

# User Manual

## ExDEGA v.5.0.0 & Data Analysis

## < 목 차 >

1.	Differentially Expressed Genes Analysis (ExDEGA)	3
2.	Functional Annotation Analysis (DAVID, ExDEGA GraphicPlus)	22
3.	Clustering Heatmap Analysis (Clustering, ExDEGA GraphicPlus)	34
4.	Principal Component Analysis (PCA, ExDEGA GraphicPlus)	39
5.	String Network Analysis (Network, ExDEGA GraphicPlus)	48
6.	Pathway Analysis (KEGG mapper)	52
7.	Gene Set Enrichment Analysis (GSEA)	55
8.	Protein-Protein Network Analysis (Cytoscape STRING)	63

# 1. Differentially Expressed Genes Analysis (ExDEGA)

(주)이바이오젠은 RNA-Seq (Quant-Seq, mRNA-Seq, Total RNA-seq)과 Microarray data 를 엑셀 기반에서 쉽게 분석할 수 있도록 분석 결과 보고 시 ExDEGA (Excel based Differentially Expressed Gene Analysis) tool 과 ExDEGA Graphic Plus 를 함께 제공한다. ExDEGA 분석 툴은 (주)이바이오젠이 연구자들이 Microarray 및 RNA-Seq 데이터를 보다 쉽게 다루고 원하는 데이터를 쉽게 얻을 수 있도록 사용자 편의를 최대한 반영한 분석 툴이고 엑셀 프로그램 안에서 다양한 분석을 직관적으로 수행할 수 있도록 개발되었다. ExDEGA 분석 툴은 사용자들의 요구사항을 지속적으로 반영하여 데이터 분석과 엑셀 사용에 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\_v(버전)\_Installer.zip 파일을 다운로드 폴더에서 압축을 풀고, setup.exe 를 실행하면 분석 툴이 설치된다(그림 1-1 A). 만약 다운로드 폴더에서 압축을 풀지 않았을 경우, 별도로 압축 킷 파일에 있는 ExDEGAGraphicPlus.exe 를 컴퓨터의 로컬 C 드라이브 아래로 복사+붙여넣기 하면 ExDEGA Graphic Plus 프로그램이 설치 완료된다(그림 1-1 B).

설치가 완료되고 ExDEGA format 의 엑셀 데이터를 열면 자동으로 ExDEGA 분석 툴이 구동된다. 참고로 ExDEGA 설치 전에 실행 중인 엑셀 파일이 있으면 종료시킨 후 다시 실행해야 ExDEGA 를 사용할 수 있다. ExDEGA 설치 및 구동에 오류가 있으면 ExDEGA 오류 해결 메뉴얼 ([Download link](#))을 확인한다.

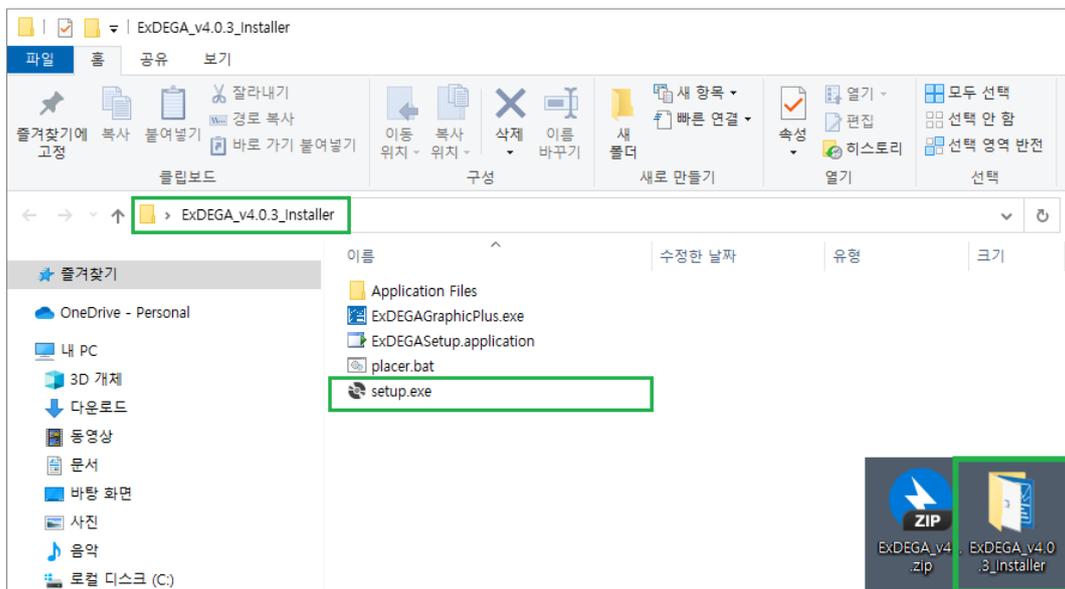


그림 1-1 A. ExDEGA set up

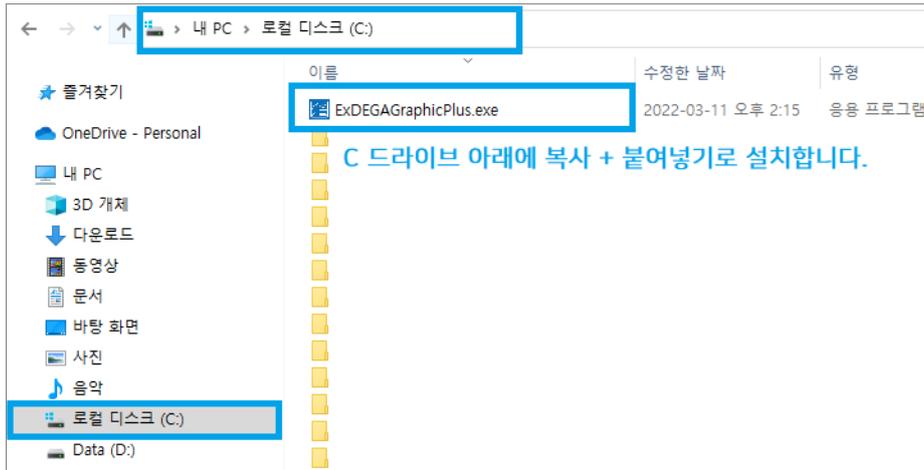


그림 1-1 B. ExDEGA Graphic Plus installation

ExDEGA format 의 엑셀 파일을 열면, 왼쪽에 Gene Category 창과 가운데에 gene expression data, 오른쪽에 DEG Analysis 창이 실행된다(그림 1-2). Gene Category 분석 창에서는 기본 설정된 Gene ontology (GO)가 있고 사용자가 원하는 대로 gene category 를 구성하여 분석할 수 있다. Gene category 창과 DEG Analysis 창은 함께 연동하여 데이터를 쉽게 얻을 수 있다. DEG Analysis 창에서는 Fold change, Normalized Data (log2), p-value 등을 선택하여 DEG 선별을 쉽게 할 수 있고 DEGs 를 gene category 별로 그래프를 작성할 수 있다. 뿐만 아니라, DEG 분석 창에서 Scatter Plot, Volcano Plot, Venn Diagram 을 직접 그릴 수 있고 선별된 유전자들을 대상으로 Clustering heatmap, KEGG 분석, DAVID 분석을 수행하기 위한 input file 을 자동으로 만들 수 있다. Gene expression graph, Gene search 기능도 이용할 수 있어 연구자가 RNA-Seq, microarray data 를 쉽게 활용할 수 있다.

Gene symbols	Fold change			p-value			Average of normalized data (log2)			Normalized data (log2)										
	B/A	C/A	C/B	B/A	C/A	C/B	A	B	C	A1	A2	A3	B1	B2	B3	C1	C2	C3	A4	
1 A1B5	0.849	0.990	0.885	0.523	0.018	0.044	1.408	1.371	0.845	1.544	1.989	1.865	1.631	1.149	1.290	0.843	0.743	0.843	1.044	0.843
2 A1B5AS1	1.049	1.907	1.487	0.742	0.014	0.013	0.954	0.422	0.946	0.629	0.196	0.221	0.514	0.562	0.159	0.944	0.844	1.044	0.629	0.196
3 A1B5P	0.988	1.061	1.079	0.373	0.238	0.163	0.917	0.900	0.101	0.900	0.900	0.900	0.900	0.900	0.900	0.900	0.900	0.900	0.900	0.900
4 ADM	1.050	1.077	1.036	0.759	0.288	0.064	0.233	0.293	0.330	0.231	0.330	0.100	0.642	0.091	0.071	0.338	0.238	0.428	0.231	0.330
5 A3M-AS1	1.659	1.405	0.871	0.290	0.062	0.089	0.034	1.122	0.925	0.300	0.799	0.111	1.771	0.917	0.294	0.933	0.823	1.023	0.300	0.799
6 A3M1	1.059	1.058	1.088	0.463	0.387	0.447	0.011	0.048	0.102	0.000	0.000	0.001	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.000
7 ADMP1	0.953	0.999	1.049	0.549	0.995	0.365	0.103	0.054	0.102	0.000	0.290	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
8 AIGAT2	1.187	0.959	0.845	0.683	0.794	0.510	0.181	0.348	0.102	0.440	0.002	0.007	0.001	0.001	0.001	0.842	0.380	0.600	0.440	0.002
9 AIGAT	0.824	0.760	0.748	0.203	0.000	0.000	0.280	0.001	1.744	0.218	0.564	0.000	0.001	0.001	0.001	0.000	1.742	1.642	1.842	0.218
10 AIGNT	0.937	1.005	1.071	0.373	0.946	0.162	0.094	0.000	0.102	0.000	0.000	0.286	0.000	0.000	0.000	0.000	0.100	0.000	0.200	0.000
11 AINS	0.902	0.968	1.071	0.374	0.779	0.162	0.148	0.000	0.102	0.406	0.001	0.000	0.000	0.000	0.000	0.000	0.100	0.000	0.200	0.406
12 AINS2	1.198	0.749	0.626	0.275	0.061	0.028	2.892	1.153	2.476	0.136	2.775	2.732	2.813	3.270	3.523	2.474	2.374	2.574	2.174	1.136
13 AINS3	0.947	2.000	2.046	0.686	0.002	0.000	1.209	1.142	2.243	1.510	1.200	0.904	1.276	1.172	0.962	2.241	2.341	2.941	1.200	1.510
14 AINS4	1.062	1.079	1.051	0.139	0.163	0.819	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.100	0.000	0.200	0.000
15 AINS5	1.000	1.073	1.073	1.000	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.100	0.000	0.200	0.000
16 AINS6	1.000	1.073	1.073	1.000	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.100	0.000	0.200	0.000
17 AINS7	1.000	1.073	1.073	1.000	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.100	0.000	0.200	0.000
18 AINS8	1.000	1.073	1.073	1.000	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.100	0.000	0.200	0.000
19 AINS9	1.000	1.073	1.073	1.000	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.100	0.000	0.200	0.000
20 AINS10	1.000	1.073	1.073	1.000	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.100	0.000	0.200	0.000
21 AINS11	1.000	1.073	1.073	1.000	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.100	0.000	0.200	0.000
22 AINS12	1.000	1.073	1.073	1.000	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.100	0.000	0.200	0.000
23 AINS13	1.000	1.073	1.073	1.000	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.100	0.000	0.200	0.000
24 AINS14	1.012	2.953	3.881	0.387	0.000	0.000	0.001	0.018	0.891	0.001	0.001	0.000	0.000	0.000	0.000	0.000	0.888	0.788	0.888	0.001
25 AINS15	0.980	0.906	0.963	0.603	0.001	0.725	2.405	1.621	1.847	2.938	3.951	2.411	1.899	1.790	1.399	1.895	1.465	1.666	2.938	2.938
26 AINS16	0.817	1.235	1.532	0.037	0.035	0.001	2.764	2.472	3.069	2.906	2.748	2.627	2.521	2.488	2.404	3.096	2.966	3.166	2.906	2.906
27 AINS17	0.918	0.921	0.897	0.688	0.208	0.349	2.305	1.979	1.811	1.758	2.341	2.341	1.814	1.683	1.839	1.719	1.919	1.919	1.719	1.919
28 AINS18	1.069	0.786	0.901	0.186	0.052	0.757	1.907	2.951	2.205	1.803	1.038	0.997	1.812	2.062	1.999	2.201	2.101	2.301	1.803	1.803
29 AINS19	0.793	0.932	0.724	0.502	0.222	0.080	4.176	3.783	3.320	4.868	3.540	3.742	3.489	3.781	4.049	3.317	3.217	3.417	3.317	3.417
30 AINS20	1.029	1.055	1.022	0.789	0.162	0.170	0.179	0.208	0.110	0.199	0.206	0.000	0.000	0.000	0.000	0.000	0.200	0.100	0.300	0.000
31 AINS21	0.741	0.938	0.930	0.080	0.003	0.013	2.875	2.442	1.885	3.039	2.943	2.609	2.209	2.408	2.873	1.843	1.783	1.983	1.843	1.783
32 AINS22	0.814	0.601	0.737	0.370	0.238	0.061	0.837	0.543	0.102	0.255	0.442	1.489	0.519	0.292	0.778	0.100	0.000	0.200	0.100	0.000
33 AINS23	1.910	1.910	1.910	1.602	0.001	0.000	2.139	2.799	2.487	1.907	1.403	1.894	2.112	0.911	1.928	1.403	1.945	1.945	1.403	1.945
34 AINS24	1.306	1.530	1.157	0.084	0.012	0.131	0.674	1.059	1.269	0.795	0.826	0.356	1.145	1.164	0.846	1.267	1.167	1.367	0.795	0.795
35 AINS25	0.889	0.810	0.702	0.407	0.000	0.144	2.354	2.351	1.840	2.539	2.307	2.394	1.789	2.340	2.591	1.898	1.738	1.938	1.898	1.938
36 AINS26	1.202	1.322	1.260	0.006	0.007	0.008	1.144	1.450	1.147	1.922	1.938	1.211	1.427	1.387	1.430	1.845	1.445	1.645	1.922	1.922
37 AINS27	1.118	1.189	1.046	0.354	0.215	0.461	2.479	2.640	2.704	2.743	2.407	2.341	2.619	2.738	2.556	2.702	2.602	2.802	2.407	2.407
38 AINS28	1.040	1.051	0.962	0.849	0.097	0.734	0.700	0.768	0.701	0.319	0.611	1.113	0.648	0.577	1.008	0.889	0.939	0.798	0.319	0.319
39 AINS29	0.912	0.668	0.730	0.802	0.357	0.050	1.074	0.941	0.492	1.781	0.631	0.409	1.157	0.911	0.698	0.490	0.900	0.590	1.781	1.781
40 AINS30	0.904	1.160	1.284	0.541	0.072	0.149	4.763	4.617	4.978	4.890	4.713	4.676	4.393	4.601	4.954	4.975	4.875	5.075	4.617	4.617
41 AINS31	1.848	1.848	1.848	1.891	0.068	0.001	0.103	2.910	3.228	3.645	1.843	2.631	2.451	3.616	1.132	1.824	1.642	1.842	3.645	3.645
42 AINS32	1.000	1.073	1.073	1.000	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.100	0.000	0.200	0.000
43 AINS33	0.879	0.946	0.821	0.642	0.007	0.124	7.493	7.250	6.983	7.522	6.878	7.366	7.402	6.370	7.689	6.561	6.461	6.661	6.878	6.878
44 AINS34	1.426	1.432	1.380	0.012	0.000	0.007	1.482	1.991	2.036	1.414	1.333	1.200	2.008	1.996	1.930	2.004	2.004	2.004	1.414	1.414

그림 1-2. RNA-seq or Microarray data in ExDEGA format

### 1-1. Gene Category 사용 방법

RNA-seq 또는 Microarray 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 분석된 종에 대한 GO 가 임의로 추가되어 있으며, 관련 유전자들을 필터링할 수 있다. 예를 들어, Aging 관련 유전자만 분석을 원할 경우, Gene Category 창에서 Aging 을 선택하면 해당 유전자 리스트만 필터링 된다. 그리고 Gene Category 의 여러 GO 들을 선택하여 동시에 해당하는 유전자 (교집합)를 필터링할 수 있는 "AND" 기능과 한 GO 만이라도 해당하는 유전자 (합집합)를 필터링할 수 있는 "OR" 기능을 갖추고 있다. (그림 1-3)

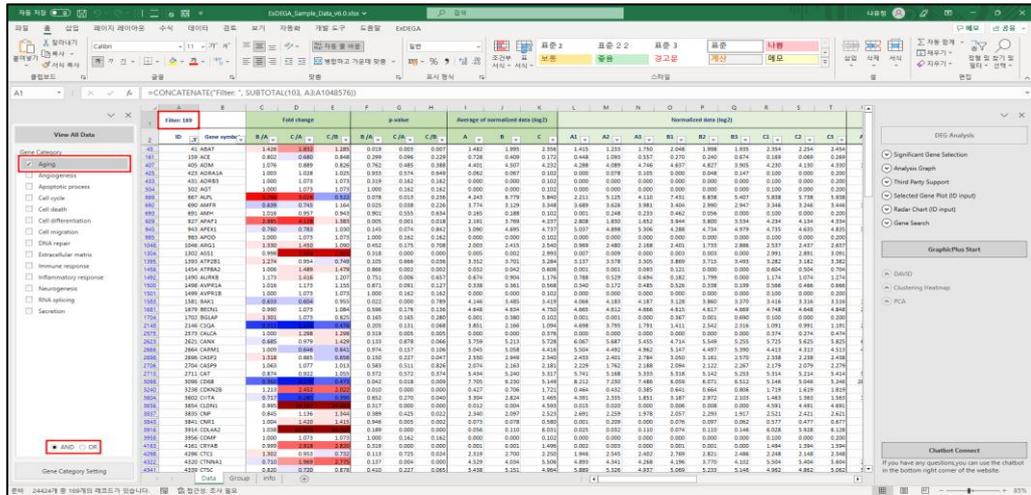


그림 1-3. Gene ontology selection

가장 왼쪽 상단에 'View All Data' 버튼을 누르면 필터를 모두 해제되어 다시 전체 결과를 볼 수 있고 기존 GO 중 관심있는 GO 가 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.

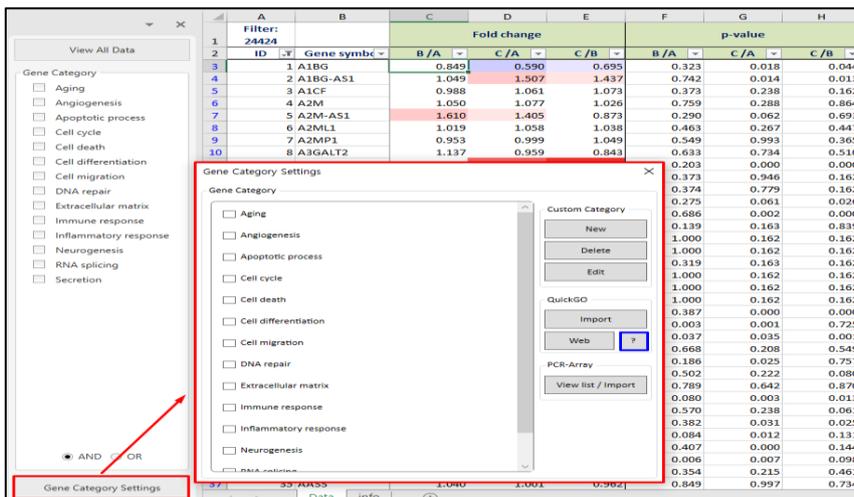


그림 1-4. Gene category settings

만약 원하는 유전자 그룹 목록을 알고 있다면, 직접 입력하여 새로운 Gene Category 를 추가할 수도 있다. Gene Category Settings 버튼을 누른 후 New 를 선택하고 원하는 Gene symbol list 입력(or 복사&붙여넣기) 한 뒤, Gene category 이름 설정 후 Setting 창을 닫아주면 새로운 Gene category 를 확인할 수 있다(그림 1-5. a, b).

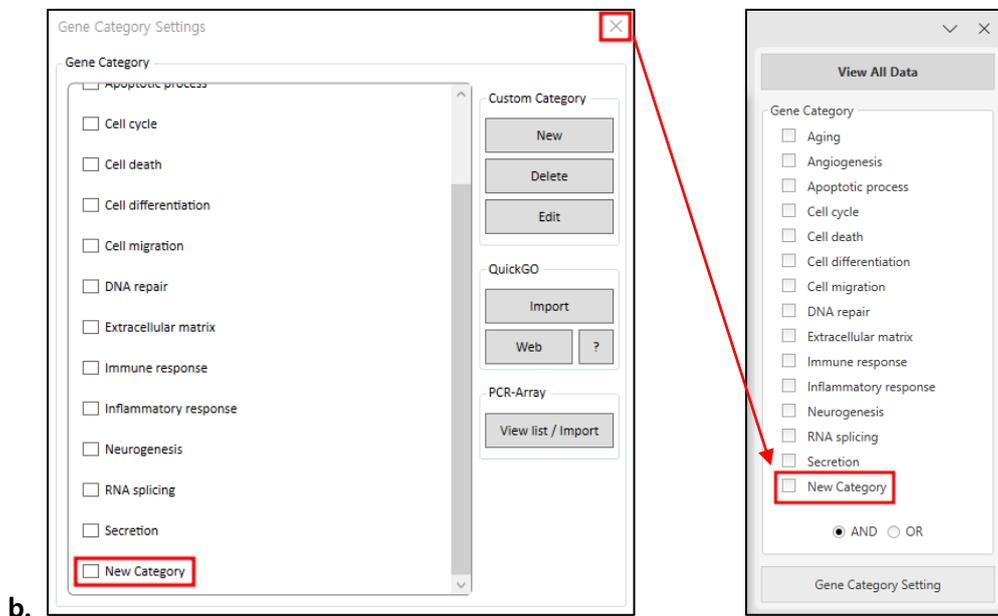
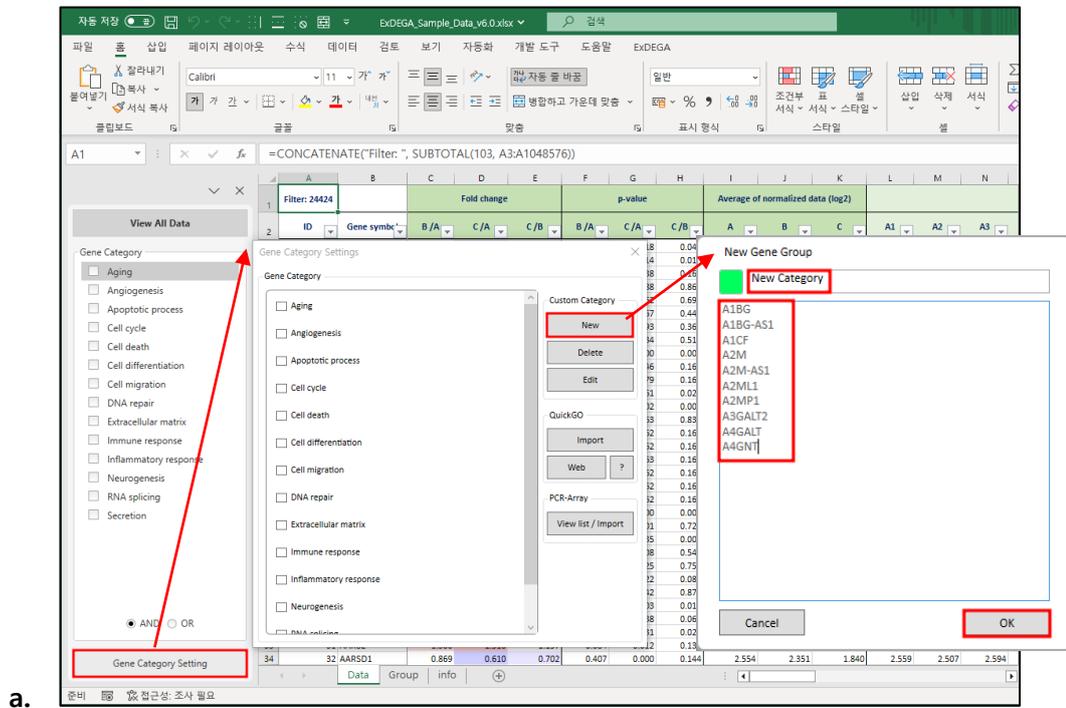


그림 1-5. Adding Genes to make a new gene category

생성된 Category 를 수정 혹은 삭제하고 싶다면 Category 선택 후 Delete 를 선택하면 삭제를, Edit 을 클릭하면 수정을 할 수 있다(그림 1-6).

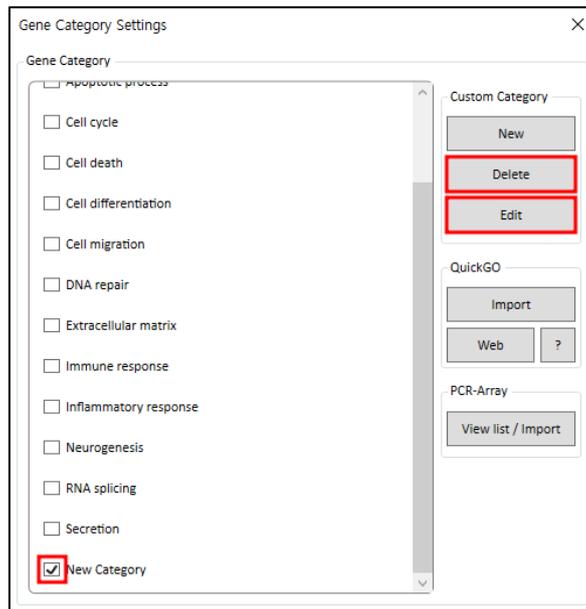


그림 1-6. Edit or Delete gene category

PCR-Array 항목의 View list / Import 를 이용하여 Pathway 별 Gene list 를 추가 할 수 있다. Gene Category Settings 버튼을 누른 후 View list / Import 버튼을 누른다. Sub Window 창에서 species 를 선택하고 Keyword 에 추가하고자 하는 Pathway 이름이나 유전자 이름을 검색하고 Check box 에 체크한 뒤 Import 버튼을 누르면 자동으로 추가된다(그림 1-7).

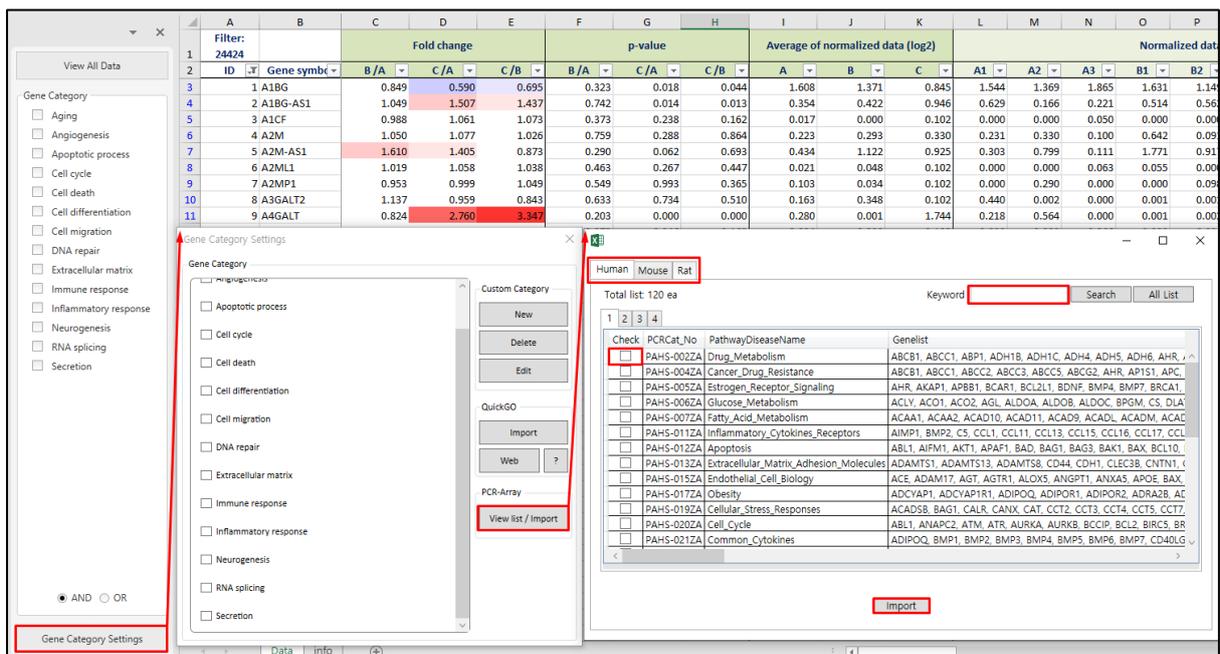


그림 1-7. PCR-Array Pathway settings

## 1-2. Significant Gene Selection 사용 방법

오른편의 DEG Analysis 부분에서 "Significant Gene Selection" 창은 전체 결과 중 대조군과 실험군을 비교한 결과에서 유의하게 발현 차이가 나는 유전자를 필터링 할 수 있도록 만들어 놓은 것이다. 예를 들어, B/A 비교조합을 선택하고 [Fold change : 2, Normalized Data (log2) : 4, p-value : 0.05] 를 선택하면, A 대비 B 에서 2 배 이상 발현이 증가 또는 감소하고, Normalized Data (log2)값이 4 이상이고, p-value 값이 0.05 이하인 유전자가 필터링 된다(그림 1-8). p-value 는 반복 실험한 데이터(N>=2)의 경우만 제공된다. 비교조합은 다중 선택할 수 있으며 "AND"나 "OR"를 기능을 이용하여 선택한 비교조합들에서 공통적인 DEGs (교집합) 또는 하나의 비교조합 이상 DEGs (합집합)을 선별할 수 있다.

유전자 선별 기준은 연구자의 데이터에 맞게 조정하여 사용할 수 있다.

ID	Gene symbol	Fold change			p-value			Average of normalized data (log2)			Normalized data (log2)						
		B/A	C/A	C/B	B/A	C/A	C/B	A	B	C	A1	A2	A3	B1	B2	B3	C1
111	ABHD3	2.437	2.076	0.853	0.002	0.000	0.176	2.778	4.063	3.834	2.936	2.645	2.739	4.281	3.902	3.979	3.832
370	ADGRE3	4.134	4.208	1.016	0.015	0.000	0.807	3.071	5.118	5.144	2.538	3.332	3.228	5.482	5.188	4.326	5.142
611	ALDH2	0.345	0.310	0.608	0.002	0.001	0.026	5.497	3.960	3.243	5.552	5.264	5.648	3.824	3.759	4.247	3.241

그림 1-8. Significant gene selection

Significant gene selection 에서 증가 또는 감소한 유전자를 각각 보고 싶다면 Up/Dn 의 selection box 에서 선택할 수 있다. Both 는 증가, 감소 유전자가 모두 필터링 되고 Up 은 증가한 유전자만 Dn 은 감소한 유전자만 따로 필터링 할 수 있다(그림 1-8).

