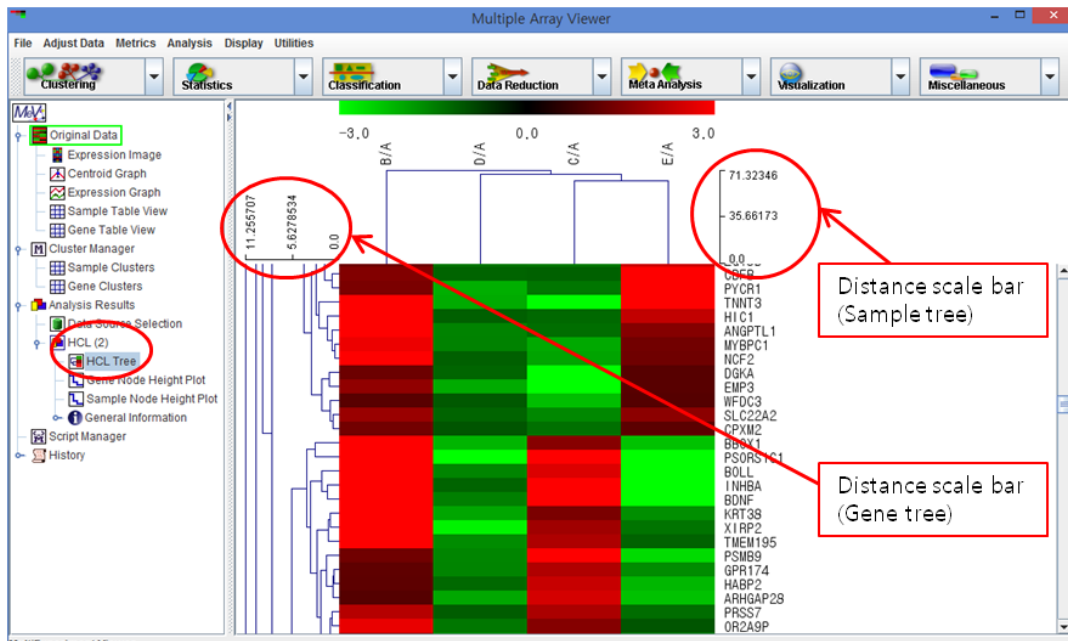


그림 4-8. Hierarchical Clustering Result

clustering 결과는 이미지의 크기와 색상을 조절하여 원하는 형태의 이미지를 만들 수 있다(그림 4-9, 4-10)

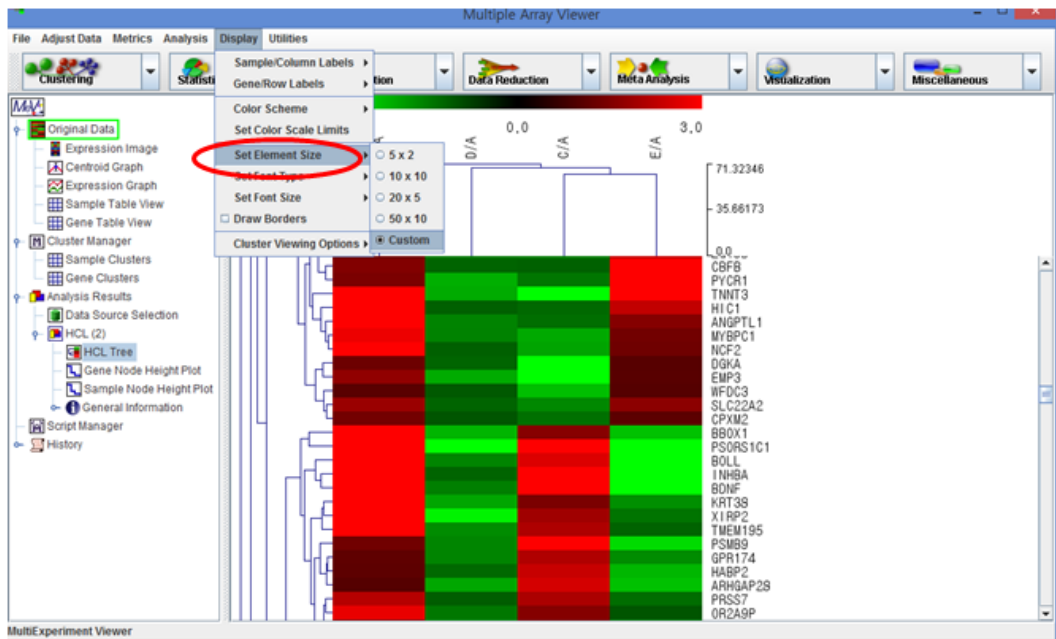


그림 4-9. Clustering image size control

Display -> Set Color Scale Limits를 누르면 color scale bar의 최소값, 중간값, 최대값을 설정할 수 있다. 보통 $\log_2(\text{fold change})$ 는 최소값과 최대값은 같은 크기에 부등호만 바꿔주고(예: min:-3, max:3) 중간값은 0으로 설정해 준다(그림 4-10). 이렇게 하면 up-regulated genes은 red, down-regulated genes은 green으로 나타나게 된다.