ID	Gene symbol	Fold change			p-value			Average of normalized data (log2)			Normalized data (log2)						
		B/A	C/A	C/B	B/A	C/A	C/B	A	B	C	A1	A2	A3	B1	B2	B3	C1
105	ABHD14B	0.697	0.698	0.433	0.090	0.004	0.006	4.404	3.883	2.692	4.032	4.596	4.521	3.532	3.959	4.085	2.690
109	ABHD17A	0.779	0.358	0.464	0.157	0.000	0.026	6.466	6.109	4.946	6.490	6.513	6.400	5.709	5.963	6.498	4.984
113	ABHD3	2.437	2.076	0.853	0.002	0.000	0.176	2.778	4.063	3.834	2.936	2.645	2.739	4.281	3.902	3.979	3.832

그림 1-8. Significant genes (separately up and down)

Gene Category와 Significant gene selection은 연동 가능하다. 그림 1-9 에서처럼 Significant Gene Selection에서 선별조건을 설정하고 Gene Category의 GO를 선택하면 모두 적용된 유전자가 선별된다. 선별된 유전자는 예시 데이터에서 선택한 GO와 관련된 유전자들 중 선별조건이 적용된 유전자를 의미한다.

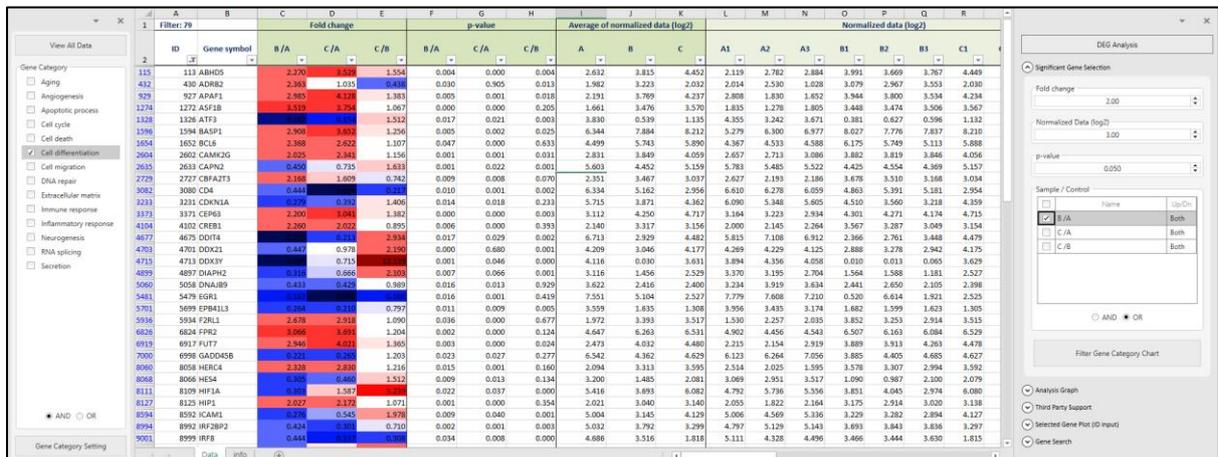


그림 1-9. Significant genes related to Cell differentiation

Gene Category Chart는 각 GO 관련 유전자 중 발현이 유의하게 차이 나는 유전자의 %와 수를 나타낸 그래프이다. 본 분석을 통해 어떤 GO의 유전자들이 상대적으로 많은 발현 변화가 있었는지를 확인할 수 있다. 전체 데이터 상태에서 Significant Gene Selection의 비교 그룹을 선택하고 "View Gene Category Chart"를 클릭하면 증가/감소한 유전자들 대상으로 GO Chart가 생성된다. 그래프의 각 영역을 클릭하면 해당 유전자들이 필터링 된다. 예를 들어 왼쪽의 Pie chart의 특정영역을 클릭하면 해당 GO의 증가/감소된 유전자가 함께 필터링 된다. 오른쪽의 bar chart에서 bar 상단의 숫자는 증가/감소에 해당하는 유전자 수이며, bar는 각각 Up significant(%)와 Down significant(%)를 의미한다. bar를 클릭했을 때 해당 유전자가 필터링 된다(그림 1-10).

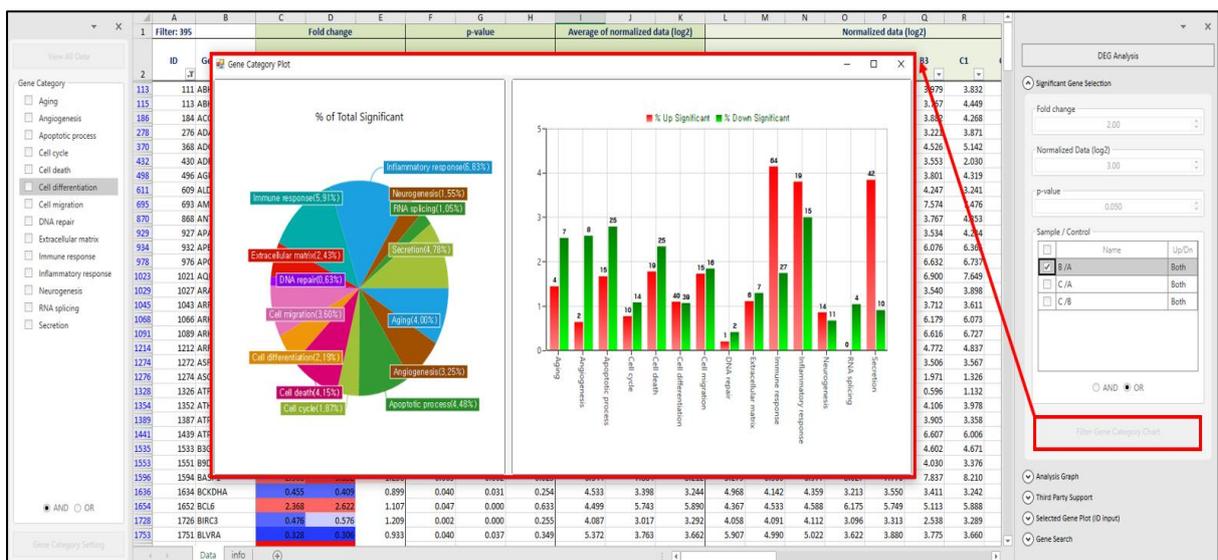


그림 1-10. View Gene Category Chart

significant chart 는 선택한 비교조합에 따른 유의한 유전자(최대 30 개)의 발현 값을 그룹 별로 확인할 수 있다. 유의한 유전자 필터링 기준은 선택한 비교조합의 p-value 순으로 표시된다. 차트의 x 축은 유전자 이름이고 y 축은 Normalized data (log2) 로 구성되어 있다. 또한 각 점(dot)은 샘플 하나를 의미하고 그룹에 따라 색이 구분된다. 마우스 커서를 각 점(dot)위에 올리면 해당되는 샘플 명도 알 수 있다(그림 1-11).

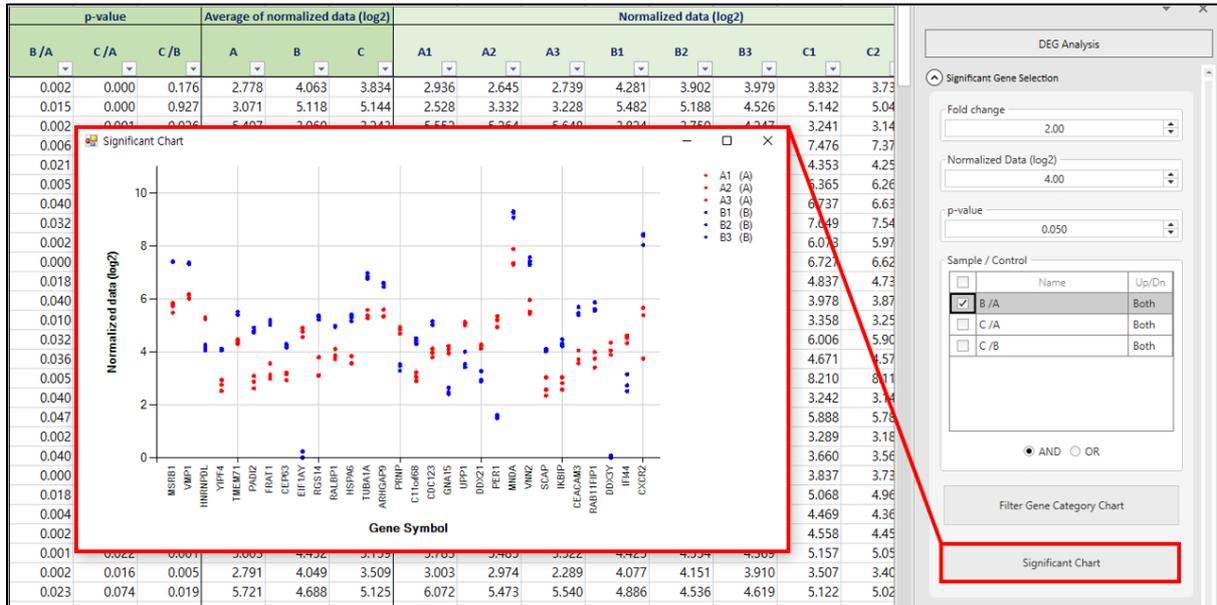


그림 1-11. Significant Chart

### 1-3. Analysis Graph 사용 방법

DEG Analysis 부분에서 "Analysis Graph" 창을 펼치면 그림 1-12 와 같이 Scatter Plot, Volcano Plot, Venn Diagram 을 쉽게 그릴 수 있다.

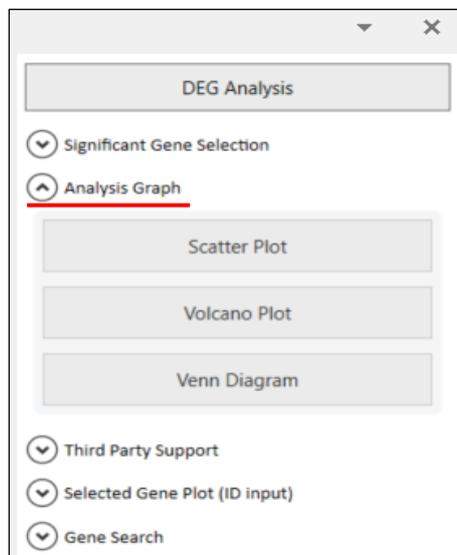


그림 1-12. Analysis Graph Tool

### 1-3-1. Scatter plot

Scatter Plot 은 대조군과 실험군의 발현양상을 확인할 수 있는 이미지이다. 오른쪽에 샘플 비교 그룹과 Fold threshold line (예시: 2fold)을 선택하고 "Graph View"를 클릭하면 왼쪽에 선택한 비교 그룹을 대상으로 Scatter Plot 이 자동 생성된다. 생성된 그래프는 마우스 스크롤로 크기 조절할 수 있다. x 축은 대조군의 normalized data (log2), y 축은 실험군의 normalized data (log2)이다. 초록색 사선 아래는 2fold 이상 감소한 유전자들, 빨간색 사선 위 2fold 이상 증가한 유전자들이다. 유의한 유전자를 식별하기 위한 색은 빨강, 파랑, 초록 중 사용자가 선택할 수 있다. Plot 에서 특정 spot 을 클릭하면 해당 유전자명이 표시되고 마우스 오른쪽을 클릭하여 지울 수도 있다. 표시된 유전자명은 마우스로 위치 조절이 가능하고 "Font Size"에서 표시된 유전자명의 글씨 크기를 조절할 수 있으며 "Hide"를 클릭하면 그래프의 grid line 을 숨길 수 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 1-13).

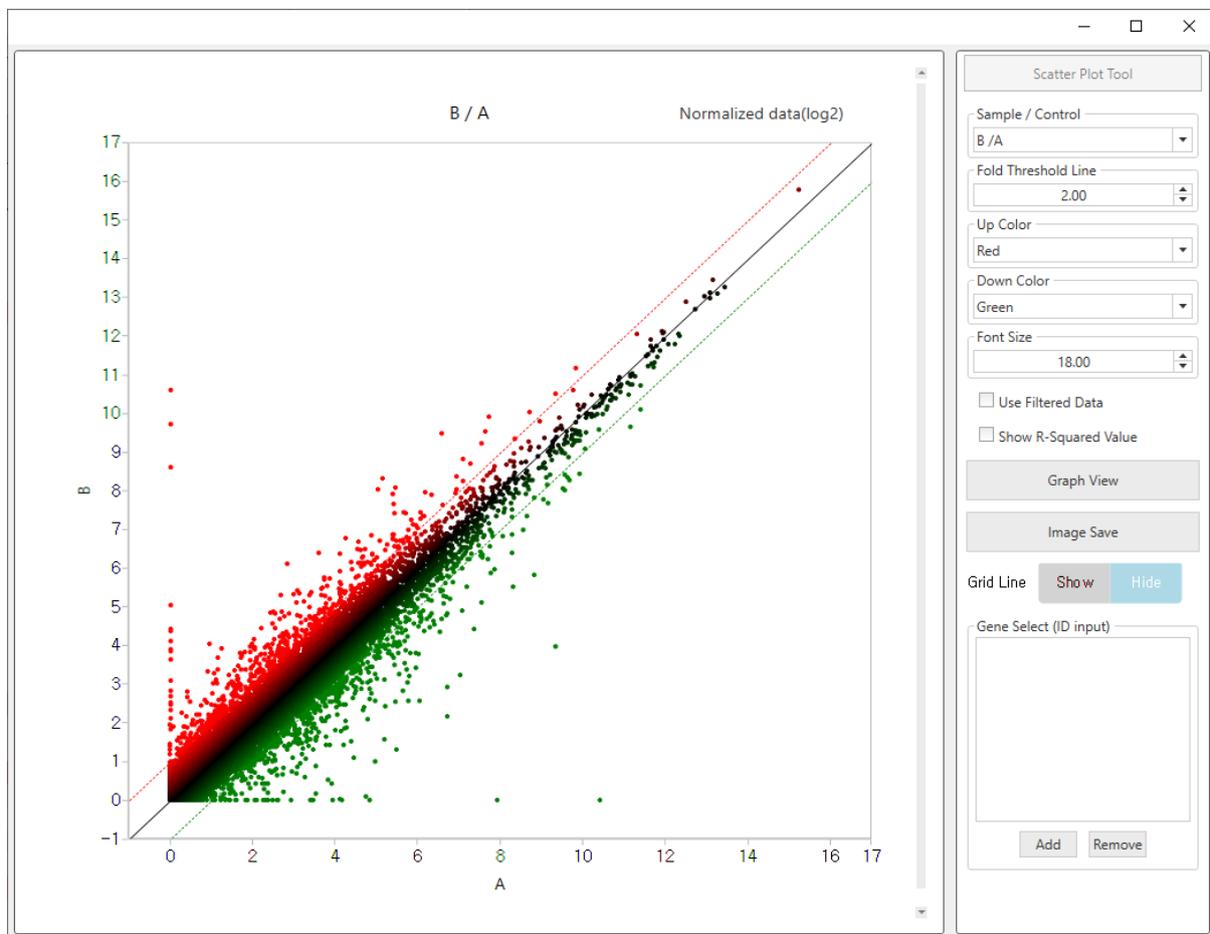


그림 1-13. Analysis Graph Tool – Scatter Plot

생성된 Scatter Plot은 전체 유전자를 기반으로 생성된 그래프이며 원하는 GO category에 대한 Scatter Plot을 그릴 수 있다. GO category 선택 후, "Use Filtered Data"를 체크하고 "Graph View"를 클릭하여 그래프를 그리면 선택한 GO Category를 기반으로 한 Scatter Plot을 생성할 수 있다(그림 1-14).

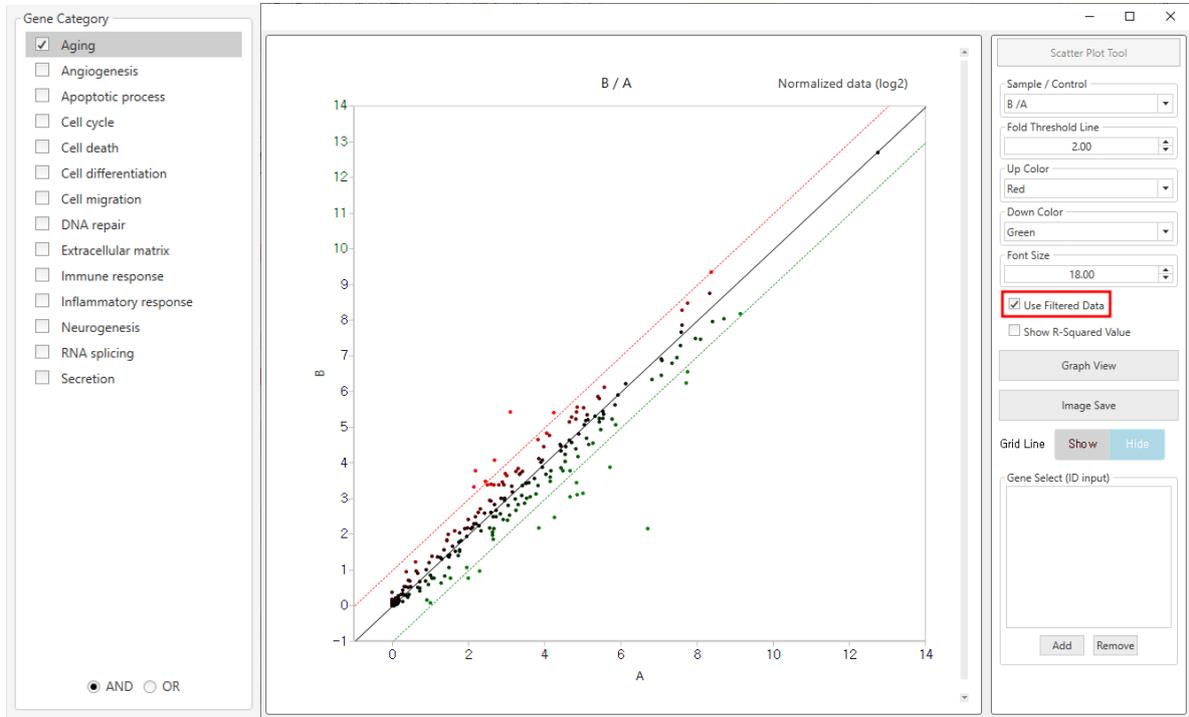
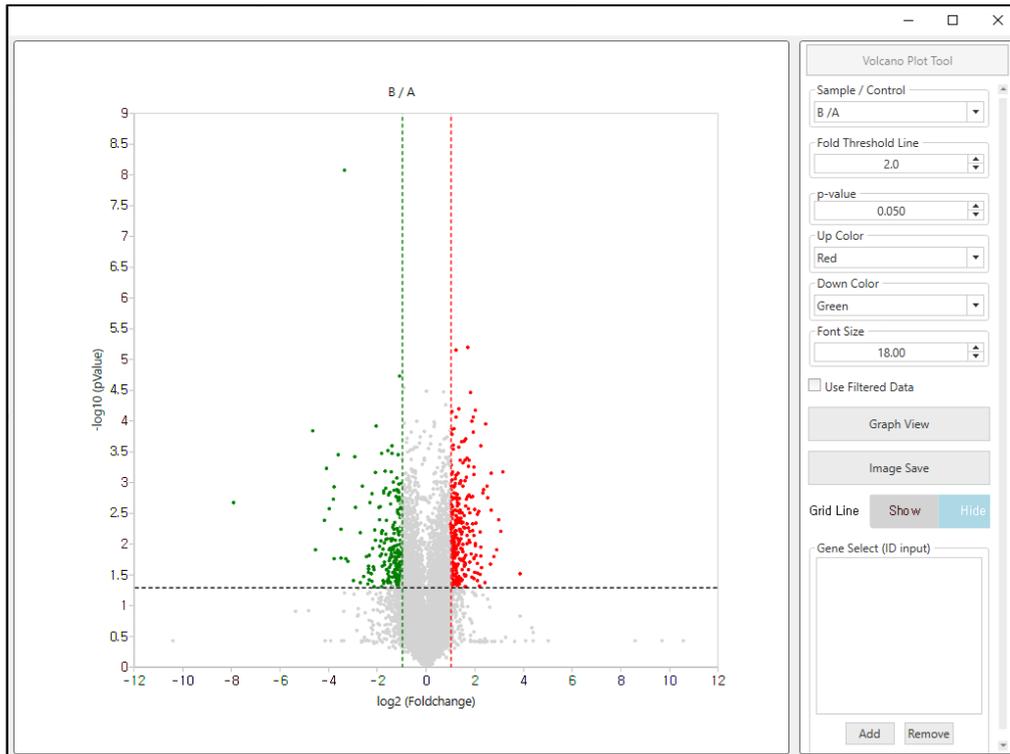


그림 1-14. Analysis Graph Tool – Use Filtered Data

### 1-3-2. Volcano plot

Volcano Plot은 반복 실험( $N \geq 2$ )이 된 경우에만 분석 가능하다. Volcano Plot은 Scatter Plot의 기능과 거의 동일하는데 오른쪽에 샘플 비교 그룹과 Fold threshold line (예시: 2fold), p-value (예시: p-value 0.05)를 선택하고 "Graph View"를 클릭하면 왼쪽에 선택한 비교 그룹을 대상으로 Plot이 자동 생성된다. 생성된 그래프는 마우스 스크롤로 크기 조절할 수 있다. 초록색 세로선 왼쪽은 2fold 이상 감소한 유전자들, 빨간색 세로선 오른쪽은 2fold 이상 증가한 유전자들, 검은색 가로선 위는 p-value 0.05 이하인 유전자들이다. 유의한 유전자를 식별하기 위한 색은 빨강, 파랑, 초록 중 사용자가 선택할 수 있다. Plot에서 특정 spot을 클릭하면 해당 유전자명이 표시되고 마우스 오른쪽을 클릭하여 표시를 지울 수도 있다. 표시된 유전자명은 마우스로 위치 조절이 가능하고 "Font Size"에서 표시된 유전자명의 글씨 크기를 조절할 수 있으며 "Hide"를 클릭하면 그래프의 grid line을 숨길 수 있다. 그리고 여러 개의 유전자를 동시에 표시하고 싶다면 "Gene Select(ID Input)" 창에 해당 유전자 ID를 복사하여 입력하고 "Add"를 클릭하면 Gene Symbol이 자동 생성된다(그림 1-15).



**그림 1-15. Analysis Graph Tool – Volcano Plot**

생성된 Volcano Plot은 전체 유전자를 기반으로 생성된 그래프이며 원하는 GO category에 대한 Volcano Plot을 그릴 수 있다. GO category 선택 후, "Use Filtered Data"를 체크하고 "Graph View"를 클릭하여 그래프를 그리면 선택한 GO Category를 기반으로 한 Volcano Plot을 생성할 수 있다 (그림 1-14).

### 1-3-3. Venn diagram

Venn Diagram 을 통해 4 개 이하의 비교조합을 대상으로 Venn Diagram 을 작성할 수 있다. Venn Diagram 을 그릴 Fold Change 와 Normalized data (log2), p-value 을 선택 후, Diagram View 를 클릭하면 결과를 확인할 수 있으며 비교조합은 최대 4 개까지 선택 가능하다. 그리고 "theme"을 통해서 Venn Diagram 의 색상 테마 선택이 가능하다. 아래의 그림은 B/A, C/A, C/B 결과 중, 2fold, Normalized data (log2) 4 이상, p-value 0.05 이하인 유전자 list 를 가지고 Venn Diagram 을 작성한 결과이다(그림 1-16).

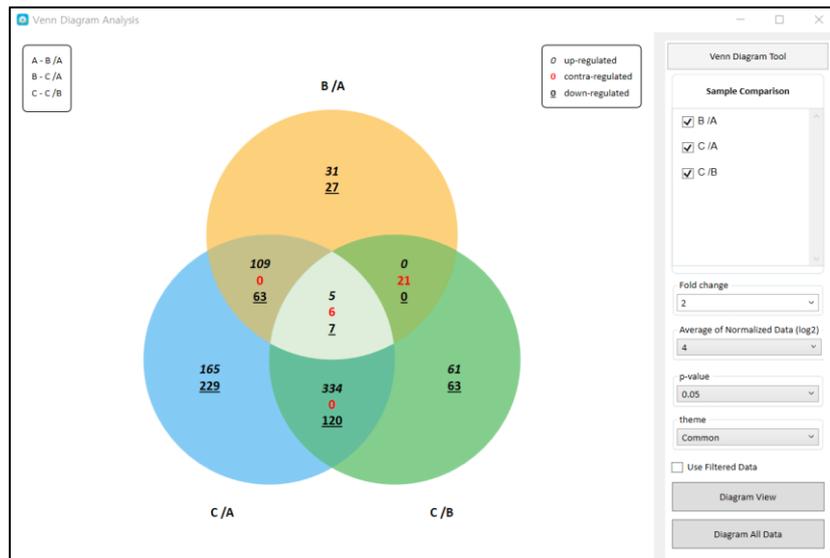


그림 1-16. Analysis Graph Tool – Venn Diagram

생성된 Venn Diagram 은 전체 유전자를 기반으로 생성된 그래프이며 원하는 GO Category 에 대한 Venn Diagram 을 그릴 수 있다. GO category 선택 후, "Use Filtered Data"를 체크하고 Diagram View 를 클릭하여 그래프를 그리면 선택한 GO Category 를 기반으로 한 Venn Diagram 을 생성할 수 있다(그림 1-17).

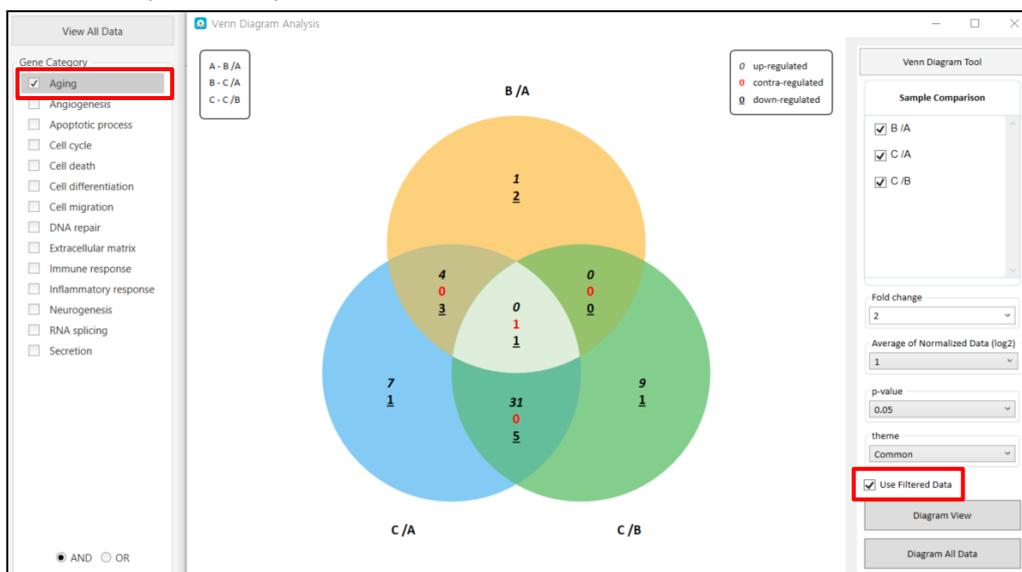


그림 1-17. Use Filtered Data

Venn Diagram 결과에서 표시되는 형식은 다음과 같다(그림 1-18).

1. **기울어진 숫자** : up-regulated 된 gene 수
2. **빨간색 숫자** : regulation 이 대조되는 gene 수  
(예: B/A 에서는 up 되고 C/A 에서는 down 되는 gene 수)
3. **밑줄 친 숫자** : down-regulated 된 gene 수

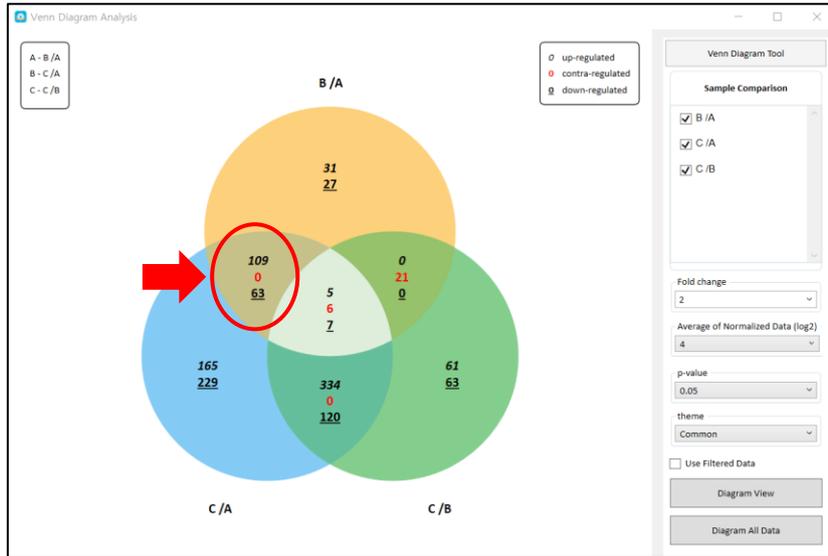


그림 1-18. For example of up, down, contra-regulated in Venn Diagram

Venn Diagram 각 영역에 어떤 유전자들이 있는지 확인할 수도 있다. 예를 들어, B/A 에서만 up 이 되는 유전자를 보고 싶으면, Venn Diagram 에서 B/A 에서만 해당되는 영역을 찾아 마우스 오른쪽 클릭하고 up-regulated 를 선택하면 증가된 유전자 list 가 엑셀 data sheet 에 filter 된다 (그림 1-19).

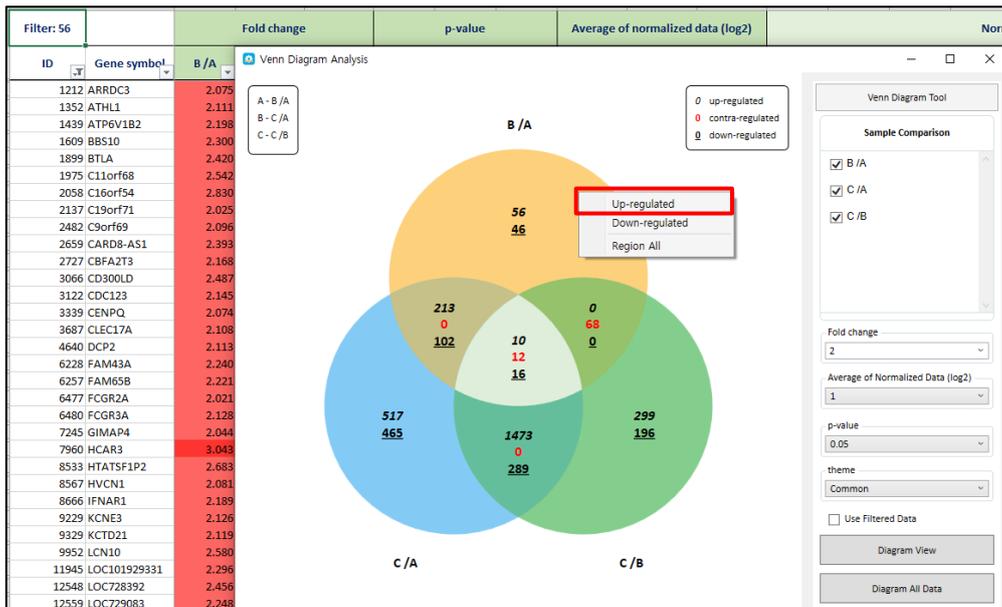


그림 1-19. Filtering 2fold up-regulated gene list in Venn Diagram

ExDEGA 에서 제공되는 모든 이미지는 오른쪽마우스를 눌러 'Save image' 버튼을 통해 저장 가능하다(그림 1-20).