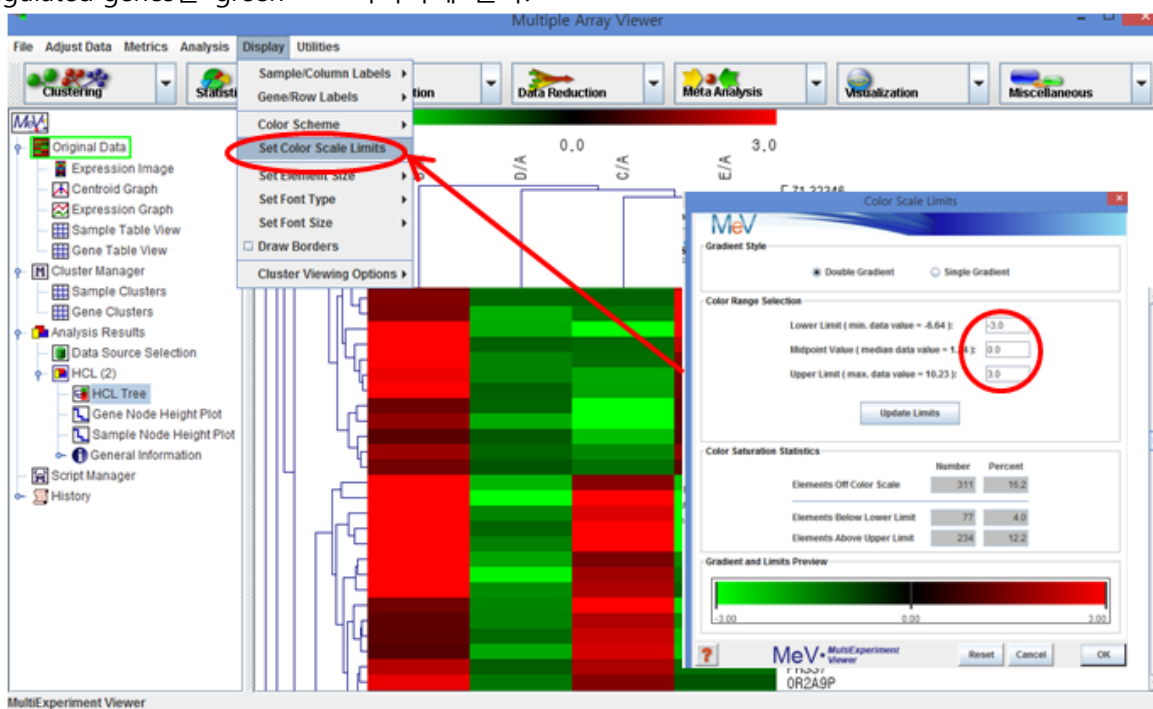


그림 4-10. Clustering image color setting

원하는 이미지 조절이 완료되면 File -> Save image를 눌러 이미지를 저장한다. 이때 파일 이름에 파일 확장자명(예: .jpg)을 꼭 기입하여야 이미지 파일로 저장이 된다(그림 4-11).

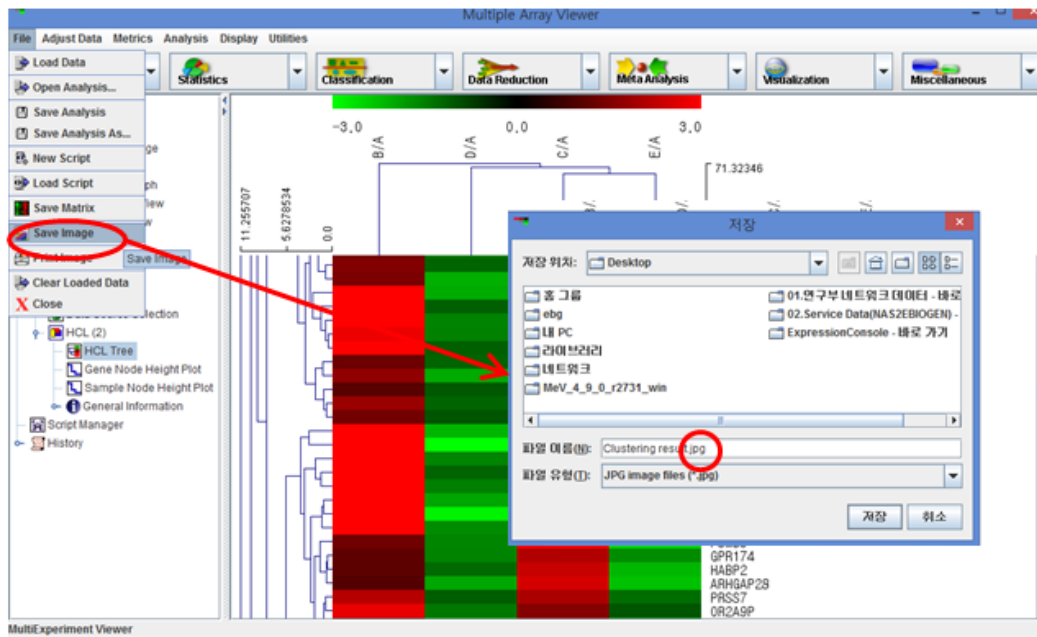


그림 4-11. Clustering image save