그림 1-20. Save image

#### 1-4. Third Party Support 사용 방법

Third Party Support 는 연구자가 선별한 유전자를 기반으로 Clustering heatmap 과 DAVID 및 KEGG 분석을 수행하기 위한 입력 데이터를 제공한다. (단, GSEA의 경우 선별하지 않고 전체 유전자 리스트를 기반으로 진행하는 것을 권장한다.) 먼저, Input File 제작에 앞서 유전자를 선별하는 것이 필요하다. 유전자 선별은 유의성 있는 DEG 분석, Gene Ontology 분석 등으로 선별할 수 있다 (그림 1-21).

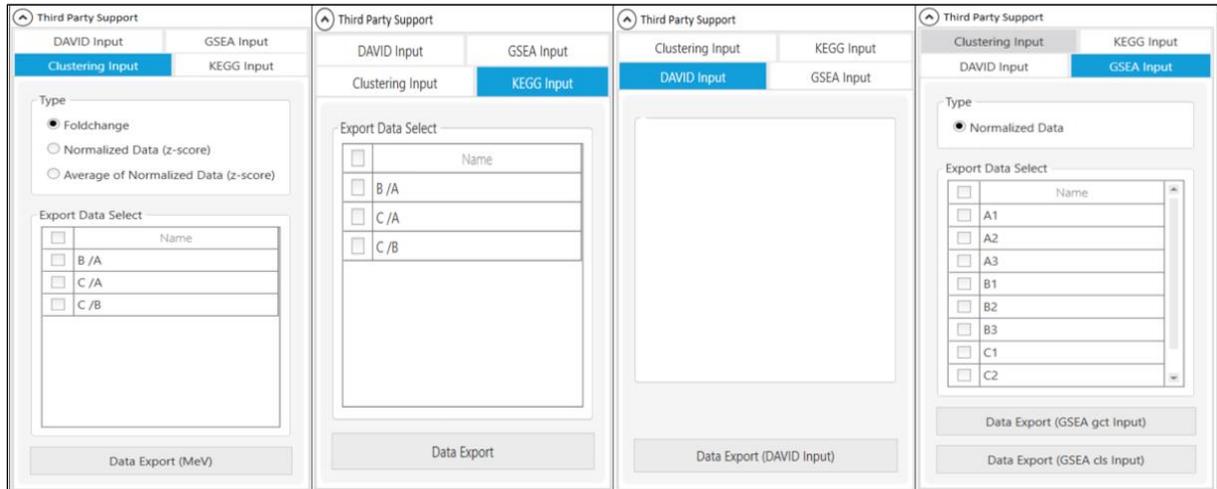


그림 1-21. Third Party Support

(왼편부터 Clustering Input, KEGG Input, DAVID Input, GSEA Input)

필터링 된 유전자 리스트를 대상으로 Clustering Heatmap Input 은 크게 세 종류의 데이터를 이용할 수 있다. Input 파일 저장 시, 파일명에는 띄어쓰기가 들어가지 않도록 주의해야 하며 저장된 Clustering Input 파일은 GraphicPlus 에서 heatmap 을 작성하거나 MeV 프로그램을 이용하여 heatmap 을 작성할 수 있다.

첫 번째, Fold change 값을 이용할 시 Type 부분에 Fold change 를 체크하고 Export Data Select 에서 Heatmap 에 표현할 비교조합을 체크한다. "Data Export"를 클릭한 후 "(input 명).txt"로 저장한다. 저장된 파일은 Fold change 가 log2 변환된 값으로 추출된다.

두 번째, 개별 샘플의 발현 값인 Normalized Data 값을 이용할 시 Type 부분에 Normalized data (z-score) 를 체크하고 이용하고자 하는 샘플을 체크한다. "Data Export"를 클릭한 후 "(input 명).txt"로 저장한다. 단, Z-score 로 그릴 때는 표준편차가 고려되기 때문에 샘플 3 개 이상에서만 가능하다

세 번째, 그룹의 발현 값인 Average of Normalized Data 값을 이용할 시 Type 부분에 Average of Normalized Data (z-score)를 체크하고 이용하고자 하는 그룹명을 체크한다. "Data Export"를 클릭한 후 "(input 명).txt" 로 저장한다. 단, Z-score 로 그릴 때는 표준편차가 고려되기 때문에 그룹 3 개 이상에서만 가능하다.

만약 이용하고자 하는 샘플 혹은 그룹의 수가 2 개 이하인 경우, z-score 변환이 불가능하므로 아래의 방법으로 input 파일을 직접 작성한다. 우선 유의성 있는 유전자를 선별한 후, 새로운 excel 에 선별된 Gene list 와 이용하고자 하는 Average of normalized data (log2) 혹은 Normalized data (log2)를 복사하여 input 파일을 만든다 (그림 1-22). input 파일 저장 시, 파일 형식은 “텍스트(탭으로 분리)(\* .txt)”을 선택한다. 저장된 input 파일은 MeV 프로그램을 이용하여 Heatmap 을 작성할 수 있다.

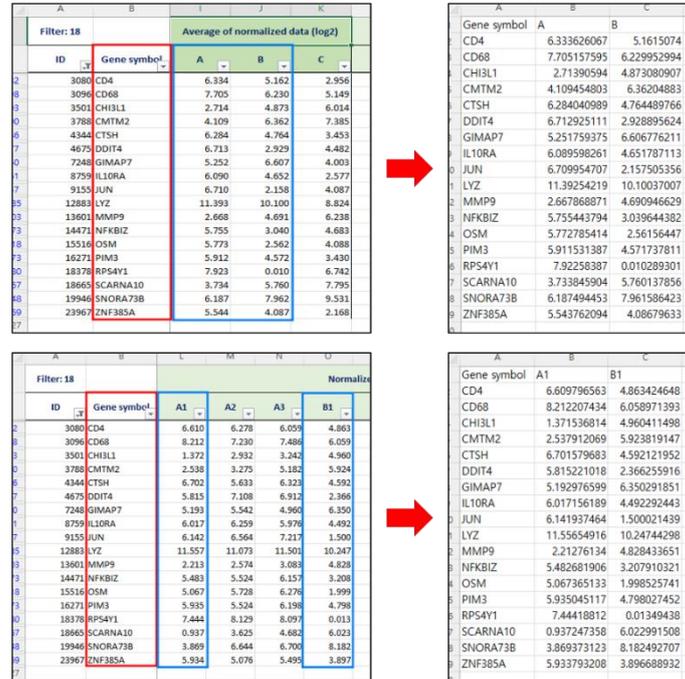


그림 1-22. Heatmap input of no more than 2 samples or groups

\* Z-score 는 일반적으로 평균으로부터 얼마만큼 떨어져 있는지를 판단하는 지표이다. 계산방식은 Normalized data 를 log10 으로 변환 후 평균값을 뺀 후 표준편차로 나누어 계산한다.  

$$Z\text{-score} = \frac{\{\text{Normalized data (log10)} - \text{average of Normalized data (log10)}\}}{\text{standard deviation of Normalized data(log10)}}$$

Clustering heatmap input file 은 Gene Symbol 과 fold change (log2) 또는 z-score 로 구성된다(그림 1-23). MeV 프로그램을 사용하여 Clustering heatmap 을 작성하는 방법은 MeV manual ([Download Link](#))에서 확인할 수 있다.

Gene Symbol	B /A	C /A	C /B
ABHD3	1.28484278864215	1.05564963600944	-0.229193276278856
ADGRE3	2.04741030637742	2.07337604787126	0.0259657102702727
ALDH2	-1.53704893440065		-2.25385030388101
AMICA1	1.38502619744398	1.41309305964294	0.0280669503482263
ANTXR2	1.19047263791619	1.27486038592147	0.0843881542963257
APBB1IP	1.43866864864493	1.45671679922983	0.0180472923397484
APOBR	1.2529266643143	1.79222275950347	0.539295814093037
AQP9	1.18587922828891	1.49957447547321	0.313695136333861

그림 1-23. Third Party Support - Clustering input

KEGG input 은 분석 결과에서 Up-/Down-regulated 된 유전자들이 어떤 Pathway 에 속하는지 확인하고자 할 때 **KEGG Mapper** 를 이용하기 위한 입력 데이터를 제공한다. KEGG Input 에서는 하나의 비교 조합만 선택 가능하다. 비교 조합 선택 후 Data Export 를 선택하면 그림 1-24 과 같이 유전자의 Entrez ID 와 비교 조합, 발현 수준에 따른 색 코드로 구성된다. KEGG Input 파일은 **KEGG Mapper** 의 입력 데이터로 사용하여 Pathway 상에 속하는 유전자와 이들의 발현 수준을 확인할 수 있다.

```

파일(F) 편집(E) 서식(O) 보기(V) 도움말(H)
Entrez ID B /A
18      #FFE4c3,black
100     #E0FFFF,black
133     .
148     .
177     #FFA07A,black
183     .
207     .
218     .
239     .
267     #B0E0E6,black
268     .
317     #FF6347,black
    
```

그림 1-24. Third Party Support - KEGG input

DAVID 는 다양한 데이터베이스를 기반으로 유전자의 상관관계를 통계적으로 분석하여 유전자의 주요 기능을 예측하는 Analysis tool 이다. DAVID 는 3 천 개 이상의 유전자는 분석할 수 없으므로 3 천 개 이하로 유전자를 선별해야 한다. Data Export (DAVID Input)는 DAVID 에서 유전자 이름을 입력하는 부분에 사용되는 파일이다. 이 파일은 유전자 이름으로 구성된다(그림 1-25). 유전자 이름을 DAVID 의 입력 데이터로 사용하여 분석하려는 GO 분석 결과 데이터를 다운로드 받는다.

```

DAVID_input.txt - Windows 메모장
파일(F) 편집(E) 서식(O) 보기(V) 도움말
ABHD3
ADGRE3
ALDH2
AMICA1
ANTXR2
APBB1IP
APOBR
AQP9
ARHGAP25
ARHGAP9
ARRDC3
ATHL1
ATP1B3
ATP6V1B2
B3GNT8
    
```

그림 1-25. Third Party Support - DAVID input

### 1-5. Selected Gene Plot 사용 방법

ExDEGA 의 기능 중에 선별한 유전자 또는 연구자가 관심있는 유전자들을 대상으로 발현 패턴을 그래프로 표현하고자 할 때는 "Selected Gene Plot" 기능을 사용할 수 있다. 선별한 유전자의 ID 을 복사하여 Selected Gene Plot 창에 붙여 넣고 "Expression Plot View"를 누르면 Normalized data (log2) 값, Fold change (log2) 값으로 line graph 가 그려진다(그림 1-26).

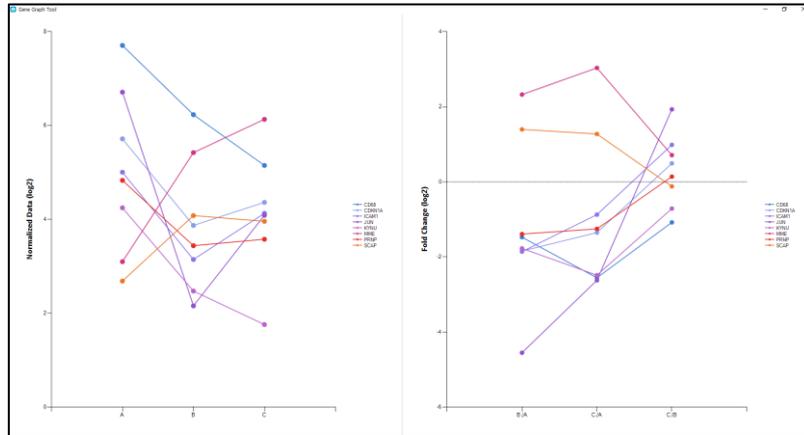


그림 1-26. Gene graph

### 1-6. Rader Chart 사용 방법

ExDEGA 의 기능 중에 선별한 유전자나 연구자가 관심있는 유전자들을 대상으로 발현 패턴을 그래프로 표현하고자 할 때는 "Radar chart" 형식의 이미지로도 표현할 수 있다. 사용방법은 선별한 유전자의 ID 을 복사하여 Radar chart 창에 붙여 넣고 옵션을 선택한다. 옵션은 Radar chart 의 꼭지점으로 활용할 항목 (유전자 or 샘플명)을 선택하고, 표현하고자 하는 발현값 (평균 발현값 or 개별 샘플 발현값)을 선택한 뒤 "Radar chart View"를 누르면 Normalized data (log2) 값, Fold change (log2) 값으로 이미지를 그릴 수 있다(그림 1-27). 단, 꼭지점에 있는 항목이 3 가지 이상일 경우 제작 가능하다.

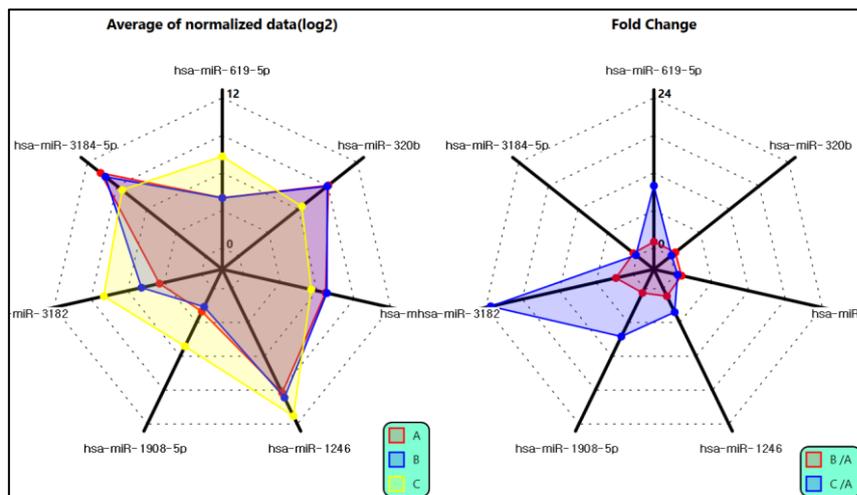


그림 1-27. Radar chart

### 1-7. Gene Search 사용 방법

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 행만 필터링 하여 확인할 수 있다(그림 1-28).

The screenshot shows a software interface with a data table and a search panel. The table has the following columns: ID, Gene symbol, Fold change (B/A, C/A, C/B), p-value (B/A, C/A, C/B), Average of normalized data (log2) (A, B, C), and Normalized data (log) (A1-A3, B1-B2). The search panel on the right has a text input field containing 'insulin' and a 'Search' button. The table shows a list of genes with their corresponding values, with some rows highlighted in red and blue.

ID	Gene symbol	Fold change			p-value			Average of normalized data (log2)			Normalized data (log)					
		B/A	C/A	C/B	B/A	C/A	C/B	A	B	C	A1	A2	A3	B1	B2	
8611	8609 IDE	0.790	1.015	1.284	0.143	0.887	0.054	1.742	1.403	1.764	1.940	1.794	1.487	1.427	1.604	
8704	8703 IGF1	1.266	1.014	0.779	0.290	0.705	0.414	0.086	0.465	0.102	0.002	0.138	0.096	0.085	1.015	
8705	8703 IGF1R	1.781	2.532	1.140	0.085	0.008	0.444	2.049	2.882	3.072	1.467	2.049	2.463	3.200	2.956	
8706	8704 IGF2	0.851	1.257	1.418	0.498	0.001	0.000	1.416	1.187	1.120	1.231	1.856	1.028	1.366	0.909	
8708	8706 IGF2BP1	0.992	1.036	1.040	0.910	0.667	0.336	0.129	0.117	0.180	0.001	0.002	0.306	0.342	0.094	
8709	8707 IGF2BP2	0.929	0.917	0.966	0.829	0.660	0.961	4.913	4.807	4.708	5.271	4.269	5.023	4.943	3.036	
8711	8709 IGF2BP3	0.759	2.549	3.360	0.507	0.009	0.001	1.219	0.822	2.547	1.884	0.447	0.949	1.180	0.902	
8712	8710 IGF2R	2.630	1.141	1.283	0.006	0.000	0.105	4.697	6.097	6.434	4.703	4.833	4.538	6.396	5.946	
8713	8711 IGF1ALS	0.914	0.957	1.047	0.130	0.635	0.383	0.166	0.036	0.102	0.115	0.359	0.000	0.104	0.001	
8714	8712 IGFBP1	1.000	1.185	1.185	0.119	0.017	0.017	0.000	0.000	0.245	0.000	0.000	0.000	0.000	0.000	
8715	8713 IGFBP2	0.639	1.157	1.749	0.512	0.000	0.000	1.284	0.638	4.322	0.186	0.374	2.285	0.741	0.711	
8716	8714 IGFBP3	0.685	0.801	1.140	0.140	0.000	0.000	1.110	0.564	5.359	0.005	1.281	1.335	0.165	0.648	
8717	8715 IGFBP4	0.675	1.147	1.698	0.298	0.000	0.000	1.709	1.343	4.146	0.034	2.184	1.811	0.750	1.023	
8718	8716 IGFBP5	0.971	4.157	4.451	0.281	0.000	0.000	0.959	0.018	2.182	0.004	0.000	0.111	0.033	0.619	
8719	8717 IGFBP6	1.059	4.809	4.541	0.755	0.000	0.000	0.588	0.671	2.854	0.475	0.378	0.885	0.236	0.788	
8720	8718 IGFBP7	0.516	1.147	1.698	0.182	0.000	0.000	4.131	3.176	7.469	4.634	3.119	4.249	2.556	2.944	
8722	8720 IGFBP11	1.000	0.565	0.565	0.997	0.857	0.794	0.187	0.166	0.165	0.126	0.051	0.363	0.169	0.118	
8911	8909 INS	1.000	1.073	1.073	1.000	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	
8914	8912 INSIG1	0.692	1.450	2.095	0.059	0.199	0.005	4.111	3.045	3.581	3.335	4.551	4.073	3.183	2.926	
8915	8913 INSIG2	1.379	1.268	0.914	0.138	0.207	0.419	1.817	2.281	2.134	1.546	2.234	1.508	2.028	2.323	
8916	8914 INS1L	0.949	0.809	0.852	0.860	0.392	0.494	0.720	0.845	0.411	0.006	0.928	1.024	0.728	1.025	
8917	8915 INS1L4	1.000	1.073	1.073	0.319	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	
8918	8916 INS1L5	1.000	1.073	1.073	1.000	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	
8919	8917 INS1L6	1.000	1.073	1.073	0.319	0.163	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	
8920	8918 INS1M1	1.000	1.073	1.073	1.000	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	
8921	8919 INS1M2	1.000	1.073	1.073	0.319	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	
8922	8920 INS1R	0.712	0.642	0.902	0.013	0.001	0.329	1.167	0.677	0.528	1.224	1.062	1.209	0.578	0.533	
8923	8921 INS1R1	1.000	1.073	1.073	0.319	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	
9007	9005 IRL1	0.920	2.508	2.718	0.610	0.001	0.000	0.505	0.396	1.838	0.636	0.342	0.390	0.460	0.481	
9008	9006 IRL2	0.451	1.367	2.918	0.023	0.067	0.000	4.513	3.410	4.984	4.086	4.648	4.776	3.395	3.029	
9009	9007 IRL4	1.000	1.073	1.073	1.000	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	
18516	18518 RXFP1	1.000	1.073	1.073	1.000	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	
18517	18515 RXFP2	0.989	1.052	1.063	0.653	0.130	0.223	0.029	0.014	0.102	0.005	0.000	0.000	0.000	0.040	
18518	18516 RXFP3	1.000	1.073	1.073	1.000	0.162	0.162	0.000	0.000	0.102	0.000	0.000	0.000	0.000	0.000	
18519	18517 RXFP4	1.199	2.078	1.733	0.731	0.025	0.126	0.515	0.777	1.570	0.006	1.191	0.000	1.380	0.202	

그림 1-28. Genes related to insulin

## 2. Functional Annotation Analysis (DAVID, ExDEGA GraphicPlus)

### 2-1. DAVID 분석 틀을 이용한 Functional Annotation 분석

DAVID는 다양한 데이터 베이스를 기반으로 유전자의 상관관계를 통계적으로 분석하여 유전자의 주요 기능을 예측하는 analysis tool이다. 분석 과정은 그림 2-1과 같다.

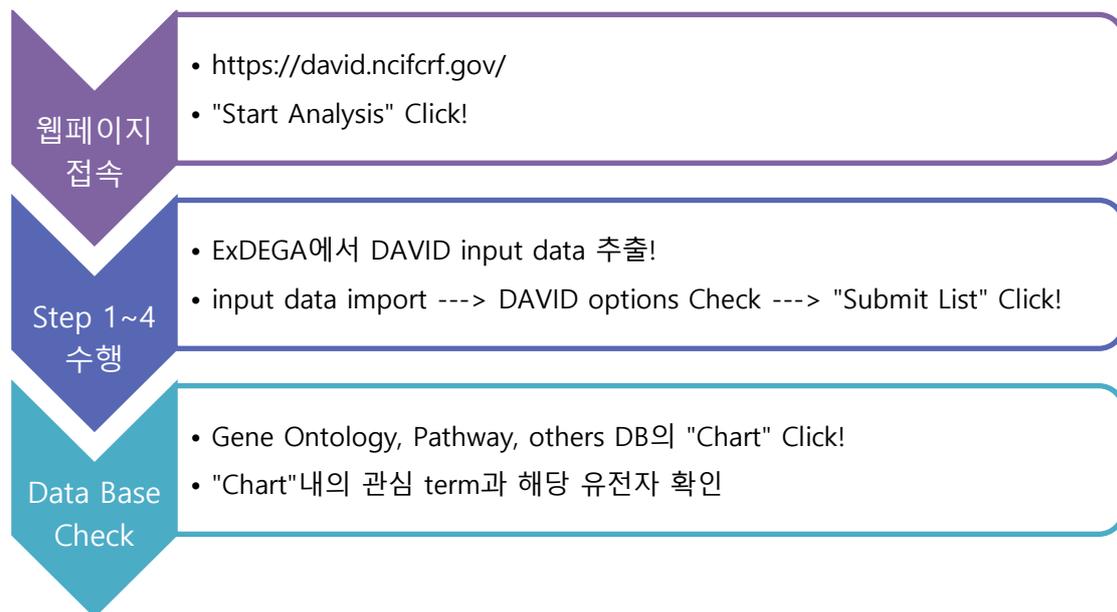


그림 2-1. DAVID tool analysis process

DAVID에서는 3천 개 이상의 유전자는 분석할 수 없으므로 3천 개 이하로 유전자를 선별해야 한다. RNA-Seq 결과에서 significant gene을 선별하여 DAVID 분석을 한다. Significant gene selection에서 Fold change, Normalized Data(log2), p-value (반복실험의 경우) 값을 지정하고, 확인하고자 하는 Fold change 조합을 선택하여 필터를 적용한다. 필터를 적용하여 선별된 유전자를 대상으로 Third Party Support를 통해 DAVID input을 추출하여, DAVID 분석에 사용한다. 반드시 하나의 비교 조합만 진행해야 한다(그림 2-2).

ID	Gene symbol	Fold change			p-value			Average of F
		B/A	C/A	C/B	B/A	C/A	C/B	
111	ABHD3	2.437	2.079	0.853	0.002	0.000	0.176	2.778
370	368 ADGRE3	4.134	4.209	1.018	0.015	0.000	0.927	3.071
611	609 ALDH2	0.345	0.210	0.608	0.002	0.001	0.026	5.497
695	693 AMICA1	2.612	2.663	1.020	0.006	0.005	0.764	6.065
870	888 ANXR2	2.282	2.420	1.060	0.021	0.000	0.718	3.081
924	932 APBB1P	2.711	2.745	1.013	0.005	0.001	0.904	4.910
978	976 APGBR	2.383	3.463	1.453	0.040	0.000	0.075	4.947
1023	1021 AQP9	2.275	2.828	1.243	0.032	0.003	0.163	6.151
1068	1066 ARHGAP25	3.604	2.757	0.765	0.002	0.001	0.068	4.612
1091	1089 ARHGAP29	2.176	2.444	1.123	0.000	0.000	0.094	5.440
1214	1212 ARHDC3	2.076	1.949	0.939	0.018	0.020	0.460	3.877
1554	1352 ATHL1	2.113	1.310	0.621	0.040	0.116	0.082	3.591
1389	1387 ATP1B3	0.454	0.412	0.834	0.010	0.002	0.323	4.640
1441	1439 ATP1B2	2.198	1.359	0.618	0.032	0.384	0.001	5.566
1535	1533 B3GNT8	2.471	3.434	1.390	0.036	0.001	0.085	2.894
1596	1594 BASP1	2.908	3.852	1.256	0.005	0.002	0.025	6.344
1636	1634 BCKDHA	0.455	0.468	0.899	0.040	0.031	0.254	4.533
1654	1652 BCL6	2.368	2.622	1.107	0.047	0.000	0.631	4.499
1728	1726 BIRC3	0.476	0.578	1.209	0.002	0.000	0.255	4.087
1753	1751 BLVRB	0.323	0.308	0.933	0.040	0.037	0.349	5.372
1977	1975 C11orf68	2.542	1.699	0.669	0.000	0.002	0.003	3.074
2060	2058 C16orf54	2.839	1.923	0.679	0.018	0.021	0.097	4.127
2350	2348 CSAR2	2.411	2.669	1.107	0.004	0.000	0.318	3.056
2622	2620 CANT1	2.312	2.566	1.110	0.002	0.000	0.278	3.201
2635	2633 CAPN2	0.450	0.735	1.633	0.001	0.022	0.001	5.603
2661	2659 CARD8-AS1	2.393	1.645	0.688	0.002	0.016	0.005	2.791
2711	2709 CAST	0.489	0.662	1.353	0.023	0.074	0.019	5.721
2983	2981 CCNH	0.404	0.677	1.677	0.019	0.112	0.001	4.239
2989	2987 CENL1	0.493	0.620	1.261	0.006	0.006	0.139	4.814
2999	2997 CCPG1	2.050	2.792	1.362	0.008	0.000	0.032	3.587
3033	3031 CD163	0.418	0.157	0.377	0.040	0.012	0.000	4.033
3082	3080 CD4	0.444	0.500	0.211	0.010	0.001	0.002	6.334
3086	3084 CD46	2.013	2.357	1.171	0.008	0.000	0.176	5.042

그림 2-2. DAVID input file generation process

DAVID 홈페이지 (<https://david.ncifcrf.gov/home.jsp>)에 접속하여 “Start Analysis”을 클릭한다(그림 2-3).

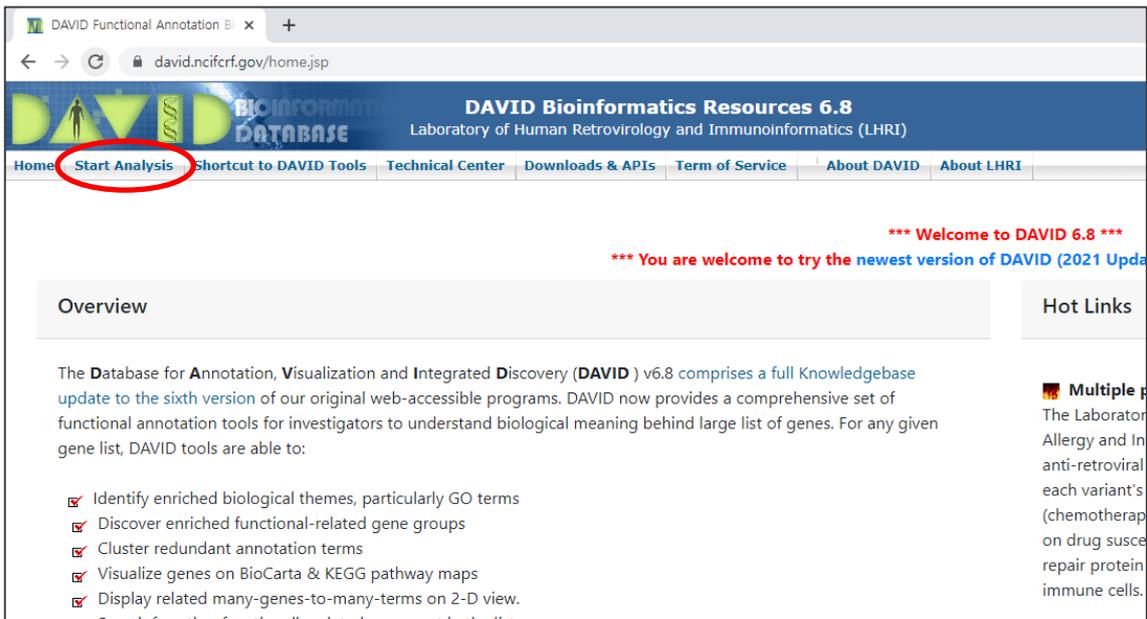


그림 2-3. DAVID tool webpage

“Upload” 탭에서 Step 1 에서 Step 4 까지 수행한다(그림 2-4). Step 1 에서 ExDEGA 에서 제작한 DAVID input 파일을 선택한다. Step 2 에서 “OFFICIAL\_GENE\_SYMBOL”를 선택한다. 만약 step 1 에서 Gene Bank No.를 넣었다면 “GENEBANK\_ACCESSION”을 선택한다. Step 2a 에서 분석하는 종의 학명을 입력한다. Step 3 에서 “Gene List”를 체크하고 Step 4 에서 “Submit List”를 누른다.

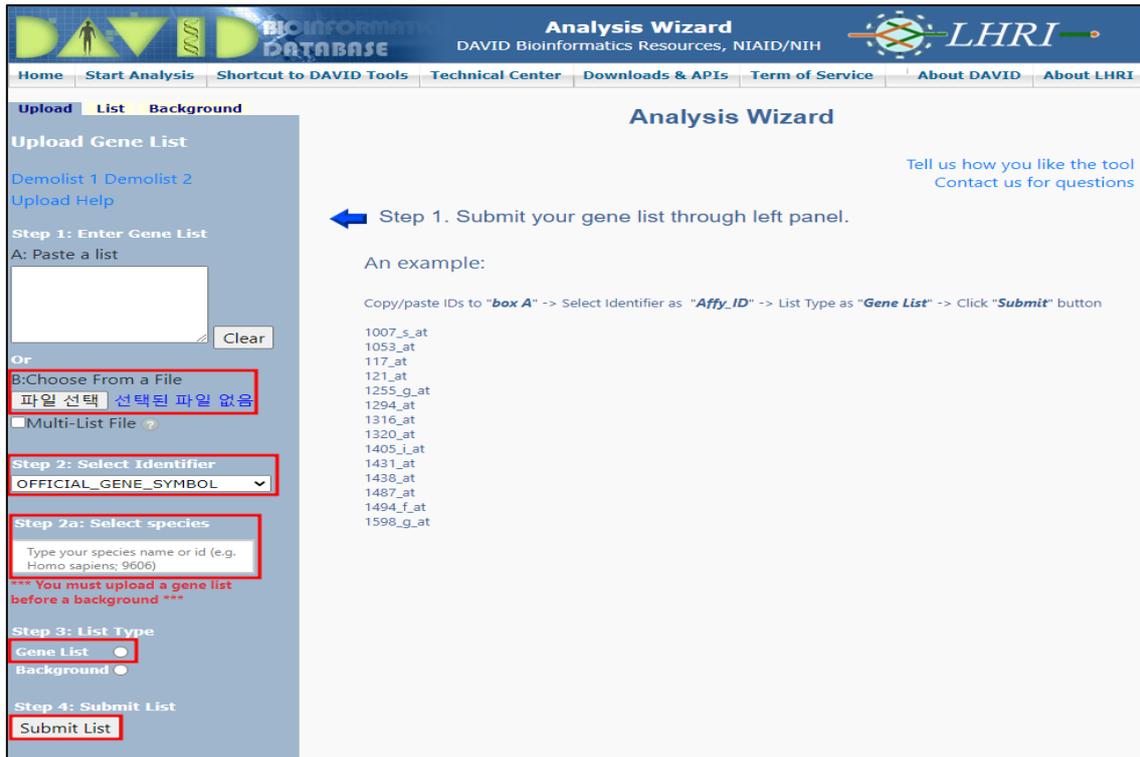


그림 2-4. DAVID tool : Step 1 ~ Step 4

List sheet 에서 분석하고자 하는 종을 선택한다(그림 2-5. a). "List" Sheet 에서 해당 종(숫자)로 표기되어 있고 가로 안의 숫자가 분석에 적용된 유전자의 개수이다. 예시에서는 269 개의 유전자 리스트를 넣었고 데이터베이스에서 기능이 밝혀진 267 개만이 Functional Annotation 분석에 이용되었다는 의미이다.

만약 Current Background 에 분석하고자 하는 종이 아닌 다른 종이 나왔다면 좌측 "Background" Sheet 에서 알맞은 종을 선택하여 "Use"를 클릭한다(그림 2-5. b).



a.



그림 2-5. DAVID tool : Select Species

확인 후, 화면에서 Functional Annotation Tool 을 클릭하면 결과가 나온다 (그림 2-6).

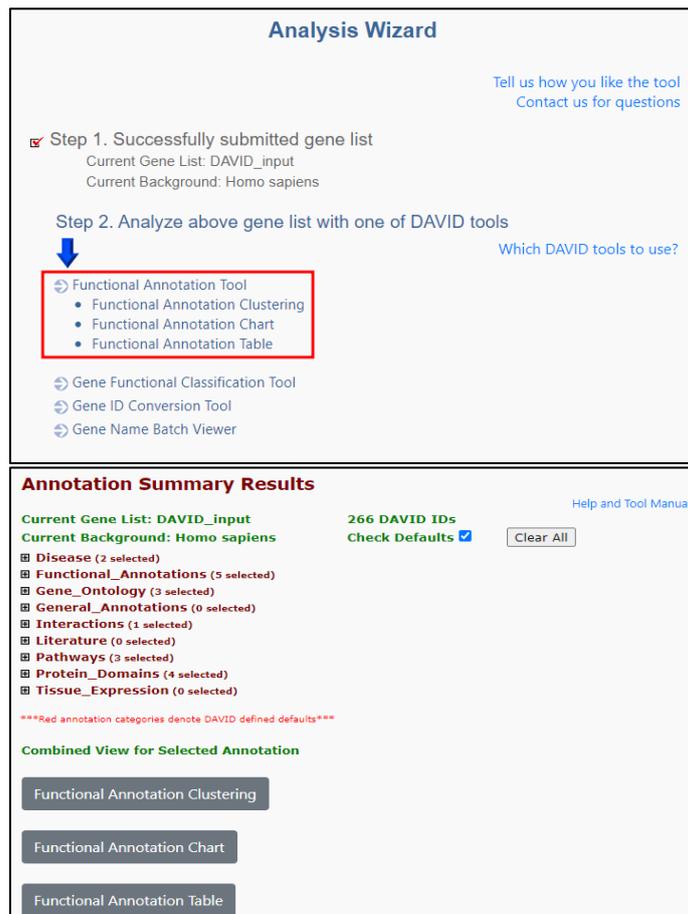


그림 2-6. DAVID results

DAVID 분석 결과 중 Gene Ontology Biological Process 결과를 확인하려면 "Gene\_Ontology"의 "+" 표시를 클릭하여 결과 창을 열고 "GOTERM\_BP\_Direct"의 "Chart"를 누른다(그림 2-7). Input 한 유전자들이 유의하게 관여하는 GO list 가 나온다. 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.

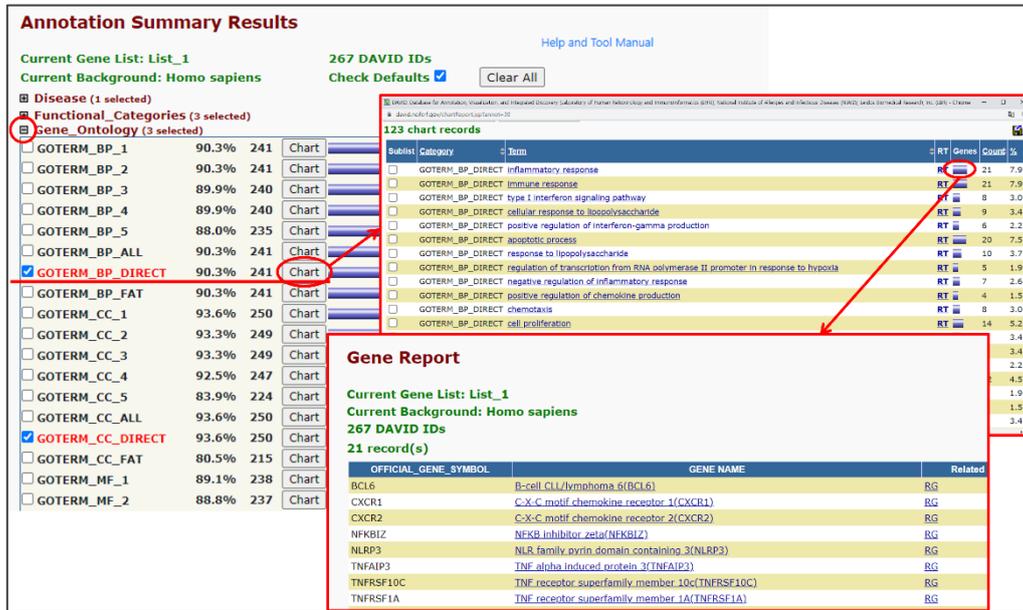


그림 2-7. DAVID tool : exploring Gene Ontology analysis result

이와 같은 방법으로 Pathway 결과를 확인해 보면 KEGG\_PATHWAY database 에서 주요 Pathway 가 나온다(그림 2-8). 각 pathway 를 누르면 pathway 그림을 확인할 수 있다. pathway 그림에서 별 표시가 되어 있는 유전자가 input 유전자 중 해당 pathway 에 관여하는 유전자이다. 유전자를 클릭하면 유전자 정보도 자세히 알 수 있다.

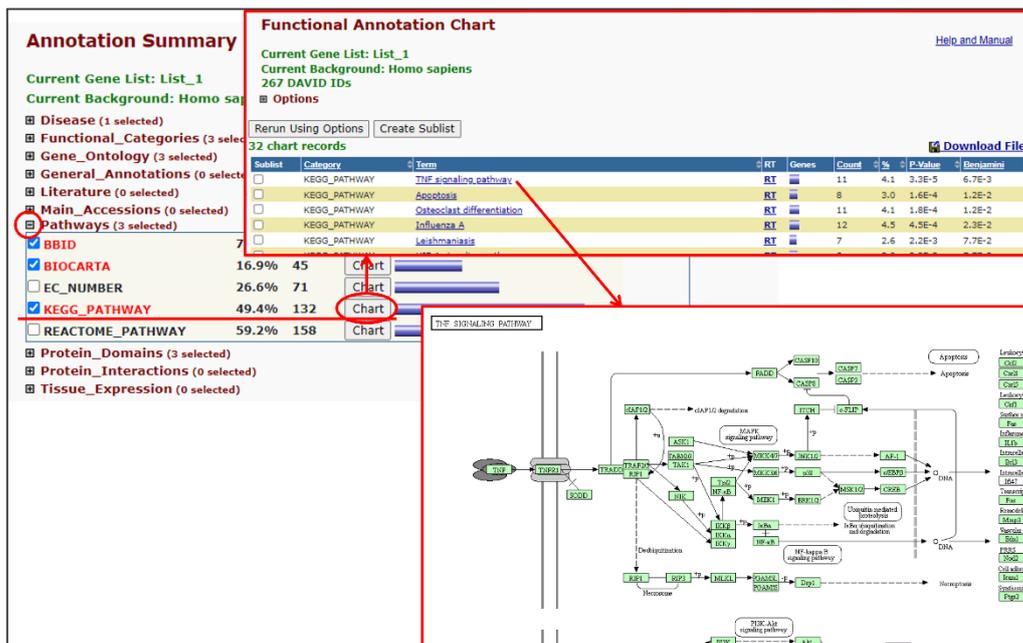


그림 2-8. DAVID tool : exploring Pathway analysis result

DAVID 분석은 input 한 유전자들이 유의하게 관련되는 GO, pathway 등을 분석하기에 유용한 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 분석 결과를 내 컴퓨터에 저장하려면, chart 를 클릭해서 나온 결과창에서 Download File 링크를 마우스 오른쪽 버튼을 클릭한 후 다른 이름으로 저장을 선택하면 DAVID results 파일을 다운로드 받을 수 있다 (그림 2-9. 다운로드 받은 DAVID 결과 파일로 그래프 작성하는 방법은 '2-2. ExDEGA GraphicPlus 를 이용한 DAVID 결과 그래프 작성'에 설명되어 있다.

\* 주의사항

: internet explorer 를 이용할 경우 다른 이름으로 저장 버튼이 보이지 않기 때문에, Chrome 을 이용하여 분석하기를 권장한다.

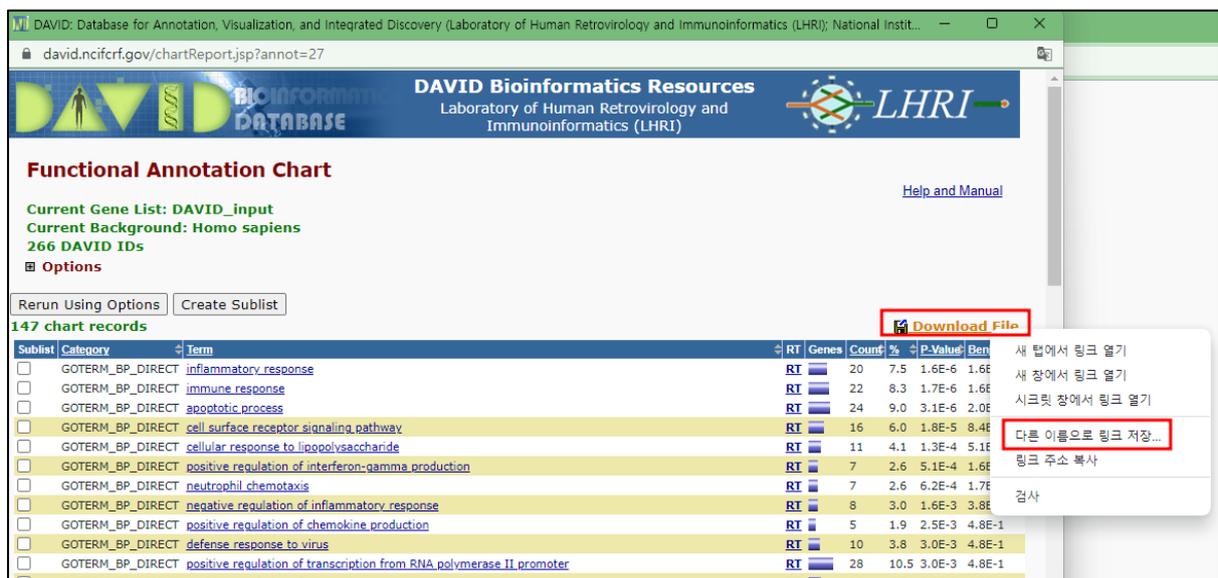


그림 2-9. DAVID data download

DAVID 에서는 유전자 2 개 이상, EASE score 0.1 이하를 default 로 분석하여 이 기준에 적합한 결과를 보여준다. option 에서 이 기준을 조정하여 리스트를 더 볼 수 있다. DAVID 분석 결과의 각 항목은 DAVID 홈페이지의 Help and Tool Manual 에 자세히 설명되어 있다(그림 2-10)



**Functional Annotation Chart**  
 Current Gene List: demolist1  
 Current Background: Homo sapiens  
 171 DAVID IDs

**Options**  
 Count Threshold: 2    EASE Threshold: 0.1    # of Records Displayed: 1000

Rerun Using Options    Create Sublist    Download File

Sublist	Category	Term	RT	Genes	Count	%	P-Value
<input type="checkbox"/>	SP_PIR_KEYWORDS	signal	RT		47	27.5%	3.0E-10
<input type="checkbox"/>	SP_PIR_KEYWORDS	glycoprotein	RT		51	29.8%	4.9E-8
<input type="checkbox"/>	GOTERM_CC_ALL	extracellular region	RT		32	18.7%	1.1E-7
<input type="checkbox"/>	SP_PIR_KEYWORDS	alternative splicing	RT		49	28.7%	6.4E-6
<input type="checkbox"/>	SP_PIR_KEYWORDS	chromoprotein	RT		7	4.1%	1.1E-5
<input type="checkbox"/>	SP_PIR_KEYWORDS	direct protein sequencing	RT		33	19.3%	1.2E-5
<input type="checkbox"/>	SP_PIR_KEYWORDS	phosphorylation	RT		31	18.1%	1.6E-5
<input type="checkbox"/>	UP_SEQ_FEATURE	signal peptide	RT		47	27.5%	3.7E-5
<input type="checkbox"/>	SP_PIR_KEYWORDS	metalloprotein	RT		8	4.7%	4.7E-5
<input type="checkbox"/>	GOTERM_BP_ALL	response to chemical stimulus	RT		14	8.2%	6.1E-5

Annotations:  
 - Gene list and population background being analyzed: points to 'Current Gene List' and 'Current Background'.  
 - Minimum number of genes for the corresponding term: points to 'Count Threshold'.  
 - Maximum EASE Score/P-Value: points to 'EASE Threshold'.  
 - Maximum number of record per page: points to '# of Records Displayed'.  
 - Original database/resource where the terms orient: points to 'Category' column.  
 - Enriched terms associated with your gene list: points to 'Term' column.  
 - Related Term Search: points to 'RT' column.  
 - Genes involved in the term: points to 'Genes' column.  
 - Percentage, e.g. 14/171=8.2% (involved genes/total genes): points to '%' column.  
 - Modified Fisher Exact P-Value, EASE Score. The smaller, the more enriched.: points to 'P-Value' column.

그림 2-10. DAVID Help and Tool Manual

## 2-2. ExDEGA GraphicPlus 를 이용한 DAVID 결과 그래프 작성

ExDEGA 레포트의 ExDEGA Graphic Plus Start를 클릭하여 프로그램을 실행한다 (그림2-11).

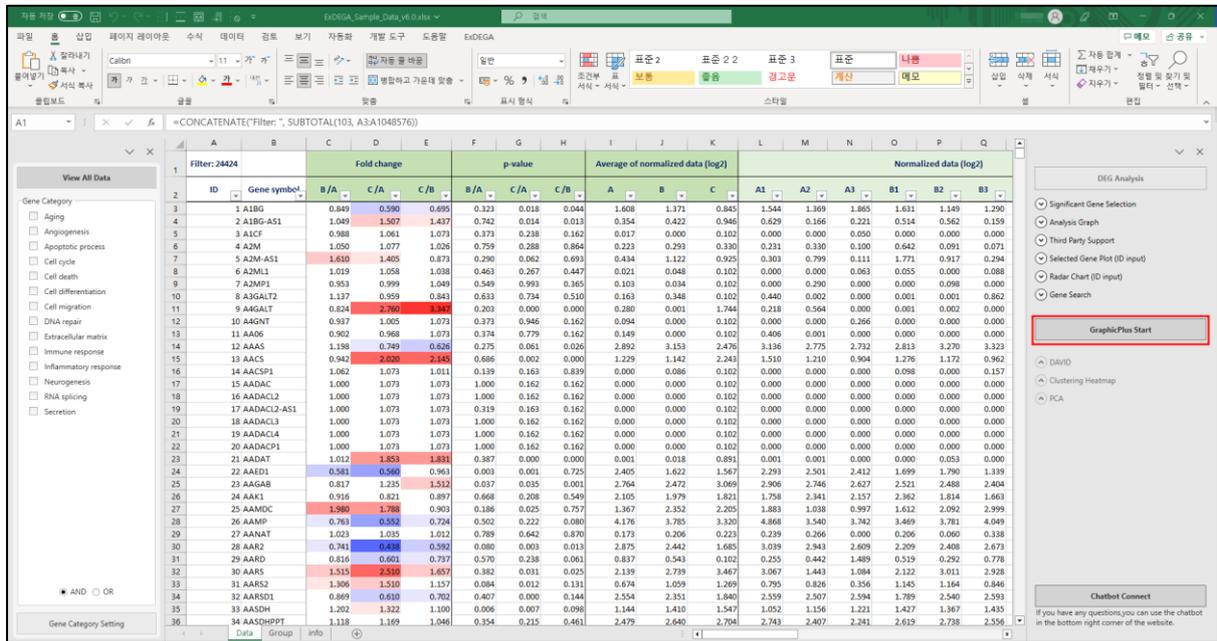


그림 2-11. Execute ExDEGA GraphicPlus

메인 화면의 4개 탭 중 'DAVID' 탭에서 DAVID Graphic 분석을 수행할 수 있다. DAVID Graphic 분석 창은 그림 2-12과 같다.

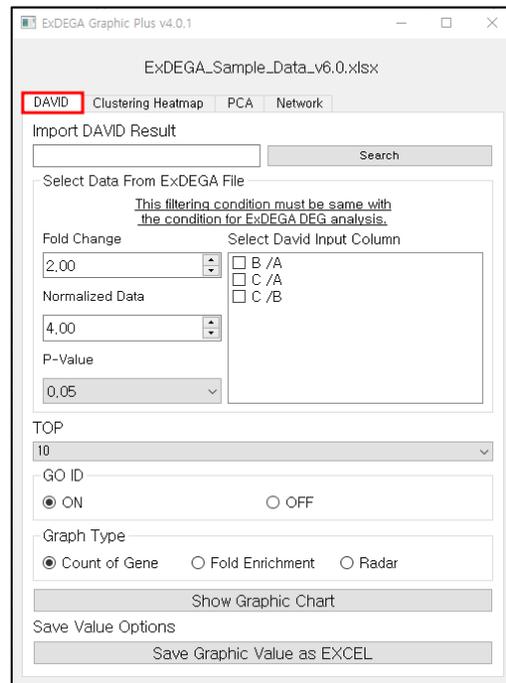


그림 2-12. Select DAVID tab to execute the DAVID Graphic Analysis

DAVID 분석 결과를 그래프로 제작하려면 DAVID result 파일과 DAVID input 이 필요하다. DAVID result 은 '2-1. DAVID 분석 틀을 이용한 Functional Annotation 분석'에서 저장한 DAVID 결과 파일이다. DAVID result 파일 안에 있는 내용은 그림 2-13 과 같다. 이 파일에 있는 Term, Count, P-value, Fold Enrichment 항목을 이용하여 그래프가 제작된다.

Category	Term	Count	%	PValue	Genes	List Total	Pop Hits	Pop Total	Fold Enrichment	Bonferroni	Benjamini	FDR
GOTERM_BP_DIRECT	GO:0055114~oxidation-reduction process	36	16.43836	2.70E-14	D3YZE4, P5	203	676	18082	4.743580027	3.05E-11	3.05E-11	4.33E-11
GOTERM_BP_DIRECT	GO:0006631~fatty acid metabolic process	19	8.675799	1.50E-13	P04117, P5	203	156	18082	10.84874321	1.70E-10	8.48E-11	2.40E-10
GOTERM_BP_DIRECT	GO:0006635~fatty acid beta-oxidation	12	5.479452	1.43E-12	Q9DC50, C	203	44	18082	24.29287953	1.62E-09	5.40E-10	2.30E-09
GOTERM_BP_DIRECT	GO:0008152~metabolic process	28	12.78539	2.10E-12	P19157, P5	203	463	18082	5.386758025	2.37E-09	5.93E-10	3.36E-09
GOTERM_BP_DIRECT	GO:0006629~lipid metabolic process	24	10.9589	1.96E-09	P51174, Q5	203	459	18082	4.657458386	2.22E-06	4.43E-07	3.14E-06
GOTERM_BP_DIRECT	GO:0006810~transport	48	21.91781	3.28E-08	P04117, P2	203	1822	18082	2.346622831	3.71E-05	6.18E-06	5.26E-05
GOTERM_BP_DIRECT	GO:0006637~acyl-CoA metabolic process	7	3.196347	1.05E-06	Q8VCT4, C	203	31	18082	20.1134594	0.001191	1.70E-04	0.001689
GOTERM_BP_DIRECT	GO:0070527~platelet aggregation	7	3.196347	3.70E-06	Q9Z1Q5, P	203	38	18082	16.40834846	0.004178	5.23E-04	0.005935
GOTERM_BP_DIRECT	GO:0006754~ATP biosynthetic process	5	2.283105	1.13E-04	Q03265, D	203	23	18082	19.36388948	0.120415	0.014155	0.181743
GOTERM_BP_DIRECT	GO:0015671~oxygen transport	4	1.826484	1.56E-04	P02089, P0	203	10	18082	35.62955665	0.161483	0.017458	0.249389
GOTERM_BP_DIRECT	GO:0006749~glutathione metabolic process	6	2.739726	2.12E-04	P48774, P1	203	49	18082	10.90700714	0.213205	0.021563	0.339389
GOTERM_BP_DIRECT	GO:0051791~medium-chain fatty acid metabol	3	1.369863	3.70E-04	Q9DC50, C	203	3	18082	89.07389163	0.342104	0.034291	0.591879

그림 2-13. DAVID output file

그래프를 제작하기 위해서는 이전 DAVID input 파일을 만들 시, 적용했던 선별조건을 그대로 설정해야 한다 (그림 2-14).

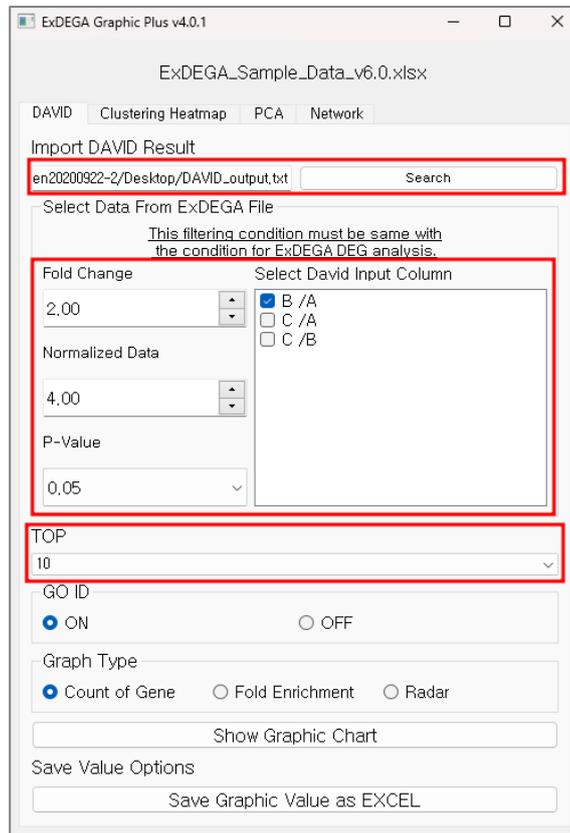


그림 2-14. Steps for Create DAVID Analysis Graph

우선, import DAVID Result 에서 DAVID result 파일을 선택한다. 그리고 아래 선별조건에서 DAVID 분석 시 사용된 input 파일과 동일하도록 ExDEGA 레포트에 적용한 Fold Change 값, Normalized Data (log2)값, P-Value 값 및 비교조합 선택을 세팅해야 한다. Input 파일과 마찬가지로 비교조합은 1 개만 선택한다. 사용자가 별도의 기준을 적용하지 않을 경우, Fold Change 2.00, Normalized Data (log2) 4.00, P-Value 0.05 로 자동 적용된다. TOP 에서는 그래프 제작에 사용될

결과 리스트의 수이다. 이는 선택한 비교조합이 DAVID result 파일의 결과값과 함께 계산되어 DAVID 분석 그래프를 그릴 때 DAVID result 파일에서 상위 n 개의 리스트를 대상으로 그래프를 제작하는 옵션이다. 2 개에서 35 개 까지의 옵션을 선택 가능하며, 사용자가 별도의 기준을 적용하지 않을 경우 이 값은 자동으로 10 으로 설정된다.

위의 과정을 거친 뒤, Graph Type 에서 3 종류의 그래프를 제작할 수 있다 (그림 2-15).

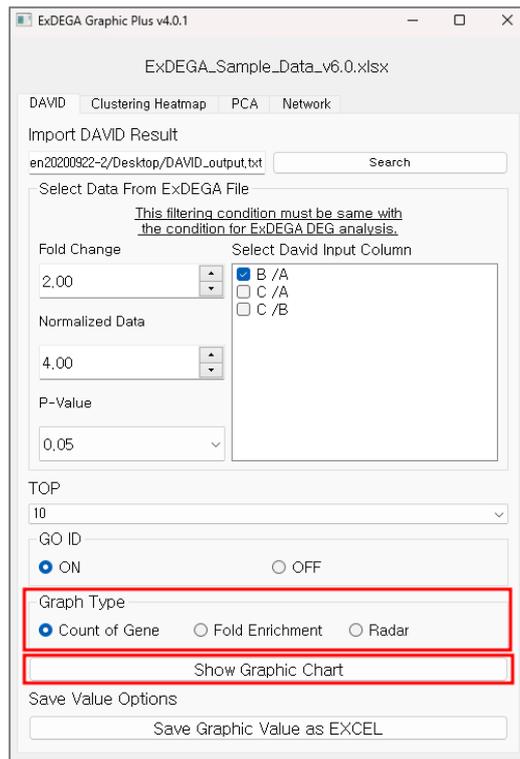


그림 2-15. Result of Graphic Chart

Count of Gene 을 선택하고 Show Graphic Chart 를 누르면 각 GO (or pathway)에서 발현이 증가하는 유전자, 감소하는 유전자 수가 그래프로 작성된다(그림 2-16).

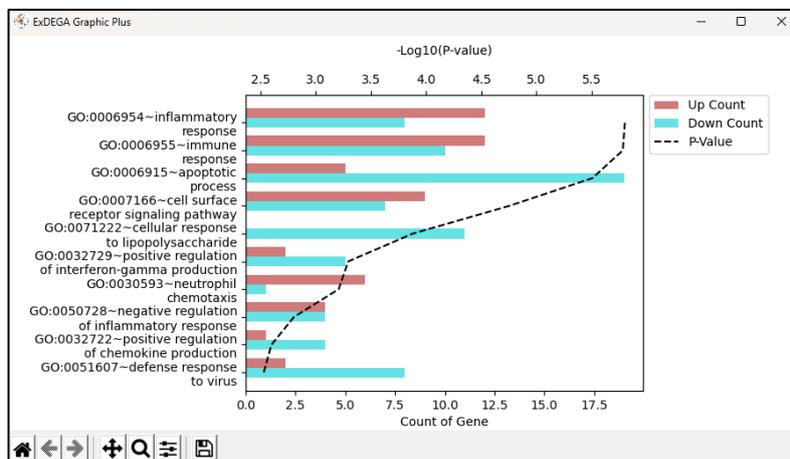


그림 2-16. Result of Graphic Chart (Count of Gene)

Fold Enrichment를 선택하고 Show Graphic Chart를 누르면 각 GO (or pathway)의 p-value, Fold enrichment 값으로 그래프가 작성된다 (그림2-17).

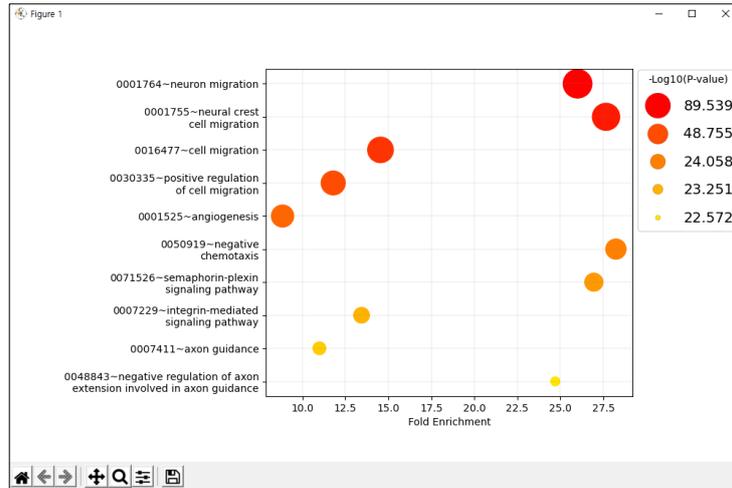


그림 2-17. Result of Graphic Chart (Fold Enrichment)

Radar를 선택하고 Show Graphic Chart를 누르면 각 GO (or pathway)의 Fold enrichment 값으로 그래프가 작성된다 (그림2-18).

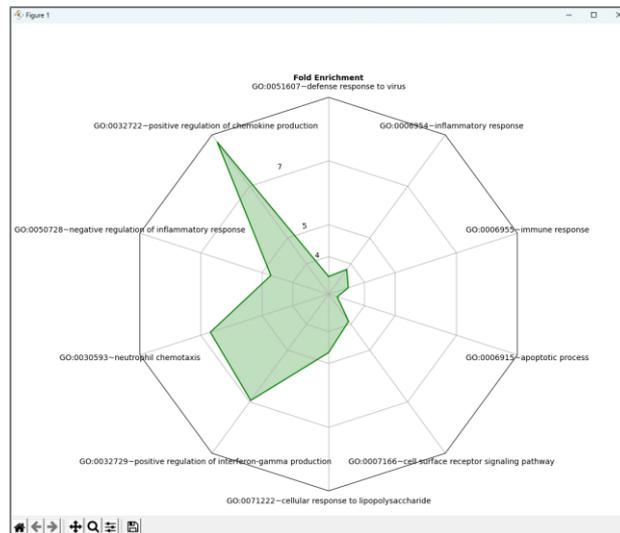


그림 2-18. Result of Graphic Chart (Radar)

GO ID에서 ON을 선택하면 DAVID의 고유번호가 앞에 붙고, OFF를 선택하면 고유번호 없이 Term만 출력된다 (그림 2-19).

GO ID

ON  OFF

그림 2-19. Result of Graphic Chart (GO ID)

마지막으로, Save Graphic Value를 누르면 그래프에 사용된 값들이 Excel형식으로 저장된다 (그림 2-20).

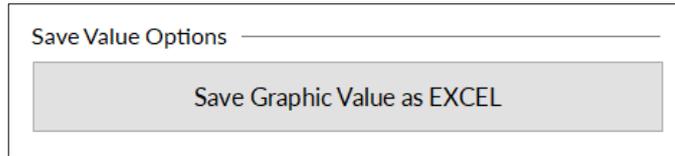


그림 2-20. Save DAVID Graphic Values as Excel

이 결과 파일은 DAVID Result (그림 2-21), Gene List & Fold Change (그림 2-22)의 2개 탭으로 이루어져 있다. DAVID Result 탭에서는 그래프 제작에 이용한 DAVID Graphic 분석에 대한 전체적인 결과값이, Gene List & Fold Change 탭에서는 설정한 Top Count에서 유의미한 결과값을 보인 Gene List와 각 gene의 fold change 결과값이 표시된다.

	A	B	C	D	E	F	G	H
		Total	Number of Molecules	Up Molecules	Down Molecules	Fold Enrichment	Log10(P-value)	
1								
2	0001764--neuron migration	64	70	31	39	26.00451467268623	99.5374064847625	
3	0001755--neural crest cell migration	44	39	20	19	27.65934742458445	50.73585351965816	
4	0016477--crest migration	118	55	29	26	14.54489803726518	48.75523653693835	
5	0030335--positive regulation of cell migration	127	48	24	24	11.794173583832496	37.14452140747271	
6	0001525--angiogenesis	134	38	20	18	8.849297530406657	24.057940473511444	
7	0050919--negative chemotaxis	21	19	8	11	28.23347307320219	23.97814634517612	
8	0071526--semaphorin-plexin signaling pathway	22	19	11	8	26.95013338805664	23.251333489938407	
9	0007229--integrin-mediated signaling pathway	65	28	11	17	13.442333738496266	23.055541647085267	
10	0007411--axon guidance	98	51	16	15	10.9928117566180998	22.571872514053652	
11	0048843--negative regulation of axon extension involved in axon guidance	24	19	11	8	24.70428893905192	22.01631832716083	

그림 2-21. Result of Saved Graphic value (Tab 1 : DAVID Result)

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z				
	0001764--neuron mig	0001764 FC	neural crest cell	0001755 FC	4477--crest migr	0016477 FC	live regulation of	0030335 FC	1525--angiogeni	0001525 FC	3--negative chel	0050919 FC	thorin-plexin sig	0001526 FC	lin-mediated sig	0007229 FC	4111--axon guid	0007411 FC	of axon exten	0048843 FC										
1	0	MF2C	0.859788263 NRP2	0.96702913 EFNA1	1.08301892 CCL3	0.14271802 NRP2	0.96702913 NRP2	0.96702913 FLNA	0.65061971 ITGB4	1.201288745 WNT5A	1.014481086 WNT5A	1.014481086	1.014481086																	
2	1	TUBB2B	0.840241504 EDN3	0.955311616 C5P4	1.005792733 DRD1	1.023817663 FGF1	0.81416482 SIRT1	1.027894843 SEMA5A	0.975121774 ITGB5	0.913818966 CDK1R1	0.825573826 NRP1	0.846690016																		
3	2	PAK6	0.950597863 NRYN1	1.02028219 JSC4	1.252164654 PDGFR	1.352585505 CAU1	1.020423852 SULT3	0.99650173 SEMA4A	1.090783042 ITGB2	0.730426075 EFNA1	1.020013992 SEMA5A	0.975121774																		
4	3	PEX5	0.991872425 ERBB4	1.00503916 RAB1A	1.076731668 FZL1	2.677534126 ACVRL1	0.850568183 APOA1	0.977555888 SEMA6C	1.178486305 ITGB3	0.630547126 PAX6	0.950597863 SEMA6A	1.098783042																		
5	4	GUJ1	0.968055913 KITLG	0.99093864 SDC2	1.143314051 EDN1	0.967391945 EFNA1	1.08301892 SEMA4G	1.145906979 SEMA6D	1.014010591 ITGB1	0.689748285 CXCL12	0.846461 SEMA6C	1.178486305																		
6	5	PGF3	1.045138762 GDNF	1.006143849 TGFBI	0.641669032 FER	1.118601907 WAP22	1.10090002 SEMA6F	1.149019447 SEMA4G	1.145906979 ITGAM	0.95845261 ITGB2	0.936951718 SEMA6D	1.014010591																		
7	6	PRKCI	0.920776102 SOXB1	0.920269904 TGFBI2	0.957342528 CXCL12	0.846461 EPHB3	1.102765533 SEMA4G	1.015347482 SEMA6F	1.149019447 ITG1	0.483001836 SHH	0.97912488 SEMA4G	1.145906979																		
8	7	CTNNA1	0.872766289 SHH	0.997912483 CUL3	0.997912483 CUL3	0.641669032 PTEN	1.564172914 ITGAV	0.88173207 SEMA3G	1.018347482 CUL3	1.303301053 ITGB2	0.97342618 SEMA6F	1.149019447																		
9	8	PEY7	0.743890745 SEMA5A	0.975121774 FAMA83D	0.980343251 SEMA5A	0.975121774 PRKX	0.867769185 SEMA6F	0.98236233 SEMA6F	0.98236233 SEMA6F	0.98236233 PTX2	0.986780604 LAMB2	0.868051993 SEMA3G	1.015347482																	
10	9	NDE1	1.536416938 EDNBB	0.933847822 PDN1	1.074524771 PTX2	0.986776902 EPHB1	1.799665001 SEMA5E	0.98170739 RAC1	1.115886798 APOA1	0.977555888 RGS10	1.10500668 SEMA6F	0.98236233																		
11	10	POMGN2	0.70209099 SEMA3G	1.018347482 NDE1	1.536416938 SEMA3G	1.015347482 NDE1	1.536416938 SEMA3G	1.218190273 SEMA4D	1.027867246 SEMA3E	0.9810739 PTX2B	1.22272732 GATA3	0.9810739																		
12	11	GATA3	0.83842688 SEMA7A	0.74064196 S1PR1	1.027926194 SEMA7A	0.74064196 PTX2	0.986780604 SEMA4C	0.910711274 SEMA3D	1.027867246 ITGB8	1.060729002 RAC1	1.027867246 ITGB8	1.060729002																		
13	12	MARGR2	1.019304094 SEMA3F	0.98236233 TRIAM1	0.958887826 SEMA3F	0.98236233 S1PR1	1.207926194 RHOA	0.806897702 SEMA4C	0.910711274 SEMA7A	0.74064196 GNB2	1.01494719 SEMA4C	0.910711274																		
14	13	NRP2	0.99113908 SEMA5E	0.9810739 RHOA	1.058887826 SEMA3F	0.88173207 ITGAV	0.88173207																							
15	14	ASPM	1.160999121 GNB2	1.01494719 PTPRX	1.20316042 SEMA3E	0.9810739 PTX2B	1.22272732 SEMA3C	0.95972267 SEMA3C	0.95972267 ITGB7	0.95972267 ITGB7	0.95972267 ITGB7	0.95972267																		
16	15	PTX2	1.03902638 SEMA3D	1.027867246 ARC	0.913978888 SEMA3D	1.027867246 ITGAV	0.88173207 ITGAV	0.960108825 SEMA3B	0.960108825 SEMA3B	0.960108825 ILK	0.944779977 USP39	1.073246406 USP39																		
17	16	VRRB2	1.029189125 SEMA3C	0.95972267 CDC88A	0.97893254 SEMA3C	0.95972267 ITGAV	1.000102039 SEMA4D	1.569040566 SEMA4D	1.569040566 SEMA4D	1.569040566 PCEN1B	0.844611124 PLR19	0.959040566 SEMA4D	1.569040566																	
18	17	TWIST1	0.94240558 SEMA3B	0.960108825 PTPRF	1.070568168 SEMA3B	0.960108825 PK3CA	0.91191364 SEMA3A	0.926246444 SEMA3A	0.926246444 SEMA3A	0.926246444 SYK	1.273158465 MATN2	0.930952064 SEMA3A	0.926246444																	
19	18	MATN2	0.98052064 SEMA3A	0.926246444 LYN	1.335519884 LAMB1	1.036702746 RHOB	0.667430259 SEMA4A	1.205043398 SEMA4A	1.205043398 SEMA4A	1.205043398 SEMA4A	1.068664471 EFNB1	1.090898862 SEMA5A	1.205043398																	
20	19	PHOX2B	1.008487428 SOD10	1.099415304 EPHB2	1.028979388 PHSR1	0.937896252 SHC1	0.970988489	None	None	ITGA2	1.048418758 EPHB2	1.028979388																		
21	20	GURS	1.189897588 RET	0.947469347 MMP14	0.66833118 EGR1	0.979826239 ANGPT1	1.565626918	None	None	ITGA4	0.9771768 VAV1	0.858406691																		
22	21	SATB2	0.9558416 SEMA3E	0.720857676 PALLD	0.855448416 ELP3	1.189397588 SOD1X1	0.999225997	None	None	ITGA4	0.638468652 PTPRO	1.265892246																		
23	22	VAV1	0.958406691 EFNB1	1.090589862 CDK5	0.971637616 RET	0.94769347 FGF2	1.039238978	None	None	VAV1	1.57432609 PDZ2	0.936392374																		
24	23	PCMY1	0.920266647 LEY1	1.480268866 AUP	0.639686974 FLT1	1.130235841 ANGPT2	1.316181839	None	None	PTPN11	0.702625411 SEMA6A	1.098783042																		
25	24	CDK5	0.971637616 LEY1	0.997824708 HES1	0.942500111 RYR	1.335519884 FN1	0.847631869	None	None	ITGAV	1.02515384 EPHA8	0.984002382																		
26	25	DDIT4	0.07282812 CORD1C	1.090236772 NCK2	0.719848425 TGFBR1	1.589486185 TGFBR1	1.589486185	None	None	ITGA5	1.157167212 SEMA6F	1.149019447																		
27	26	ACLY1	1.01048192 SEMA4A	1.098783042 SDC1	1.074789832 PCOLCE	1.18712543 ACOR3	1.196800016	None	None	LAMB4	1.179186308 CNTN2	1.028627292																		
28	27	NDEL1	1.795200611 SEMA6C	1.178486305 NDEL1	1.795200611 PK3CD	2.067214458 MYH9	1.043748236	None	None	ITGB1BP1	0.833151231 RELN	1.007248142																		
29	28	PSEN1	1.255145272 HIF1A	0.802931584 SH3BP1	1.058986232 SUN2	0.881667121 LEP	1.07092241	None	None	None	APPB2	0.920874729																		
30	29	CTNNA	0.93950018 SEMA4G	1.145906979 LAMC1	1.068478383 LEF1	1.483038666 VEGFC	0.938504698	None	None	None	CHL1	1.015587263																		
31	30	HELN	1.007248142 SEMA4D	1.014010591 S12B	1.00242763 HGF	1.192310905 HIF1	0.922919524	None	None	None	MYH10	1.445797351																		
32	31	CHL1	1.015587263 LAMB4	1.179186308 LAMB4	1.179186308 MMP14	0.86833138 VEGFA	0.97186658	None	None	None	None	None																		
33	32	BR3A	0.937120933 SEMA4F	1.145906979 UPPRX	1.160889157 ZNFX1	0.233511999 PECAM1	1.806447711	None	None	None	None	None																		
34	33	CDK1R1	0.825573826 SEMA4C	0.910711274 PEAK1	0.940020021 NOTCH1	2.85167016 AKAP1	0.729179237	None	None	None	None	None																		
35	34	JDOCA9B	0.861983043 SEMA4B	2.409734569 PTPC7	1.096392148 SEMA6C																									

### 3. Clustering heatmap analysis (ExDEGA GraphicPlus)

Hierarchical Clustering Heatmap은 연구자가 선택한 유전자의 발현 유사성을 기반으로 Sample 간의 유사성, 유전자 간의 유사성을 판단할 때 사용한다. 본 챕터에서는 Heatmap (개별 색상의 직사각형 데이터 행렬)과 Dendrogram (계층적 클러스터링)을 합쳐 Hierarchical Clustering Heatmap을 그리는 방법을 설명한다.

Hierarchical Clustering Heatmap은 ExDEGA GraphicPlus를 이용하여 분석할 수 있다. ExDEGA GraphicPlus를 이용하기 위해서 먼저 해당 프로그램을 연다. ExDEGA 레포트의 ExDEGA GraphicPlus Start 버튼을 클릭하여 프로그램을 동작시킨 뒤, Clustering Heatmap 탭을 클릭하여 준비한다(그림 3-1).

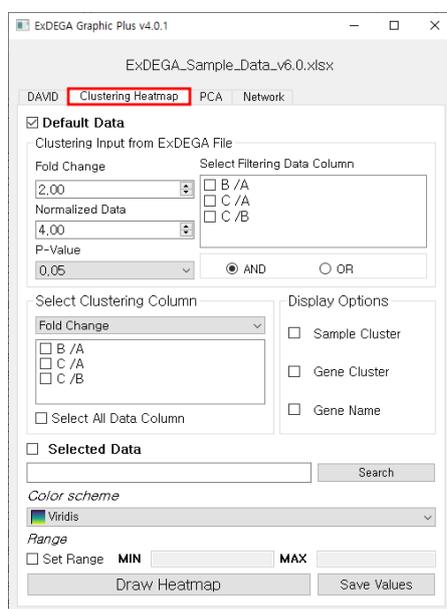


그림 3-1. Select Clustering Heatmap tab

Clustering Heatmap을 제작을 위해 input Data는 크게 두가지 형식으로 이용할 수 있다. GraphicPlus가 인식한 ExDEGA report의 전체 data를 이용하는 Default Data와 Third Party Support에서 export한 Selected Data가 있다.

Default Data를 이용하는 경우, 1~7의 과정을 거쳐 Clustering Heatmap을 그리게 된다.(그림 3-2) 1에서는 인식된 ExDEGA 레포트를 바탕으로, DEG 분석 기준 및 DEG 분석 기준을 적용할 Fold Change 샘플 그룹을 지정하게 된다. 사용자가 별도의 기준을 적용하지 않을 경우, Fold Change 2.00, Normalized Data 4.00, P-Value 0.05로 자동 적용되며 Fold Change 샘플 그룹은 AND/OR 기능을 통해 여러 개의 샘플에 기준을 복수 적용할 수 있다.

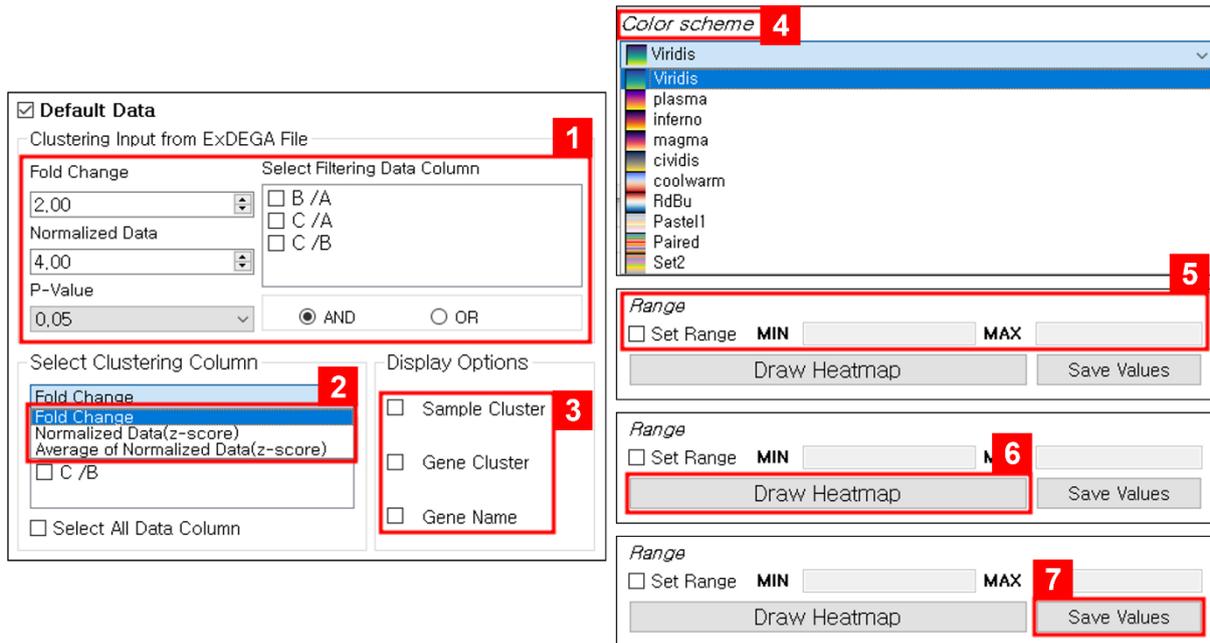


그림 3-2. Steps for Create Clustering Heatmap

2에서는 1에서 적용한 기준을 바탕으로, 실제 Clustering Heatmap을 그리고자 하는 데이터 컬럼을 선택한다. ExDEGA 레포트 유형(Group, Single)에 따라 Group 레포트에서는 Fold Change, Normalized Data(z-score), Average of Normalized Data(z-score) 타입이, Single 레포트에서는 Fold Change, Normalized Data(z-score) 타입이 제공된다. 반영하고자 하는 타입을 선택한 후, 그림 3-3와 같이 타입에 해당하는 샘플 그룹을 선정할 수 있다. 이 때, 하단의 'Select All Data Column'을 선택하면 해당하는 타입에 존재하는 모든 샘플 그룹을 동시에 선택/선택 해제할 수 있다.

단, 사용자가 Normalized Data(z-score), Average of Normalized Data(z-score) 옵션을 선택하더라도, 3에서 고른 샘플 그룹의 수가 2개 이하인 경우 z-score가 적용되지 않고 ExDEGA 레포트에 표현된 데이터 값을 그대로 적용하여 Clustering heatmap이 작성되니 이 점을 유의해야 한다.

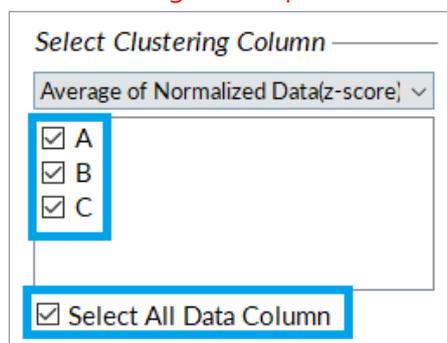


그림 3-3. Select specific sample group for composition of clustering heatmap

3에서는 Clustering Heatmap 옵션에 적용할 Display option을 선택할 수 있다(그림 3-4). Sample cluster를 선택하면 발현이 유사한 비교조합 또는 샘플 간의 dendrogram이 작성된다. Gene cluster를 선택하면 발현이 유전자 간의 dendrogram이 작성된다. Gene name을 표시하면 Raw에 해당하는 gene symbol이 표시된다. 단, 입력한 유전자가 80개 이상일 경우에는 gene symbol을 표시할 수 없다.

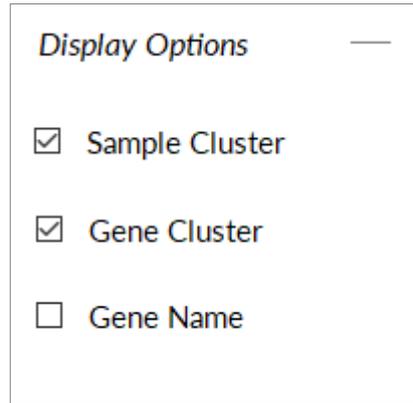


그림 3-4. Set a Display option for clustering heatmap

4에서는 Color scheme는 Heatmap의 색상을 설정할 수 있다. 10가지의 기본 옵션을 제공한다. 사용자가 별도로 옵션을 선택하지 않을 경우, 'Viridis' 옵션이 기본으로 설정된다.

5의 Set Range는 데이터 표현 범위를 설정하는 것으로 선택하지 않으면 input file의 최소값과 최대값으로 자동 설정된다. Set Range 앞의 체크박스를 선택하면 MIN(최소값), MAX(최대값) 입력 박스가 활성화되며, 사용자가 원하는 최소값/최대값을 설정할 수 있다(그림 3-5).



그림 3-5. Set a min-max option for clustering heatmap

6의 Draw Heatmap 버튼을 누르면 Hierarchical Clustering Heatmap 결과 창이 생성된다(그림3-6). 각 heatmap의 위쪽에 표시된 dendrogram은 비교조합 또는 샘플 간의 발현 유사성(sample cluster)을 표시한 결과이다. 왼쪽에 표시된 dendrogram은 유전자 간의 발현 유사성(gene cluster)을 표시한 결과이다. 가깝게 묶일수록 발현이 유사한 것이다.

또한 그림 3-6. a heatmap은 정제한 전체 발현값을 대상으로 제작된 heatmap이며, 그림 3-6. b heatmap은 min-max 값을 적용하여 제작된 heatmap이다. (각 그래프 왼쪽 상단의 legend 값에서 차이점을 확인할 수 있다.)

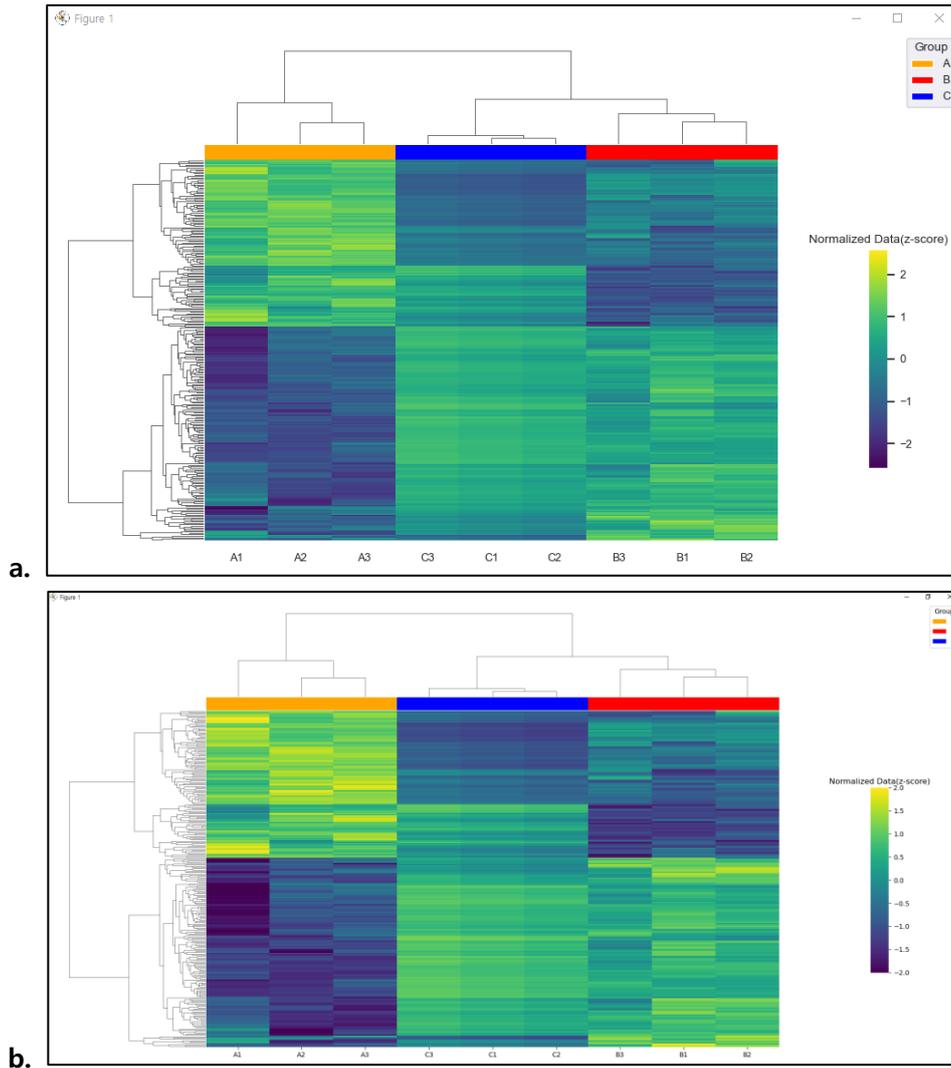


그림 3-6. Clustering heatmap result

7의 Save Values 버튼을 누르면 Clustering Heatmap을 만들 때 사용한 value값이 엑셀로 저장된다. 이 때, Save Values는 Clustering Heatmap을 제작하지 않았더라도 필터링 기준, 필터링 기준을 적용할 샘플그룹, Clustering Heatmap을 적용할 샘플 그룹들을 선택하면 결과값을 받아볼 수 있다. 각각의 값을 지정하고 Save Values 버튼을 누르는 경우, 조건을 확인하기 위한 결과창이 뜨며(그림 3-7) 이 결과창에서 'OK'를 누를 때 값이 사용자가 정한 위치에 .xlsx 형식으로 저장된다.

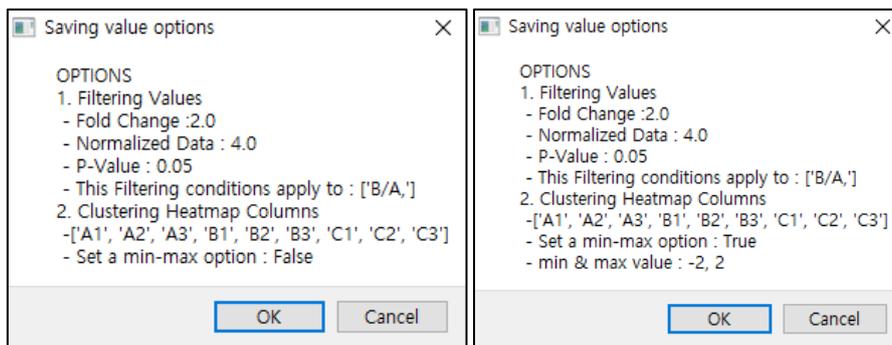


그림 3-7. Confirm all conditions for saving values used to create clustering heatmap



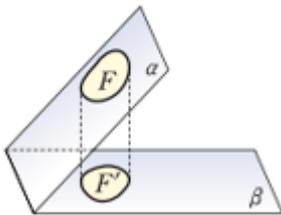
## 4. Principal component analysis (ExDEGA GraphicPlus)

본 챕터에서는 PCA 에 대한 이론과 PCA 2D/3D 를 그리는 방법에 대해 설명한다. PCA 는 Sample 간의 발현 유사성을 확인하기 위한 목적으로 Clustering Heatmap 과는 다르게 Sample 내의 유전자 전체 발현값을 기반으로 분석된다.

### 4-1. PCA (Princial Component Analysis) 이론

#### PCA Essential Description.

PCA 는 주성분 분석의 준말로 고차원 데이터를 정사영(구조 유지, 차원 감소) 시켜 저차원 데이터로 차원을 축소하는 알고리즘이다.



$F'$ 은  $F$ 의  $\beta$  위로의 정사영(=  $F'$ 은  $F$ 를  $\lambda$  만큼 표현한다)

PCA 는 앞서 설명한 바와 같이 Sample 에 속해 있는 전체 유전자 발현을 대상으로 Sample 간의 유사성을 확인하려는 목적으로 분석을 진행한다. Human RNA-seq 기준으로 각 Sample 에는 약 24,000 개 이상의 유전자가 포함된다. 24,000 개의 변수가 생긴다는 말과 동일하기 때문에 Sample 간의 유사성을 파악하기 힘들다. PCA 알고리즘을 통해 방향 변화 없이 variance 최대(선형 변환)가 되는  $\lambda$  값(eigenvalue)을 계산한다.  $\lambda$  값은 sample 의 개수만큼 나오게 되고 주성분 비율은  $\lambda$  값의 전체 합에서 해당  $\lambda$  값이 차지하는 비율이다.

예시 : i. PC1( 60%), PC2(30%) => 주성분 1, 2 로 데이터의 90% 표현

ii. PC1(50%), PC2(30%), PC3(10%) => 주성분 1,2,3 으로 데이터의 90% 표현

**PCA additional Info.**

다음은 공분산을 이용하여 PCA 를 계산하는 방법을 기술한다.

2 개의 Sample(열)에서 5 개의 유전자(행)에 발현을 관찰하여 행렬(D)로 만든다.

$$D = \begin{bmatrix} 100 & 111 \\ 152 & 45 \\ 19 & 33 \\ 22 & 31 \\ 27 & 10 \end{bmatrix} \in \mathbb{R}^{5 \times 2}$$

유전자 각각의 변동을 알기 위해 행렬(D)에서 평균 값(기준점)을 뺀으로써 행렬(X) 구할 수 있다. (평균 값을 원점으로 데이터가 확장된 정도)

$$X = D - m_{\text{ean}}(D) = \begin{bmatrix} 100 & 111 \\ 152 & 45 \\ 19 & 33 \\ 22 & 30 \\ 27 & 10 \end{bmatrix} - \begin{bmatrix} 64 & 46 \\ 64 & 46 \\ 64 & 46 \\ 64 & 46 \\ 64 & 46 \end{bmatrix} = \begin{bmatrix} 36 & 65 \\ 88 & -1 \\ -45 & -13 \\ -42 & -16 \\ -37 & -36 \end{bmatrix}$$

Transpose X 와 X 를 곱함으로써 데이터의 내적 값을 구할 수 있고 (유전자의 개수-1)을 나눔으로 데이터 공분산 행렬을 구할 수 있다.

$$X^T X = \begin{bmatrix} 36 & 88 & -45 & -42 & -37 \\ 65 & -1 & -13 & -16 & -36 \end{bmatrix} \begin{bmatrix} 36 & 65 \\ 88 & -1 \\ -45 & -13 \\ -42 & -16 \\ -37 & -36 \end{bmatrix} = \begin{bmatrix} 14198 & 4841 \\ 4841 & 5947 \end{bmatrix}$$

$$\frac{X^T X}{n-1} = \begin{bmatrix} 14198 & 4841 \\ 4841 & 5947 \end{bmatrix} / 4 = \begin{bmatrix} 3549.5 & 1210.25 \\ 1210.25 & 1486.74 \end{bmatrix}$$

공분산은 공통되게 움직이는 정도를 표현한 것으로 이제 방향(Eigenvalue)를 계산해야 한다. Eigenvalue(  $\lambda$  )는 nonzero solution vector K (Eigenvector) 가 존재해야 한다는 조건이 있다. 따라서,  $AK = \lambda K$  를 만족하여야 하며  $K(A - \lambda I) = 0$  과 같다. 이때  $K(A - \lambda I)$  이 역 행렬이 존재 한다면  $K = 0$  이기 때문에 조건에 모순이 된다.

$$\therefore \det(A - \lambda I) = 0$$

위의 예제를 공식에 따라 대입해보면 다음과 같다.

$$\det \left( \begin{bmatrix} 3549.5 - \lambda & 1210.25 \\ 1210.25 & 1486.74 - \lambda \end{bmatrix} \right) = 0 \rightarrow (3549.5 - \lambda)(1486.74 - \lambda) - 1210.25^2 = 0$$

$$\therefore \lambda_1 = 275.101, \lambda_2 = -5311.341$$

## 4-2. ExDEGA GraphicPlus 를 이용한 PCA 분석 방법

ExDEGA GraphicPlus를 이용하려면, ExDEGA 레포트 우측 하단의 ExDEGA Graphic Plus Start 버튼을 클릭하여 ExDEGA Graphic Plus를 활성화시킨다.

메인 화면의 3 개 탭 중 'PCA' 탭에서 PCA 분석을 수행할 수 있다. PCA 분석 창은 그림 4-2 과 같다.

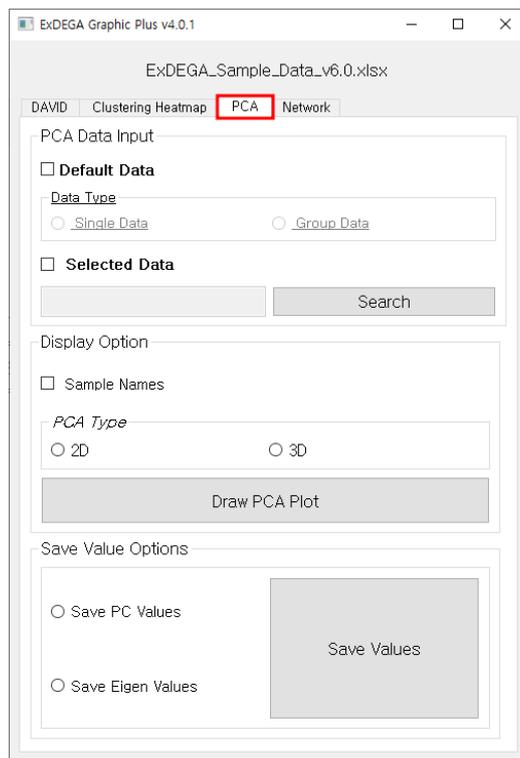


그림 4-2. Select PCA tab

1~3 의 과정을 거쳐 PCA 분석을 수행한다(그림 4-3).

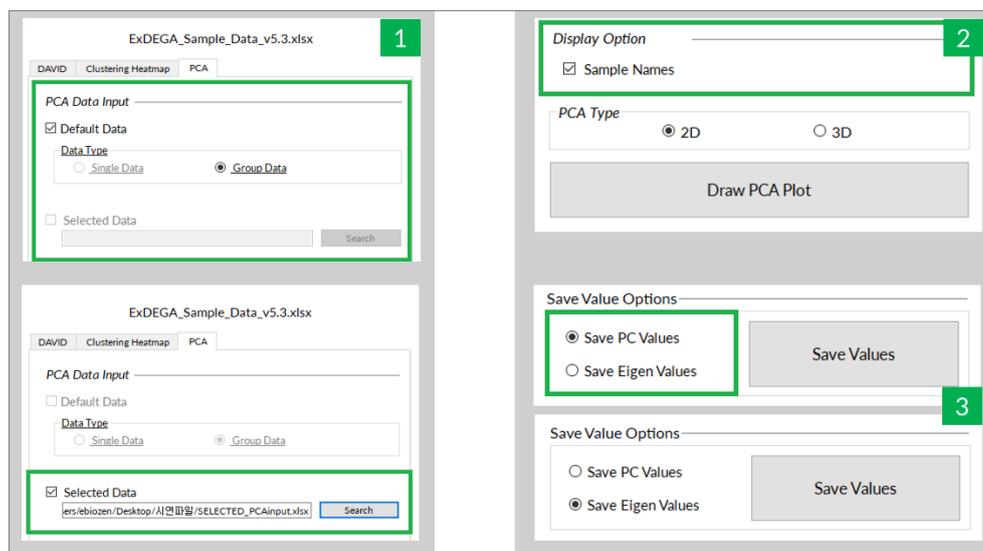


그림 4-3. Steps for Create PCA Graph

먼저, 1 에서 PCA 그래프 제작에 사용될 Data input 타입을 선택해야 한다. 현재 ExDEGA Graphic Plus 에 인식된 ExDEGA 레포트를 그대로 사용하는 경우 Default Data, 별도의 값을 이용하여 PCA 그래프를 제작할 경우 Selected Data 를 선택한다. 이 때, Default Data 를 선택하면 인식된 ExDEGA 레포트의 단일/반복 실험 유형에 따라 자동으로 Single/Group 으로 인식된다. Selected Data 의 경우, 반드시 정해진 양식의 .xlsx 형식의 파일만 입력 가능하니 주의가 필요하다. \* Array Data 의 경우 Default Data 옵션을 이용할 수 없다. Selected Data 옵션을 체크하고, input 파일을 별도로 만들어 이용해야 한다.

Default Data 이면서 Group Data 인 경우, 반복실험 결과( $N \geq 2$ )로 PCA 분석을 하는 경우에 해당한다.

Default Data 이면서 Single Data 인 경우, 반복실험 하지 않은 실험 결과( $N=1$ )로 PCA 분석을 하는 경우에 해당한다.

Selected Data 의 PCA input 파일 만드는 방법은 본 매뉴얼의 '**4-3. PCA Plot input 파일 작성방법**'에 설명되어 있다.

2 에서는 PCA 그래프의 표현 옵션을 제공한다. Display Option 에서 Sample Names 를 체크하면, PCA plot 에 Sample 의 인덱스(샘플 순서대로 부여하는 숫자)를 표시한다. 이 Sample Names 는 그래프 내의 범례로도 표현된다.

마지막으로, 3 에서 2D 를 누르면 PCA 2D (2 차원 평면) 분석 결과가 나오고 3D 를 누르면 PCA 3D (3 차원 공간) 분석 결과가 나온다(그림 4-4). PCA 2D 는 x 축이 PC1, y 축이 PC2 로 작성된 결과이다. PCA 3D 는 x 축이 PC1, y 축이 PC2, z 축이 PC3 로 작성된 결과이다.

각각 좌측은 Display Option 에서 Sample Names 옵션을 선택한 경우, 우측은 선택하지 Sample Names 옵션을 선택하지 않은 경우에 해당한다.

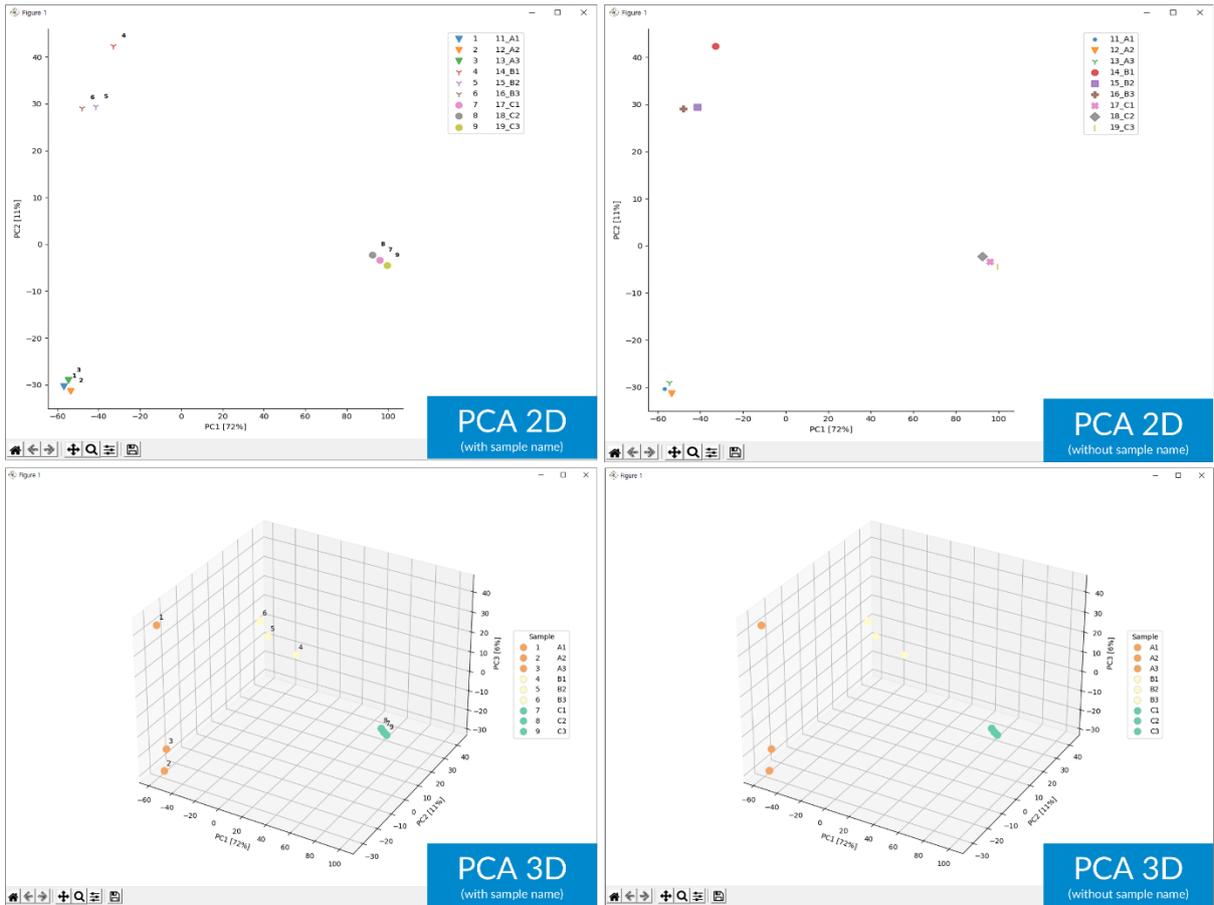


그림 4-4. PCA results 2D, 3D

3 의 Save Value Options 에서 Save PC Values 를 체크한 후 Save Values 버튼을 누르면 각 주성분(PC)이 전체 데이터를 얼마만큼 설명할 수 있는 지 분산 비율로 계산되어 Excel 파일 형식으로 저장된다(그림 4-5). 이 때, 저장하기 전 PC 값들을 그래프로 먼저 확인할 수 있다(그림 4-6).

	A	B
1	Principal Components	Variance Explained
2	PC1	0.716810026
3	PC2	0.10768151
4	PC3	0.056442647
5	PC4	0.047627244
6	PC5	0.038690301
7	PC6	0.024854022
8	PC7	0.007894249
9	PC8	4.87531E-30
10	PC9	3.9363E-31
11		
12		
13		
14		
15		
16		
17		
18		
19		
20		
21		

그림 4-5. Save PC Values as Excel

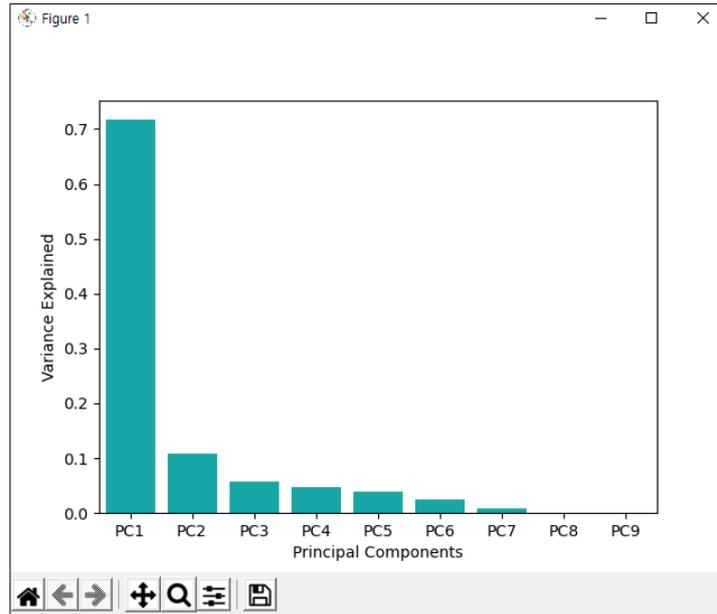


그림 4-6. Check pc values before saving the data as excel

Save Eigen Values 를 체크한 후 Save Values 버튼을 누르면 계산된 Eigenvalue 의 결과를 Excel 파일 형식으로 저장한다(그림 4-7). Eigen values 는 PCA plot 에서 각 샘플의 좌표이다.

	A	B	C	D	E	F	G	H	I	J
1	Sample	PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8	PC9
2	A1	-56.8864	-30.3457	43.08954	-13.297	7.524544	-4.4288	0.033604	1.73E-13	1.17E-14
3	A2	-53.6851	-31.329	-28.1938	-13.4698	-22.2992	-14.1345	0.252672	1.73E-13	1.17E-14
4	A3	-54.6259	-29.076	-18.6736	21.70531	22.02369	16.76101	-0.15929	1.73E-13	1.17E-14
5	B1	-32.9178	42.36758	-12.7574	-19.6184	26.06745	-11.3873	-0.07946	1.73E-13	1.17E-14
6	B2	-41.5794	29.39304	3.757813	-14.5291	-19.7461	25.75895	-0.5247	1.73E-13	1.17E-14
7	B3	-48.145	29.10471	10.89999	36.80833	-12.3007	-12.5653	-0.29896	1.73E-13	1.17E-14
8	C1	95.94653	-3.37155	0.625813	0.8002	-0.42324	-0.00138	0.258714	7.22E-14	1.32E-13
9	C2	92.43661	-2.23767	0.83776	1.080566	-0.51381	0.500988	15.43273	2.24E-13	-4.9E-14
10	C3	99.45645	-4.50543	0.413867	0.519833	-0.33267	-0.50374	-14.9153	2.24E-13	-4.9E-14
11										
12										
13										
14										
15										
16										
17										
18										
19										
20										
21										

그림 4-7. Save Eigen Values as Excel

이 값으로 엑셀에서 분산형 차트를 그려 직접 PCA 2D 를 작성할 수도 있다(그림 4-8).

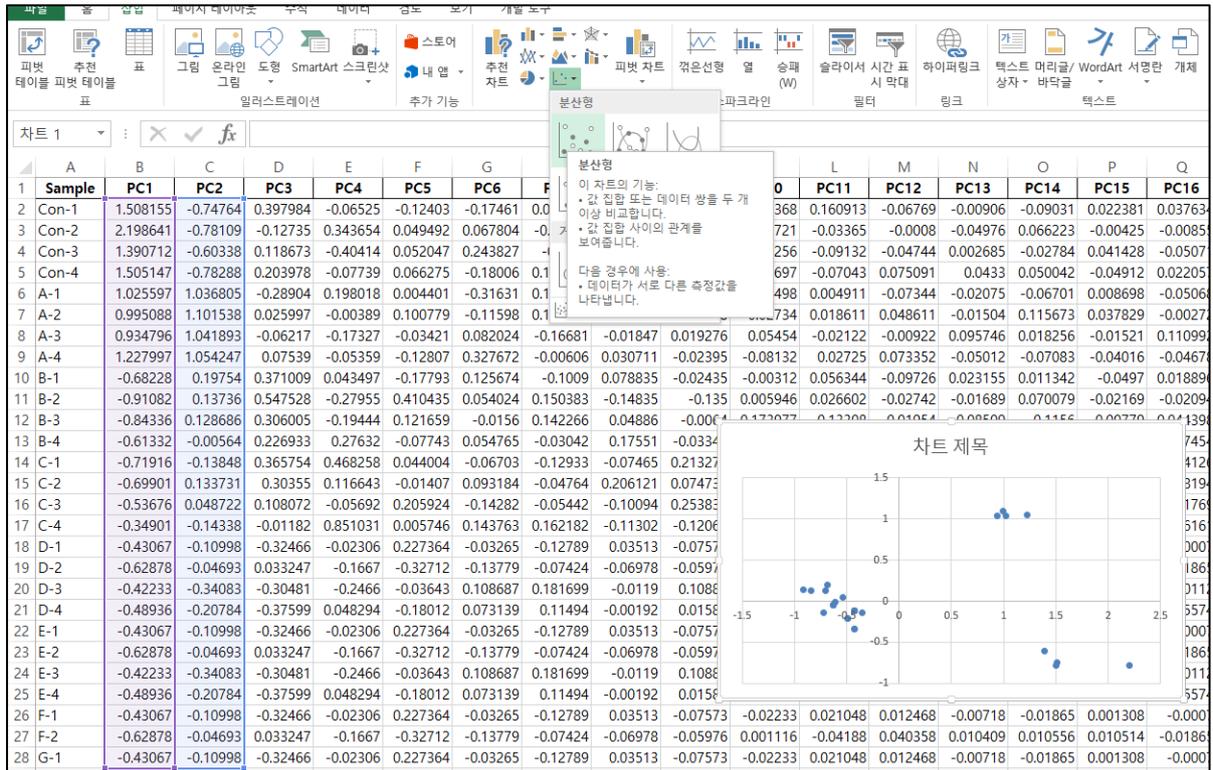


그림 4-8. Create PCA 2D Graph using saved eigen values

### 4-3. PCA Plot input 파일 작성방법

#### Selected Data – Group Data 의 경우

같은 Group 의 Sample 들을 묶어 같은 색상으로 표시하여 PCA 를 표현할 수 있다. Input File 은 Excel 에서 작성한다.

<작성 순서>

1) Normalized data(log2)항목을 복사한다(그림 4-9).

Normalized data (log2)									
	A1	A2	A3	B1	B2	B3	C1	C2	C3
5	1.544	1.369	1.865	1.631	1.149	1.290	0.843	0.743	0.943
6	0.629	0.166	0.221	0.514	0.562	0.159	0.944	0.844	1.044
2	0.000	0.000	0.050	0.000	0.000	0.000	0.100	0.000	0.200
0	0.231	0.330	0.100	0.642	0.091	0.071	0.328	0.228	0.428
5	0.303	0.799	0.111	1.771	0.917	0.294	0.923	0.823	1.023
2	0.000	0.000	0.063	0.055	0.000	0.088	0.100	0.000	0.200
2	0.000	0.290	0.000	0.000	0.098	0.000	0.100	0.000	0.200
2	0.440	0.002	0.000	0.001	0.001	0.862	0.100	0.000	0.200
4	0.218	0.564	0.000	0.001	0.002	0.000	1.742	1.642	1.842
2	0.000	0.000	0.266	0.000	0.000	0.000	0.100	0.000	0.200
2	0.406	0.001	0.000	0.000	0.000	0.000	0.100	0.000	0.200
6	3.136	2.775	2.732	2.813	3.270	3.323	2.474	2.374	2.574
3	1.510	1.210	0.904	1.276	1.172	0.962	2.241	2.141	2.341
2	0.000	0.000	0.000	0.098	0.000	0.157	0.100	0.000	0.200
2	0.000	0.000	0.000	0.000	0.000	0.000	0.100	0.000	0.200

그림 4-9. Copy normalized data

참고) 전체 데이터 선택 방법은 A1(초록색 네모) 클릭 후 Shift 를 누른 상태에서 C3(파란색 네모)를 클릭한다. A1~ C3 까지 선택 후에 Ctrl + Shift + 아래 방향키(↓)를 입력하면 전체 데이터가 한번에 선택된다.

2) 새 엑셀 파일(기존 파일에 sheet 추가가 아닌 엑셀 새로 만들기)을 열어서 **2번째** 열에 붙여넣기 한다(그림 4-10).

	A	B	C	D	E	F	G	H	I
1									
2	A1	A2	A3	B1	B2	B3	C1	C2	C3
3	1.544	1.369	1.865	1.631	1.149	1.290	0.843	0.743	0.943
4	0.629	0.166	0.221	0.514	0.562	0.159	0.944	0.844	1.044
5	0.000	0.000	0.050	0.000	0.000	0.000	0.100	0.000	0.200
6	0.231	0.330	0.100	0.642	0.091	0.071	0.328	0.228	0.428
7	0.303	0.799	0.111	1.771	0.917	0.294	0.923	0.823	1.023
8	0.000	0.000	0.063	0.055	0.000	0.088	0.100	0.000	0.200
9	0.000	0.290	0.000	0.000	0.098	0.000	0.100	0.000	0.200
10	0.440	0.002	0.000	0.001	0.001	0.862	0.100	0.000	0.200
11	0.218	0.564	0.000	0.001	0.002	0.000	1.742	1.642	1.842
12	0.000	0.000	0.266	0.000	0.000	0.000	0.100	0.000	0.200
13	0.406	0.001	0.000	0.000	0.000	0.000	0.100	0.000	0.200
14	3.136	2.775	2.732	2.813	3.270	3.323	2.474	2.374	2.574
15	1.510	1.210	0.904	1.276	1.172	0.962	2.241	2.141	2.341
16	0.000	0.000	0.000	0.098	0.000	0.157	0.100	0.000	0.200
17	0.000	0.000	0.000	0.000	0.000	0.000	0.100	0.000	0.200
18	0.000	0.000	0.000	0.000	0.000	0.000	0.100	0.000	0.200
19	0.000	0.000	0.000	0.000	0.000	0.000	0.100	0.000	0.200

그림 4-10. Paste normalized data in new excel file

3) 1 번째 열에는 그룹명 입력 후 병합한다(그림 4-11).

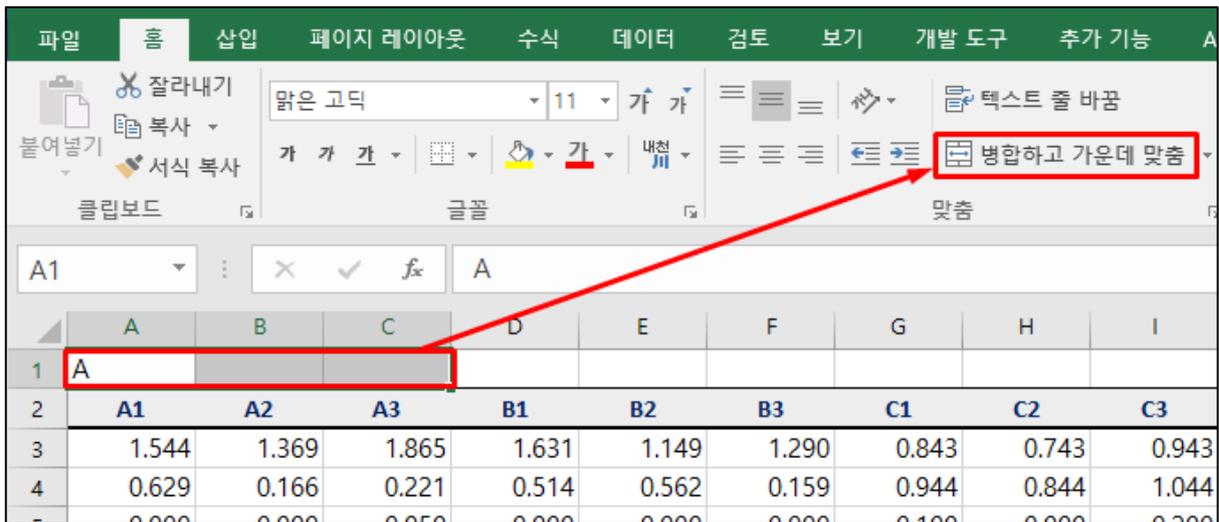


그림 4-11. Group information

4) 완성 후 파일형식은 엑셀로 저장한다(그림 4-12). Input file 명에는 띄어쓰기가 들어가지 않도록 주의한다.

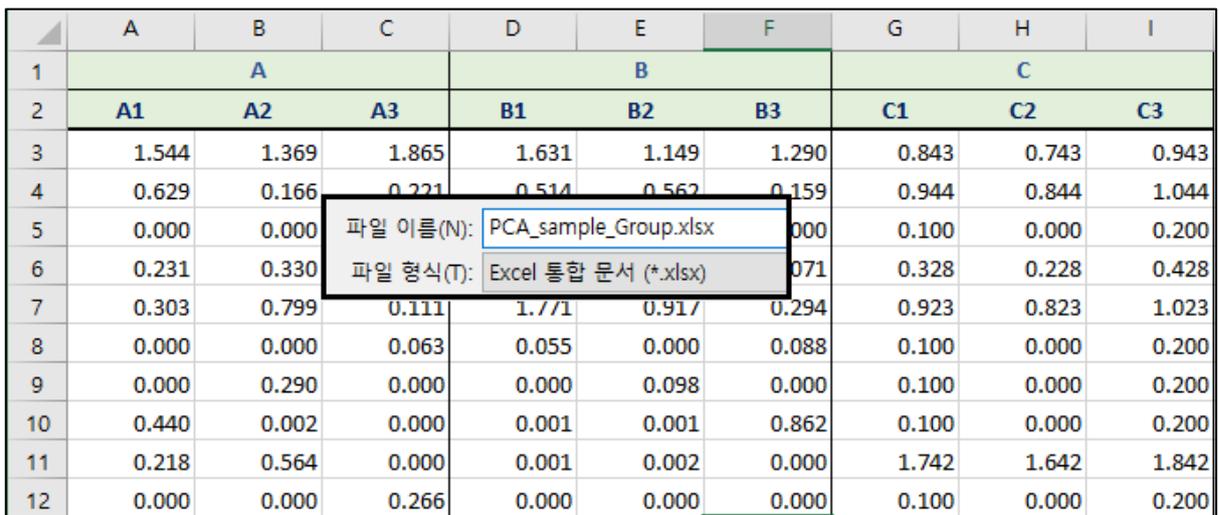


그림 4-12. Save PCA input (group data) file

## 5. String Network Analysis (ExDEGA GraphicPlus)

STRING tool 은 Protein-Protein Interaction 데이터 베이스를 기반으로 유전자의 상관관계를 통계적으로 분석하여 Network 을 작성해주는 분석 툴이다. Graphic plus 를 이용하면 관심있는 유전자 리스트를 넣고 바로 STRING network 이미지를 얻을 수 있다. 이미지 편집이 필요하다면 매뉴얼 뒤쪽의 Cytoscape-STRING 분석 툴을 활용하여 분석이 가능하다.

분석 과정은 그림 5-1 과 같다.

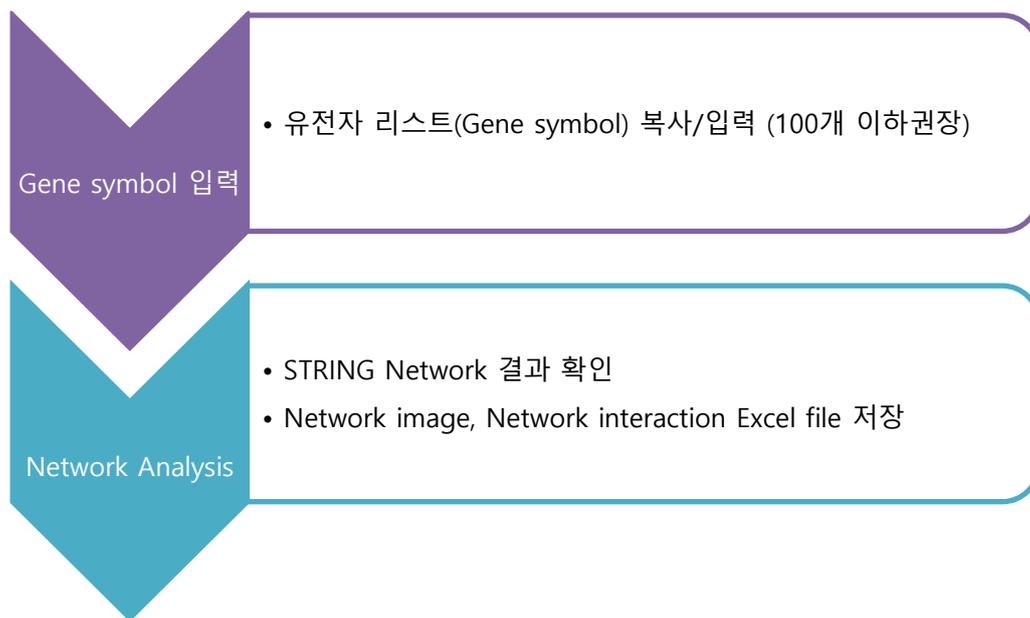


그림 5-1. STRING analysis process

분석을 진행하기 위해서는 먼저 관심있는 유전자를 선별해야 하는데, 이때 유전자 개수는 100개 이내로 선별하는 것을 권장한다. 유전자 개수가 너무 적으면 network가 잘 형성이 되지 않고, 개수가 너무 많으면 복잡도가 높아져서 이미지를 직관적으로 보기 어렵기 때문이다.

유전자를 선별한 뒤 그림 5-2와 같이 빨간색 상자에 유전자 리스트를 복사 붙여넣기 해준다. 만약 유전자 리스트를 엑셀데이터 새창으로 저장했을 경우, 상단의 Import 버튼을 통해서도 불러올 수 있다.

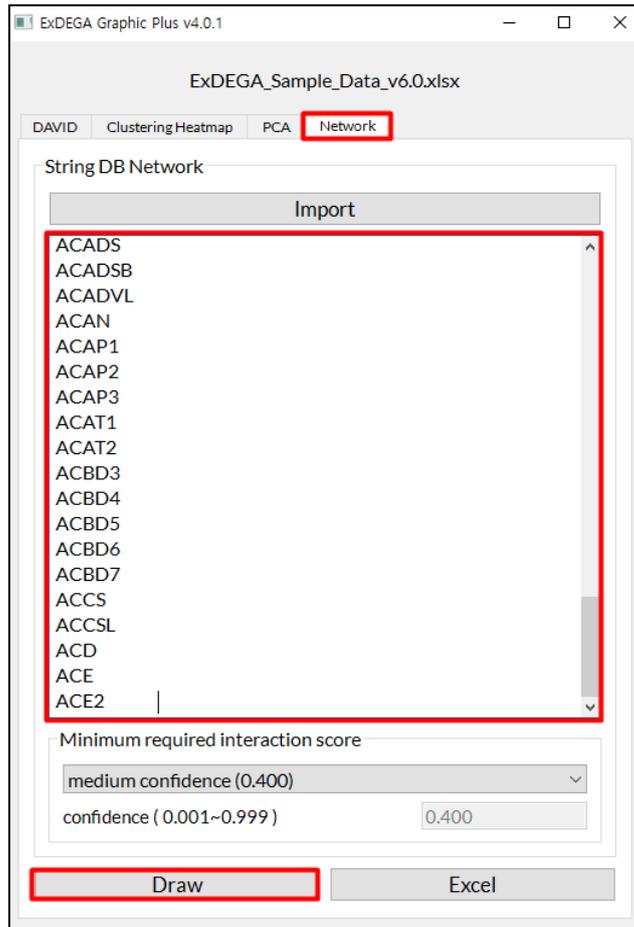


그림 5-2. Graphic plus STRING network 분석 창

그리고 Draw 버튼을 클릭하면 바로 STRING network 분석이 진행되어 이미지를 저장할 수 있는 창이 뜬다. 이미지는 벡터형식인 .svg 형식으로 추출되며, 저장하고자 하는 경로를 설정한 뒤 저장버튼을 클릭하면 이미지 저장이 완료된다. (그림5-3)

오른편의 Excel 버튼을 클릭하면 유전자 간의 interaction score 값을 엑셀형식으로 저장하여 확인할 수 있다. 결과자료는 그림 5-4와 같이 3개의 열로 구성이 되는데, 앞쪽 두개의 열(node1, node2)은 서로 연관 있는 유전자명을 나타내며 세번째 열은 STRING 분석을 통해 계산된 두 유전자 간의 interaction score 값을 나타낸다.

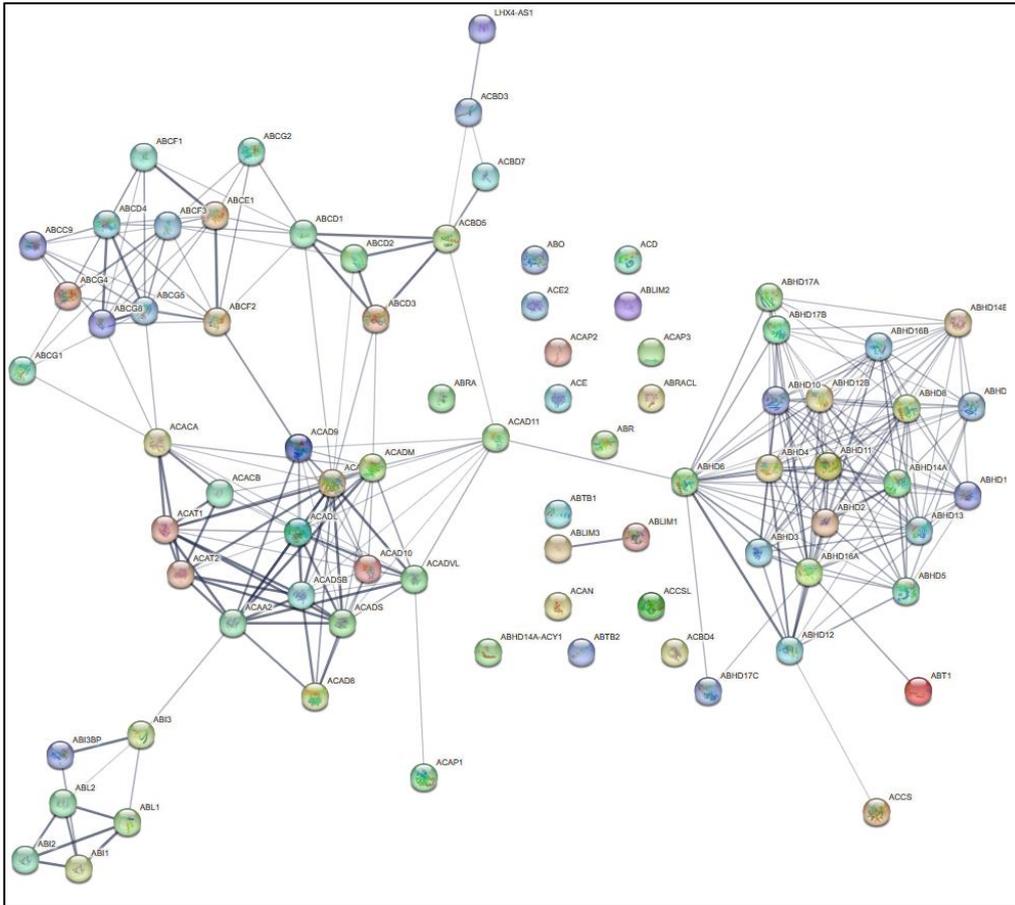


그림 5-3. STRING network 이미지

	A	B	C
1	node1	node2	score
2	ACAP1	ACADV L	0.453
3	ABCD1	ABCF2	0.451
4	ABCD1	ABCF1	0.46
5	ABCD1	ACAA1	0.499
6	ABCD1	ABCE1	0.535
7	ABCD1	ABCF3	0.539
8	ABCD1	ABCG2	0.587
9	ABCD1	ABCD2	0.902
10	ABCD1	ACBD5	0.947
11	ABCD1	ABCD3	0.959
12	ABCF2	ABCC9	0.462
13	ABCF2	ABCF3	0.487
14	ABCF2	ABCG8	0.529
15	ABCF2	ABCG2	0.59
16	ABCF2	ABCD4	0.604
17	ABCF2	ABCG5	0.609
18	ABCF2	ABCG4	0.616
19	ABCF2	ACAD9	0.725
20	ABCF2	ABCE1	0.964

그림 5-4. STRING network interaction score 값

STRING network 옵션으로 그림 5-5와 같이 interaction score 임계점을 설정할 수 있는데, 기본 값은 0.4 로 분석이 진행되며 0~1 사이 값을 설정할 수 있다. Score 값이 1에 가까워질수록 연관도가 높은 유전자 간의 edge(선)이 형성되며, 0에 가까워질수록 연관도가 적은 유전자까지도 edge가 형성된다. 만약 gene list가 너무 많아서 network가 복잡할 경우에는 score 값을 올리고, gene list가 적어서 단순할 경우에는 score값을 내리는 방식으로 이미지를 제작할 수 있다.

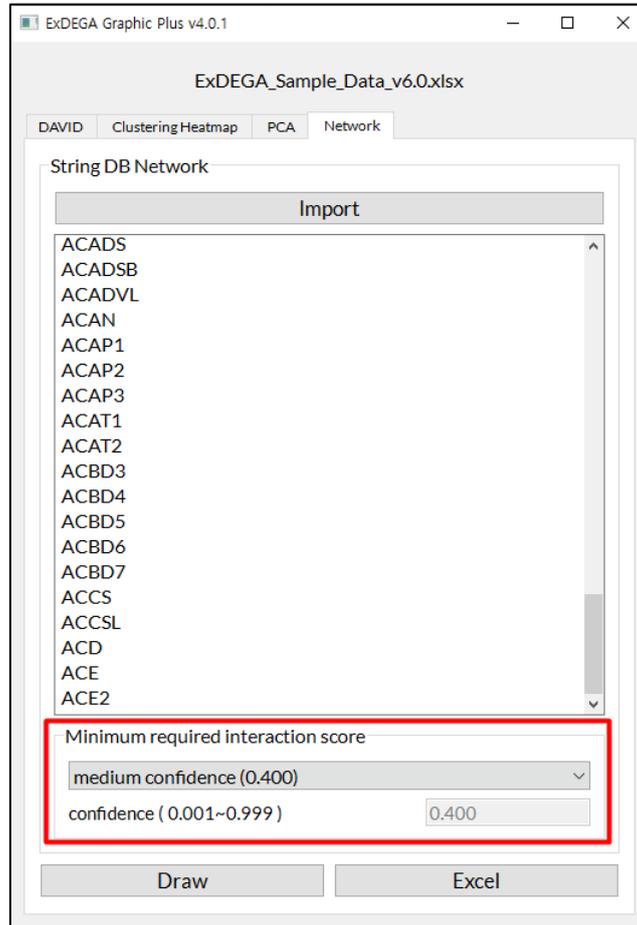


그림 5-5. STRING network score 옵션 설정

## 6. Pathway analysis (KEGG mapper)

RNA-Seq 분석 결과에서 up/down-regulated genes들이 어떤 Pathway에 속하는지 확인하고자 한다면 KEGG에서 제공하는 KEGG Mapper를 이용한다. 사용방법은 그림 6-1과 같은 순서로 진행된다.

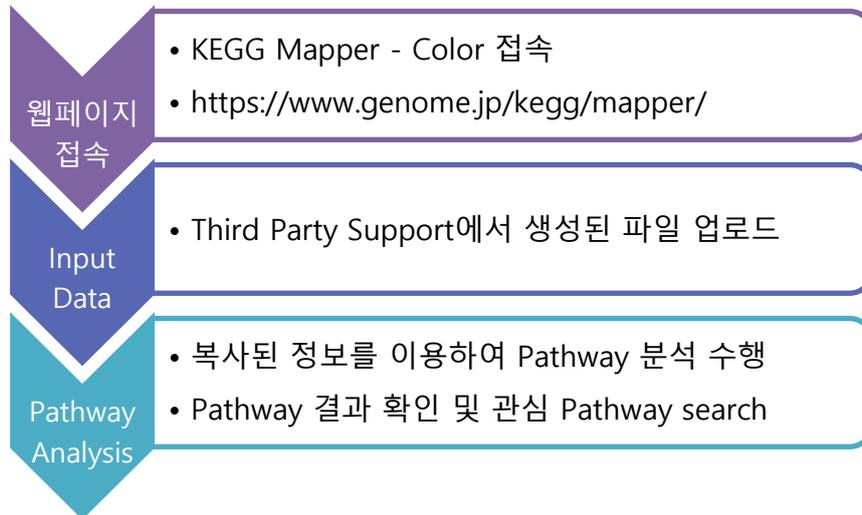


그림 6-1. KEGG Mapper tool analysis process

그림 6-2는 mRNA-Seq report에서 Fold change 2, normalized data(log2) 4, p-value<0.05을 기준으로 선별한 유전자를 KEGG 분석하는 과정이다.

Significant gene selection에서 Fold change, Normalized Data(log2), p-value (반복실험의 경우) 값을 지정하고, 확인하고자 하는 Fold change 조합을 선택하여 필터를 적용한다.

필터를 적용하여 선별된 유전자를 대상으로 Third Party Support를 통해 KEGG input을 추출하여, KEGG 분석에 사용한다. 반드시 하나의 비교 조합만 진행해야 한다.

ID	Gene symbol	Fold change			p-value			Average of r
		B/A	C/A	C/B	B/A	C/A	C/B	
113	111 ABHD3	2.437	2.075	0.853	0.002	0.000	0.176	2.778
370	368 ADGRE3	4.194	4.289	1.018	0.015	0.000	0.927	3.071
611	609 ALDH2	0.345	0.233	0.668	0.002	0.001	0.026	5.497
695	693 AMICA1	2.612	2.663	1.020	0.006	0.005	0.764	6.065
870	868 ANXR2	2.282	2.420	1.060	0.021	0.000	0.718	3.081
934	932 APBB1P	2.711	2.745	1.013	0.005	0.001	0.904	4.910
978	976 APOBR	2.383	3.483	1.453	0.040	0.000	0.075	4.947
1023	1021 ADP9	2.275	2.828	1.243	0.022	0.003	0.169	6.151
1068	1066 ARHGAP25	3.604	2.757	0.765	0.002	0.001	0.068	4.612
1091	1089 ARHGAP9	2.176	2.444	1.123	0.000	0.000	0.094	5.440
1214	1212 ARRD3	2.075	1.949	0.939	0.018	0.020	0.460	3.877
1354	1352 ATHL1	2.111	1.310	0.621	0.040	0.116	0.082	3.591
1389	1387 ATP1B3	0.494	0.438	0.834	0.010	0.002	0.329	4.640
1441	1439 ATP9B2	2.198	1.359	0.618	0.022	0.304	0.001	5.566
1535	1533 B3GNT8	2.471	3.434	1.390	0.036	0.001	0.085	2.894
1596	1594 BASP1	2.908	3.652	1.256	0.005	0.002	0.025	6.344
1638	1634 BCKDHA	0.455	0.409	0.899	0.040	0.031	0.254	4.533
1654	1652 BCL6	2.368	2.622	1.107	0.047	0.000	0.633	4.499
1728	1726 BIRC3	0.476	0.576	1.209	0.002	0.000	0.255	4.087
1753	1751 BLVRA	0.326	0.288	0.933	0.040	0.037	0.349	5.372
1977	1975 C11orf68	2.542	1.659	0.669	0.000	0.002	0.003	3.074
2060	2058 C16orf54	2.830	1.923	0.679	0.018	0.021	0.097	4.127
2350	2348 CSAR2	2.411	2.669	1.107	0.004	0.000	0.318	3.056
2622	2620 CANT1	2.312	2.968	1.110	0.002	0.000	0.278	3.201
2635	2633 CARM2	0.458	0.735	1.638	0.001	0.022	0.001	5.603
2661	2659 CARD6-AS1	2.393	1.645	0.688	0.002	0.016	0.005	2.791
2711	2709 CAST	0.489	0.662	1.353	0.023	0.074	0.019	5.721
2983	2981 CCNH	0.404	0.677	1.677	0.019	0.112	0.001	4.239
2989	2987 CCNL1	0.492	0.620	1.261	0.006	0.006	0.159	4.814
2999	2997 CCPG1	2.050	2.792	1.362	0.008	0.000	0.032	3.587
3033	3031 CD163	0.418	0.157	0.377	0.040	0.012	0.000	4.033
3082	3080 CD4	0.444	0.117	0.272	0.010	0.001	0.002	6.334
3086	3084 CD46	2.013	2.357	1.171	0.008	0.000	0.176	5.042

그림 6-2. KEGG Mapper input file generation process

그림 6-3과 같이 KEGG Mapper 웹페이지(<https://www.genome.jp/kegg/mapper/>)에 접속하고 Color 항목에 들어가면 아래와 같은 화면을 볼 수 있다.

- (1) 분석하고자 하는 유전자의 species를 선택 (Human이면 hsa 선택)
- (2) Or upload file의 파일 선택 버튼을 클릭하고 ExDEGA에서 생성된 입력 데이터 선택
- (3) "Use uncolored diagram"과 "Include aliases" 항목에 체크를 한 후
- (4) Exec 버튼을 누른다.

**KEGG Mapper - Color**

KEGG2 About Reconstruct Search Color Join Convert ID Assign KO Taxonomy

### Color tool

The Color tool searches various KEGG objects, including genes, KOs, EC numbers, metabolites and drugs, against KEGG pathway maps. Found objects may be marked in any combination of background and foreground colors.  
See new article: [KEGG mapping tools for uncovering hidden features in biological data](#)

**Search mode:**  Reference  hsa  other org

**Enter KEGG identifiers followed by color specification**

Examples:  
Select ▼

Or upload file:  Kegg input.txt

Default bgcolor:

Use uncolored diagrams  
 Include aliases (for hsa and other org modes)

그림 6-3. KEGG Mapper tool analysis process

분석결과, 입력한 유전자들이 관여하는 pathway list가 나온다(그림 6-4). 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/blue)이 색으로 표시되어 있다. Pathway 이미지는 “다른 이름으로 저장”이 가능하고 “html”으로 저장하면 이미지에 링크된 항목을 그대로 유지해서 저장이 가능하다.

**\*참고사항**

만약 오른쪽마우스 버튼을 클릭했을 때 다른 이름으로 저장이 보이지 않을 경우, Internet explorer 대신 chrome 창을 이용하면 확인할 수 있다.

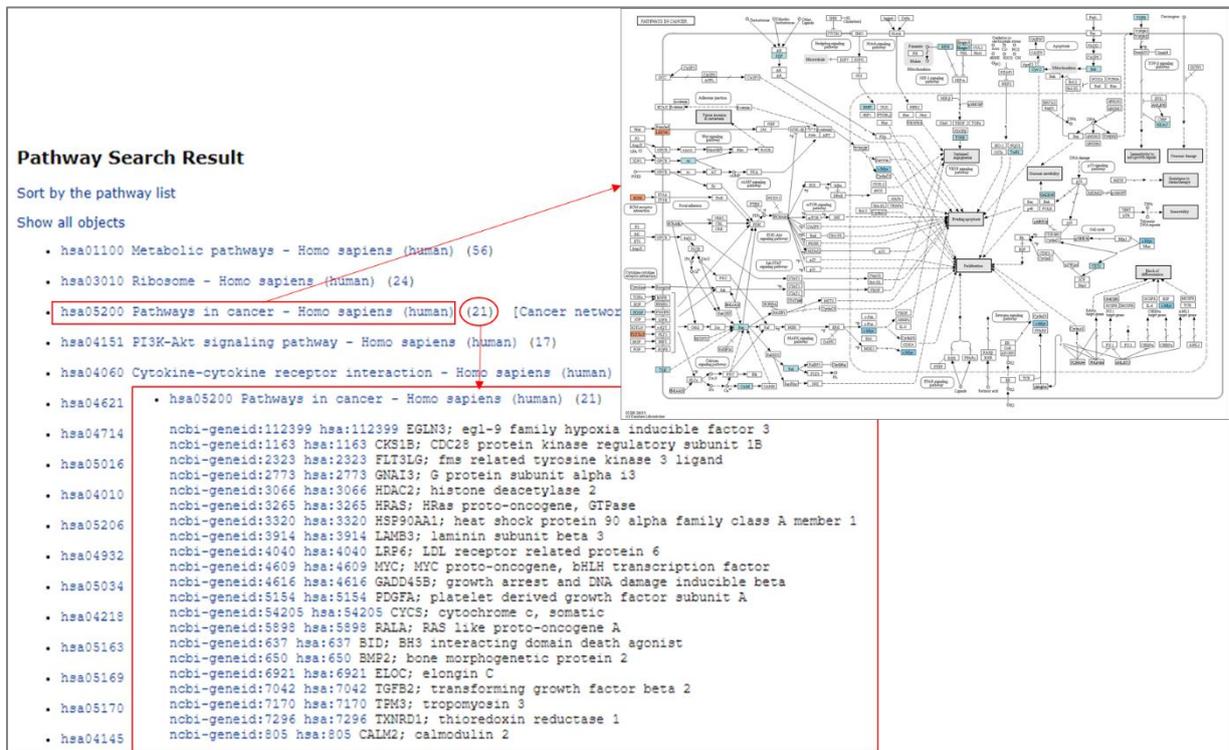


그림 6-4. KEGG Mapper tool analysis result

## 7. Gene set enrichment analysis (GSEA)

Gene set enrichment analysis (GSEA)는 Microarray 또는 RNA-seq data 를 넣어 대조군, 실험군 에서 유의한 gene set 을 분석하는 프로그램이다. GSEA 는 Human, Mouse, Rat 만 지원하고 또한 그룹 비교 (반복 실험) 데이터만 분석이 가능하며, 그룹 당 최소 3 반복 이상이어야 진행이 가능하다.

MSigDB 에 있는 gene set (GO, pathway 등)을 기반으로 분석한다. 분석 과정은 그림 6-1 과 같은 순서로 진행된다.

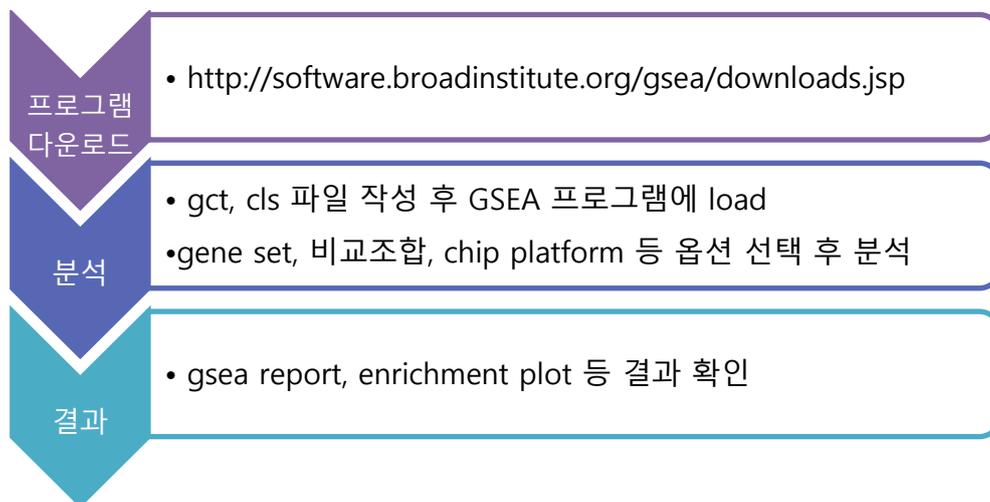


그림 7-1. GSEA tool analysis process

GSEA 홈페이지(<http://software.broadinstitute.org/gsea/downloads.jsp>)에 들어가 회원가입 후 로그인 하여 GSEA 프로그램을 다운로드 받는다 (그림 7-2).

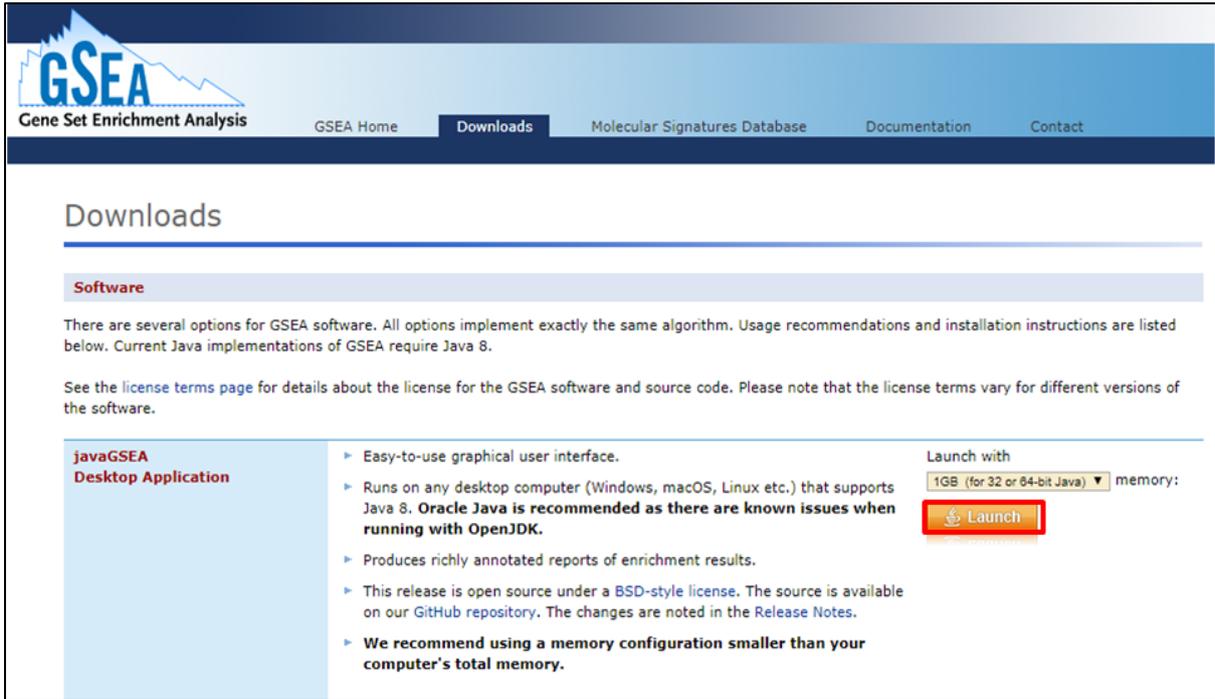


그림 7-2. GSEA program download

GSEA 분석을 위해서는 유전자 발현값 정보가 포함되어 있는 gct 파일과 샘플 정보가 포함되어 있는 cls 파일이 필요하다. Input data는 Third party support에서 추출하여 사용할 수 있다 (그림 7-3).

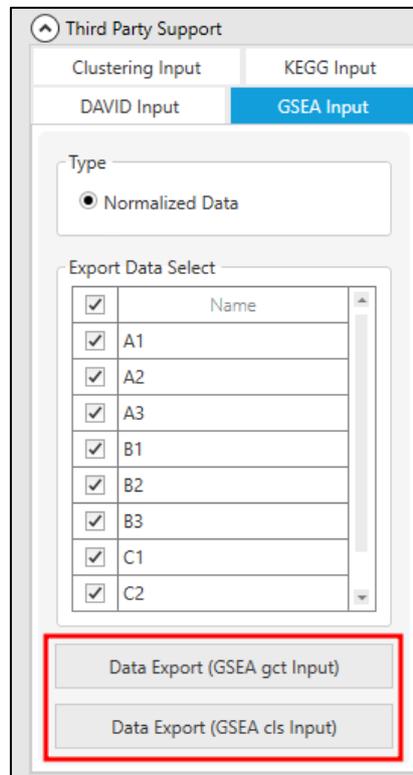


그림 7-3. GSEA Input data

gct 및 cls 파일은 각각 그림 7-4, 7-5와 같은 형식으로 되어있으며, 전체 샘플 및 전체 유전자 리스트에 대해 분석하는 것을 권장한다.

파일 저장할 때는 파일명 뒤에 “.gct” 또는 “.cls”를 반드시 붙이고 파일 형식은 “텍스트 (탭으로 분리) 파일”로 저장한다.

	A	B	C	D	E	F	G	H	I	J	K
1	#1.2										
2	24424	9									
3	Gene sym	Gene Title	A1	A2	A3	B1	B2	B3	C1	C2	C3
4	A1BG	NA	1.544002	1.369196	1.864898	1.630973	1.149183	1.289642	0.842837	0.742837	0.942837
5	A1BG-AS1	NA	0.629423	0.166494	0.220651	0.514093	0.562111	0.158501	0.943667	0.843667	1.043667
6	A1CF	NA	4.39E-05	5.74E-05	0.050181	1.67E-05	2.14E-05	0	0.1	0	0.2
7	A2M	NA	0.230622	0.330027	0.100307	0.641994	0.090967	0.071063	0.328091	0.228091	0.428091
8	A2M-AS1	NA	0.303073	0.798804	0.111377	1.771126	0.91748	0.29448	0.922872	0.822872	1.022872

파일 이름(N): gsea\_input.gct  
 파일 형식(T): 텍스트 (탭으로 분리)

그림 7-4. gct file

	A	B	C
1	9 3 1		
2	#A B C		
3	A A A B B B C C C		
4			

파일 이름(N): gsea\_input.cls  
 파일 형식(T): 텍스트 (탭으로 분리)

그림 7-5. cls file

Microarray data의 경우, gct 파일은 해당 파일을 열어 A열에 probe ID, B열에 Gene symbol로 수정이 필요하다.

GSEA 프로그램을 열어 Load data 버튼을 누르고 Browse for files 버튼을 누른 후 gct, cls 파일을 연다(그림 7-6). gct, cls 파일은 파일의 경로가 길면 input 파일을 잘 인식하지 못하므로 되도록 바탕화면에 두고 진행한다. 문제없이 완료가 되면 NO errors 라는 메시지창을 확인할 수 있다(그림 7-7). 메시지창 확인 후, Run GSEA 버튼을 누른다.

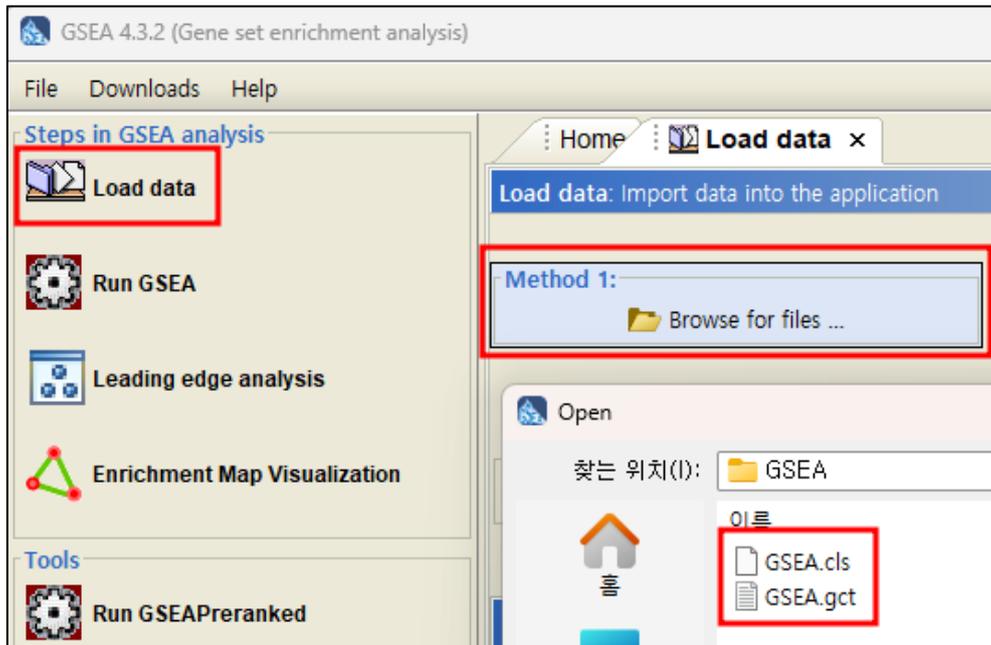


그림 7-6. Load data in GSEA program

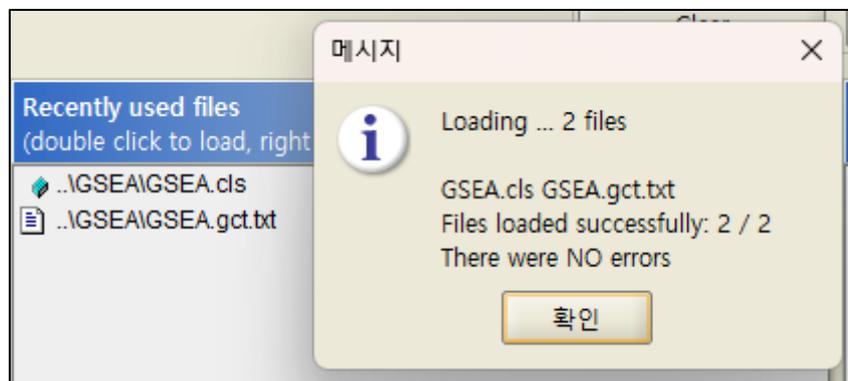


그림 7-7. Load data in GSEA program

Run GSEA 를 누르고 Expression dataset 는 gct 파일명을 선택, gene sets database 는 분석하고자 하는 gene set 을 선택한다(그림 7-8). pathway 분석을 하고자 하면 c2 에서 선택, gene ontology 분석을 하고자 하면 c5 에서 선택한다. Gene set 에 대한 자세한 설명은 GSEA 홈페이지 (<http://software.broadinstitute.org/gsea/msigdb/collections.jsp>)에 있다.

gene sets database 는 현재 Human, Mouse 탭만 있으며, 해당하는 종으로 선택하여 진행하면 되는데, Rat 의 경우, 아직 탭이 존재하지 않아 Human, 또는 Mouse 탭으로 진행한 후, 다음 chip 선택 과정에서 Orthologs 로 진행해야 한다.

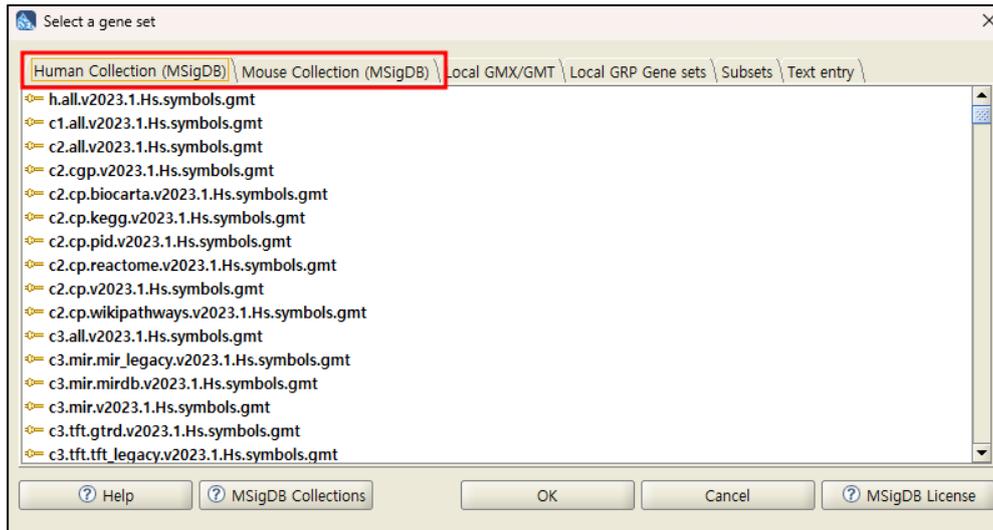


그림 7-8. Select options in GSEA program

Number of permutations은 기본값인 1000으로 기입하고, Phenotype labels은 분석하고자 하는 비교조합을 선택한다. 비교조합 선택 시, 순서는 Test (실험군) versus Control (대조군)인 점을 유의해야 한다.

Collapse/Remap to gene symbols 은 기본값인 Collapse 을 선택하고, permutation type 은 gene\_set 을 선택한다.

Chip platform 은 RNA-seq 의 경우 분석하려는 종의 탭을 선택하여 Human (or Mouse)\_Symbol\_with\_Remapping\_MSigDB~.chip 을 선택한다.

Rat 은 해당하는 탭이 없으므로 Human 또는 Mouse 탭에서 Rat\_Gene\_Symbol\_Remapping\_Human (or Mouse)\_Orthologs\_MSigDB~.chip 을 선택하여 진행한다 (그림 7-9). 여기서 주의할 점은, 이전 gene set 선택 과정 (그림 7-8)에서 선택한 탭과 같은 탭으로 진행해야 한다. Microarray 의 경우에는 실험한 chip 을 선택하여 진행한다.

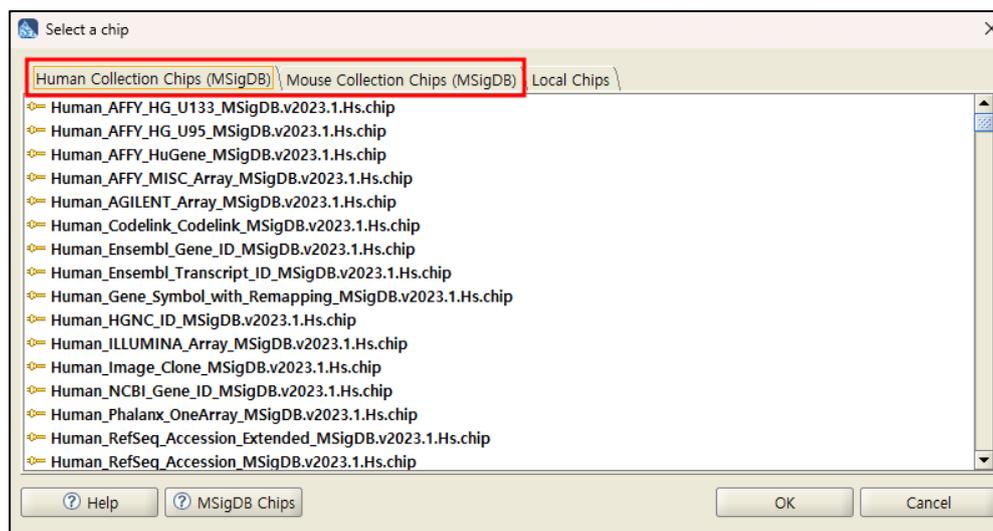


그림 7-9. Select options in GSEA program

모든 옵션이 선택이 완료되고 Run 버튼을 누르면 분석이 시작된다 (그림 7-10).

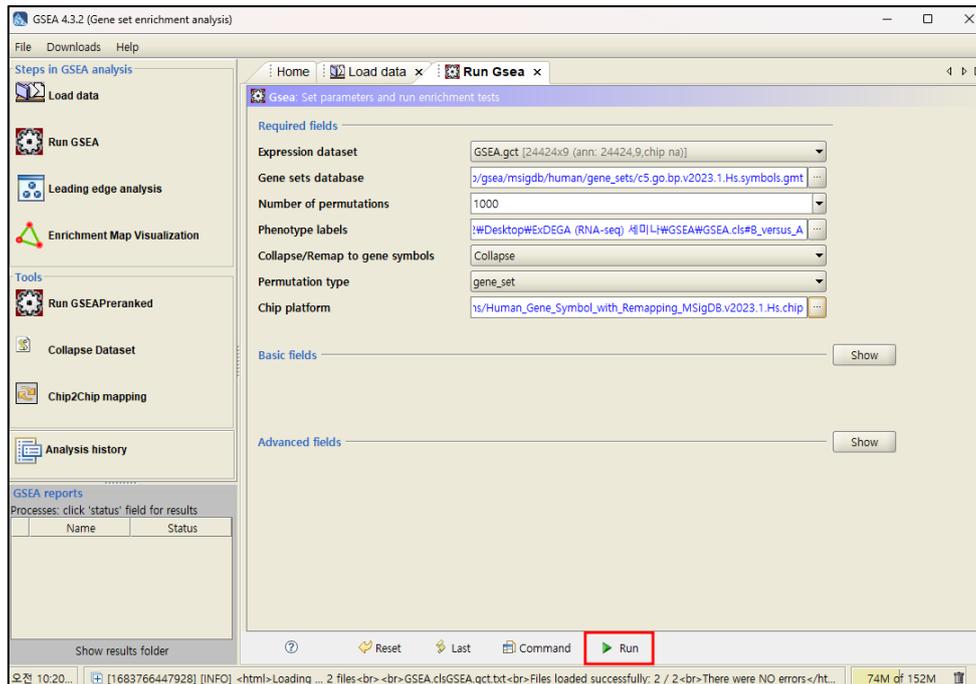


그림 7-10. Run GSEA program

분석이 완료되면 GSEA 왼쪽 아래 GSEA reports 창에 status가 Success로 바뀐다. Show results folder를 누르면 GSEA 분석 결과 창이 열린다(그림 7-11).

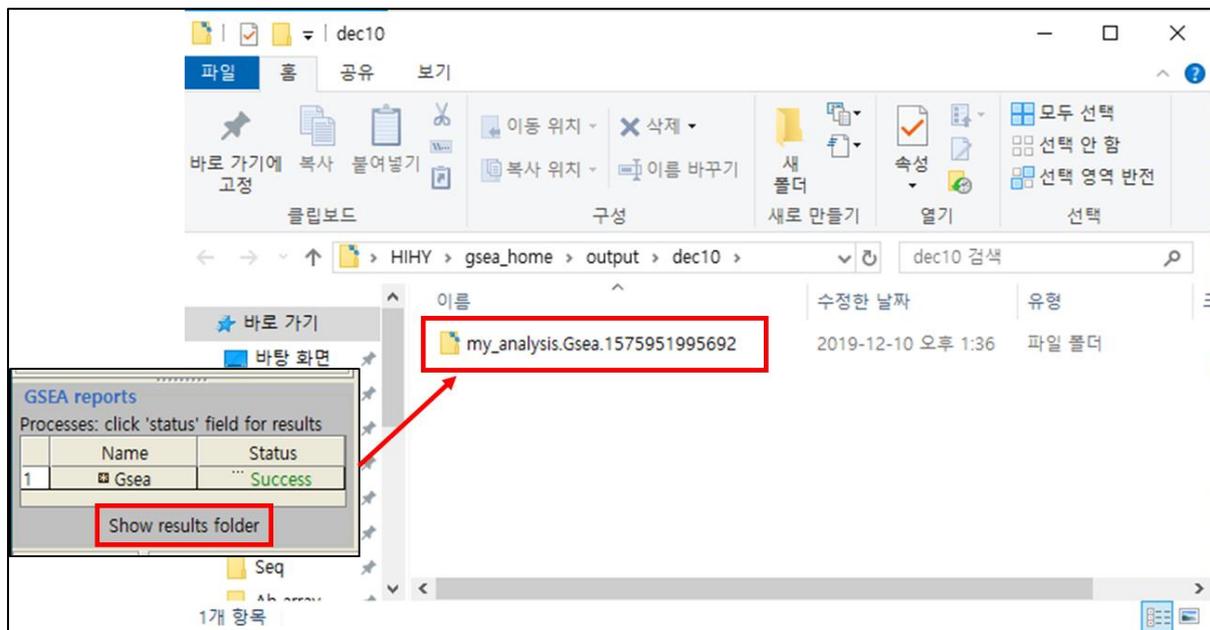


그림 7-11. GSEA results folder

GSEA 결과 중 중요 파일은 'gsea\_report\_for'로 시작하는 엑셀 파일이다. \_for 대조군 파일은 대조군에서 유의한 gene set, \_for 실험군 파일은 실험군에서 유의한 gene set 이다(그림 7-12).

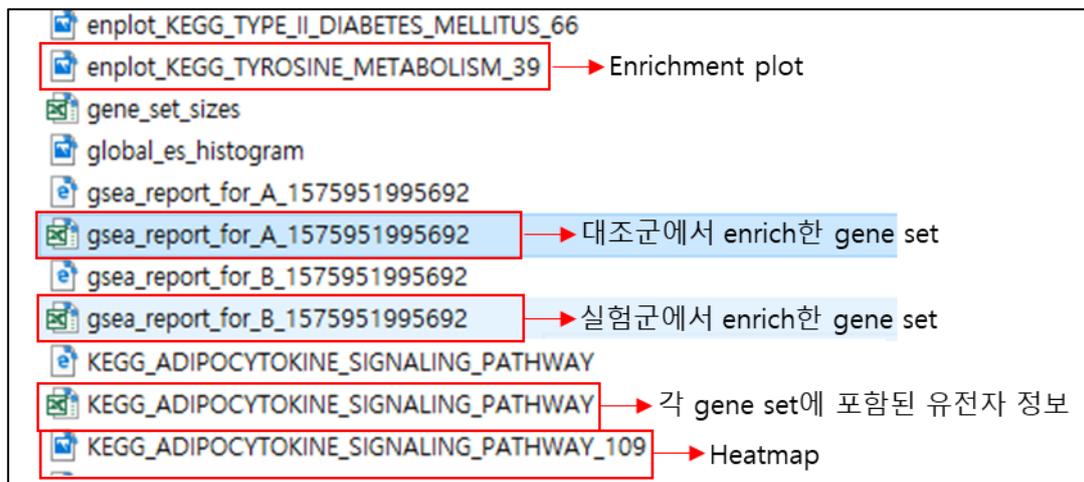
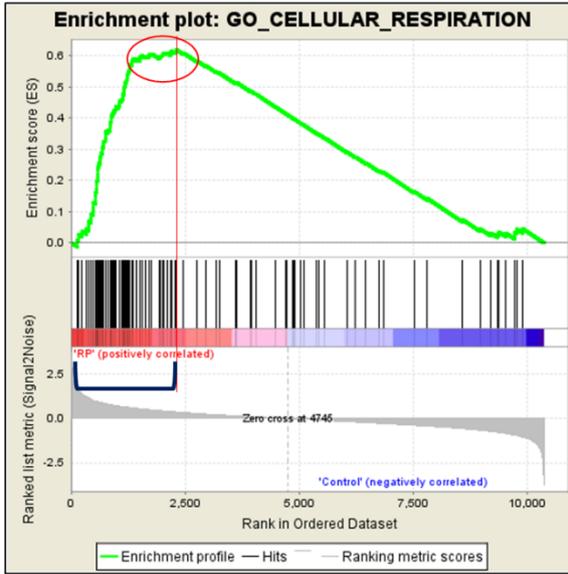


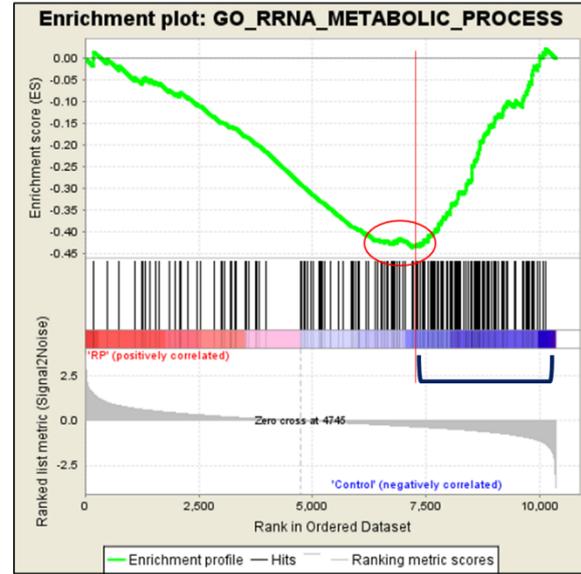
그림 7-12. GSEA result files

\_for 대조군 파일에는 enrichment score (ES)와 Normalized enrichment score (NES)가 음수, \_for 실험군 파일에는 ES와 NES는 양수다. 음수 양수와 관계없이 NES의 절대값이 큰 순서로 ranking 되어 있다. 음수는 DOWN (ranking 하위)에서 core gene의 밀집도가 있다는 것을, 양수는 UP (ranking 상위)에서 core gene의 밀집도가 있다는 것을 의미한다. NES 절대값이 높을수록 유의한 gene set이다. 상위 20개 gene set은 enrichment plot, heatmap, 각 gene set에 포함된 유전자들의 정보가 담긴 excel file이 있다. GSEA 분석 결과 중 Enrichment plot이 논문에 많이 실린다. Enrichment plot 이미지에서 세로 선이 해당 gene set에 포함된 유전자들이며 fold change 순으로 나열된다(그림 7-13). Peak가 왼쪽에 생기면 대조군 대비 실험군에서 up된 유전자들이 많다는 의미이고, peak가 오른쪽에 생기면 down된 유전자가 많다는 의미이다.

▼\_for test 파일내 gene set/ ES, NES 양수



▼\_for control 파일내 gene set/ ES, NES 음수



┌└┐ core enrichment(=core gene) 영역, 관련된 유전자영역이 밀집되어 있는 곳

그림 7-13. GSEA enrichment plot

GSEA 분석과정 및 결과에 대한 카테고리의 자세한 의미는 GSEA user guide

(<https://software.broadinstitute.org/gsea/doc/GSEAUserGuideFrame.html>)에서 확인할 수 있다.

## 8. Protein-Protein Network Analysis (Cytoscape STRING)

STRING tool 은 Protein-Protein Interaction 데이터 베이스를 기반으로 유전자의 상관관계를 통계적으로 분석하여 Network 을 작성해주는 분석 툴이다. 분석 과정은 그림 8-1 과 같다.

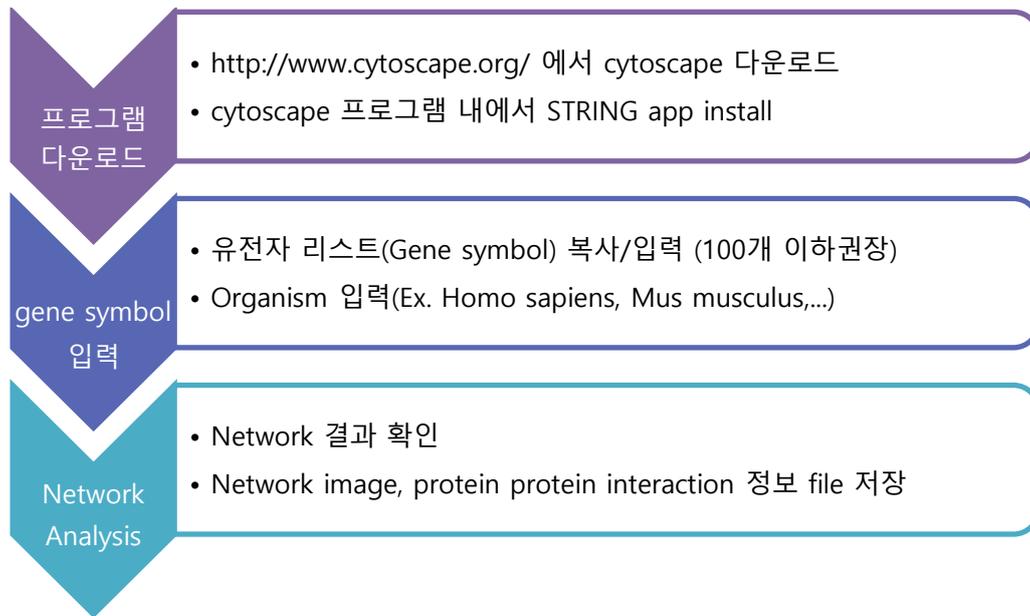


그림 8-1. STRING analysis process

Cytoscape 홈페이지 (<http://www.cytoscape.org/>)에서 cytoscape 프로그램을 다운로드 받아 설치한다(그림 8-2).



그림 8-2. Cytoscape download

Cytoscape 프로그램을 열어 상위에 있는 메뉴 중 [Apps] > [App Manager]로 들어간다(그림 8-3). StringApp 을 선택 후 Install 버튼을 누른다.

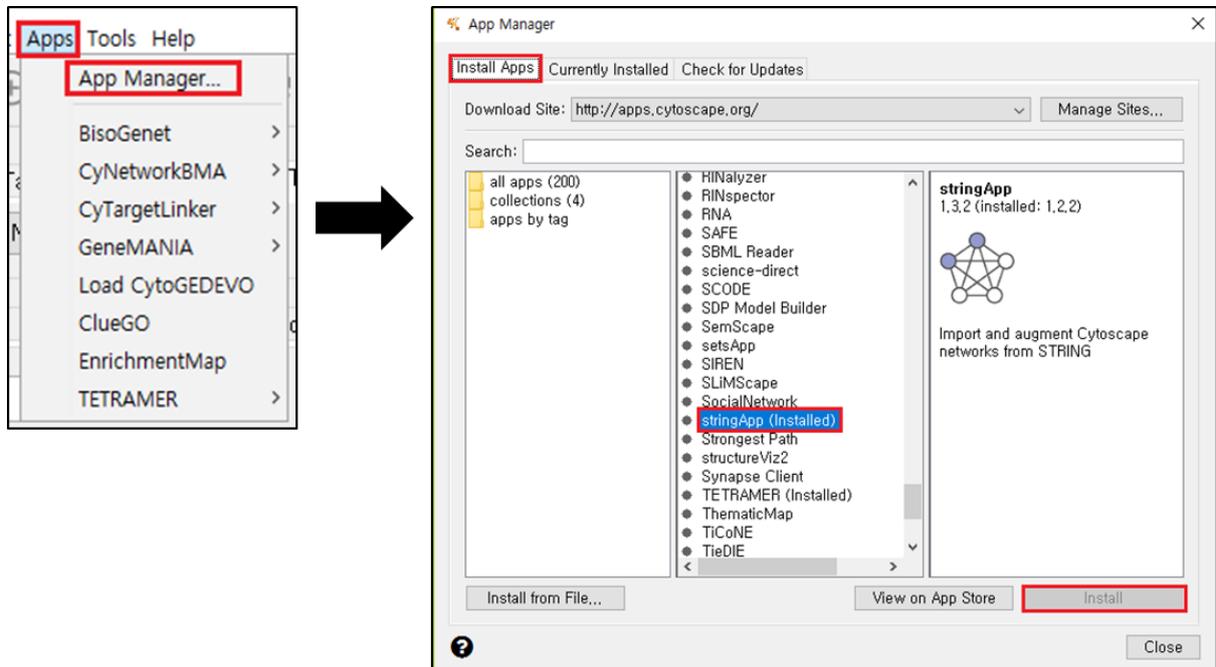


그림 8-3. STRING app installation in Cytoscape

Cytoscape 상위 메뉴 중 [File] > [Import] > [Network from Public Databases]로 들어간다(그림 8-4).

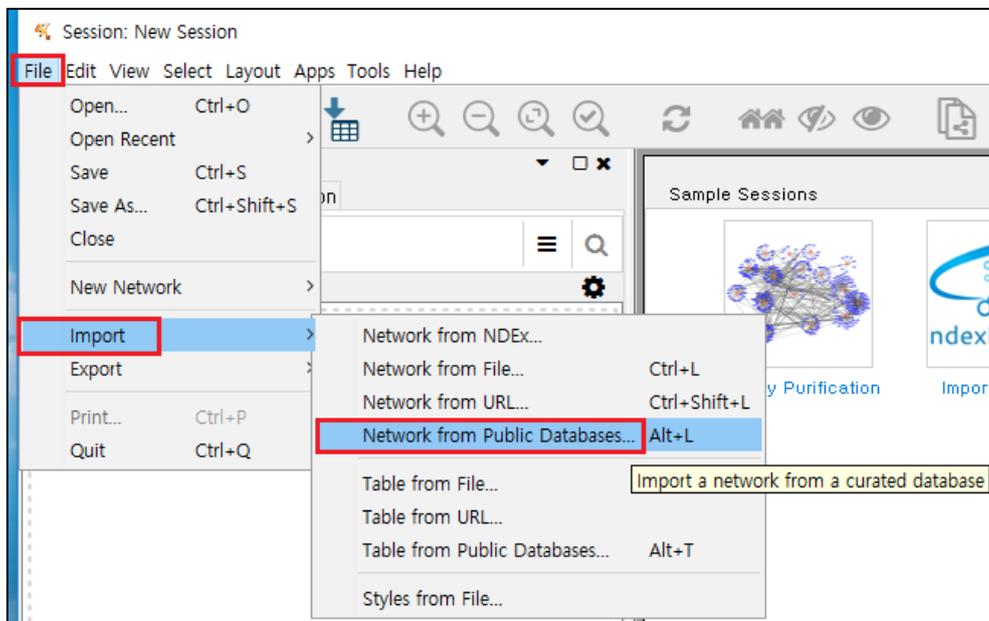


그림 8-4. STRING analysis process 1

Data Source 를 "STRING : protein query" 선택하고 Species 를 선택한다(그림 8-5). 분석하고자 하는 유전자들의 gene symbol 을 입력한다. Confidence (score)는 Protein-Protein Interaction 강도를 뜻하는 것으로 0 부터 1 까지이고, 1 로 갈수록 Interaction 이 강함을 의미한다. Maximum additional interactors 를 0으로 하면 input 한 유전자 안에서만 network 이 그려지고 숫자를 높이면 input 하지 않은 neighborhood protein 까지 network 이 그려진다.

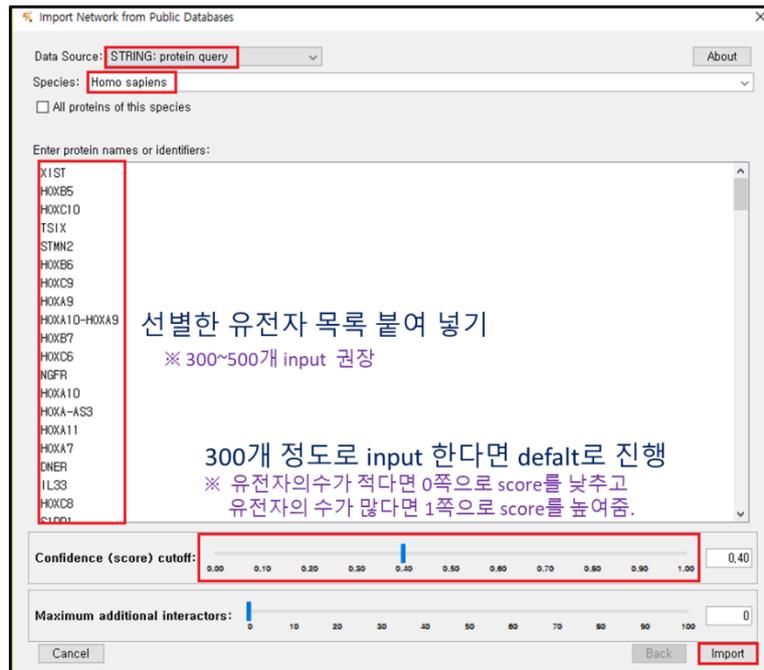


그림 8-5. STRING analysis process 2

Input 한 gene symbol 과 match 가 되지 않는 protein 이 있으면 그림 8-6 과 같은 화면이 나온다. 두 개 이상의 protein 이름이 나타나는 경우는 유사 protein 을 확인하라고 한다. 연구자의 선택에 따라 모두 check 또는 해지한다. Import 버튼을 누르면 분석이 진행된다.

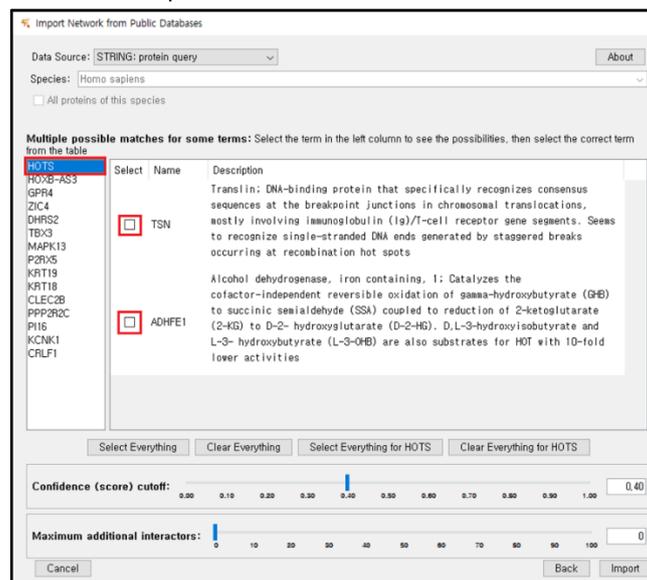


그림 8-6. Not matched proteins in STRING

분석이 완료되면 network image 가 나온다(그림 8-7).

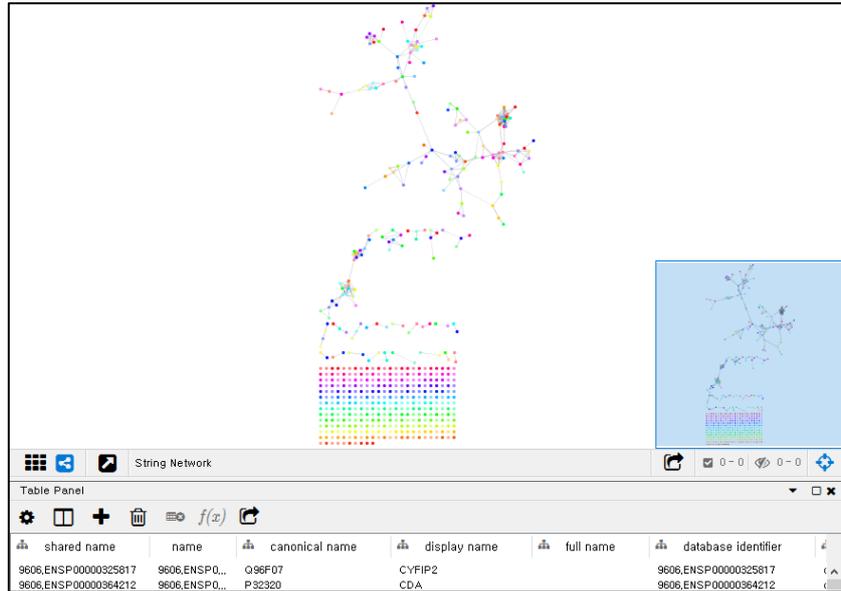


그림 8-7. Network result

[File] > [Export] > [Network to Image]를 눌러 이미지를 저장한다(그림 8-8). PDF 파일형식으로 저장하는 것을 권장한다. Pdf 파일로 저장하면 확대를 하여도 이미지가 깨지지 않는다.

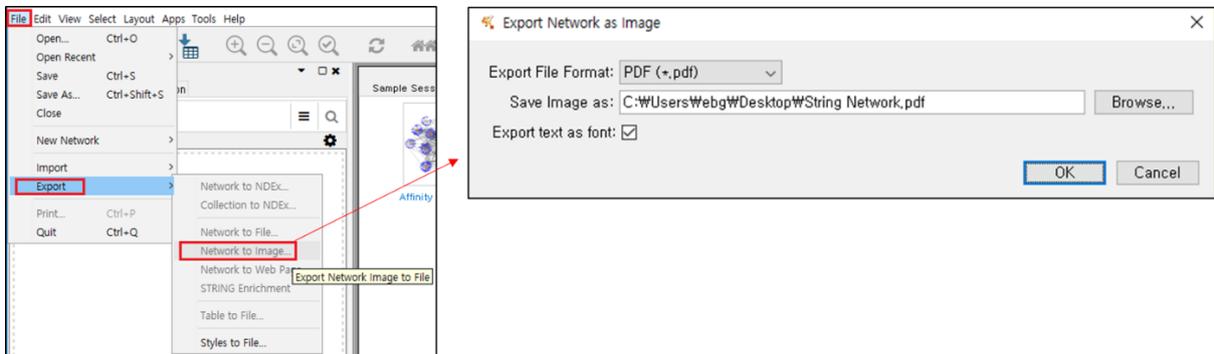


그림 8-8. Save network image

어떤 유전자들이 protein-protein interaction 을 하는지 정보를 저장하고 싶으면 [File] > [Export] > [Table to File...]로 들어가 String Network default edge 파일을 저장한다(그림 8-9).

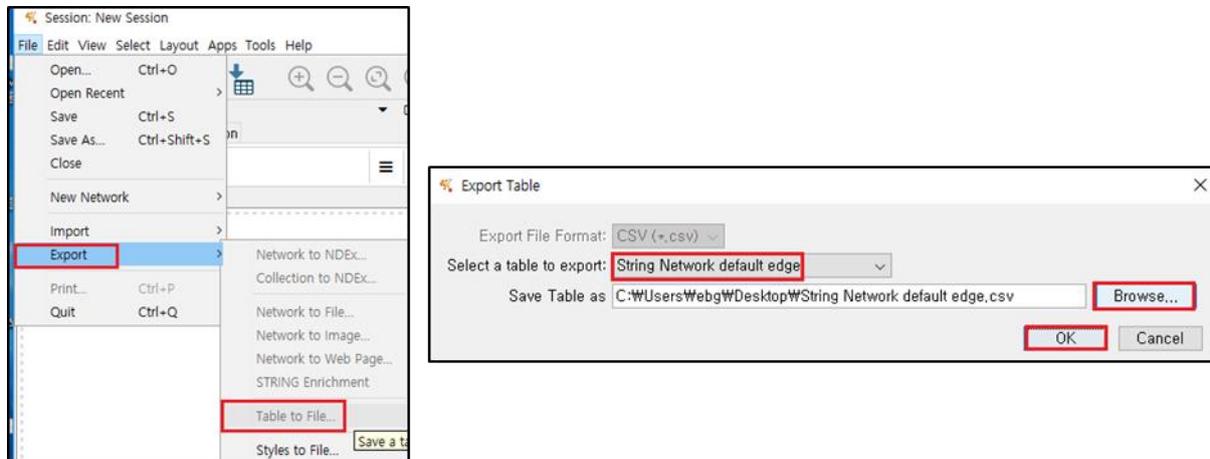


그림 8-9. Save edge table

String Network default edge 파일에서 name 에 interaction 정보, score 에 confidence score 가 나와있다(그림 8-10). Name 에 A (pp) B 라고 적혀있으면 A 유전자와 B 유전자가 Protein-Protein Interaction 한다는 것이고 score 값이 1 에 가까울수록 interaction 이 강한 것이다.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
1	SUID	coexpri	cooccu	databa	experir	fusion	interac	intersp	name	neighb	score	selecte	shared	shared	textmit
2	701	0.397		0.9	0.961		pp		EFNB2 (pp) EPHA4		1.000	FALSE	pp	EFNB2 (pp)	0.926
3	702	0.397		0.9	0.817		pp		EFNB2 (pp) EPHB1		0.999	FALSE	pp	EFNB2 (pp)	0.887
4	961	0.929		0.9			pp		IFIT1 (pp) MX1	0.064	0.998	FALSE	pp	IFIT1 (pp)	0.722
5	700	0.397		0.9	0.76		pp		EFNB2 (pp) EPHA5		0.997	FALSE	pp	EFNB2 (pp)	0.768
6	680	0.878		0.9	0.292		pp		OAS3 (pp) IFIT1		0.996	FALSE	pp	OAS3 (pp)	0.519
7	683	0.888		0.9	0.132		pp		OAS3 (pp) MX1		0.995	FALSE	pp	OAS3 (pp)	0.506
8	935	0.892		0.9			pp		IFI6 (pp) MX1		0.995	FALSE	pp	IFI6 (pp) N	0.503

그림 8-10. Interaction information in edge table

Network image 에서 색이나 모양을 변경하고 싶은 경우에는 STRING Manual ([Download link](#))에서 image 수정 방법을 확인할 수 있다.