



Public Training: Advanced

October 8, 2010



Table of Contents

Logging In to tranSMART Training	1
Lesson 1: Remove Parts of a Search String	3
Lesson 2: Create a Gene Signature	7
Lesson 3: Search for Studies Using a New Gene Signature as a Filter	19
Lesson 4: Use a Heat Map to Compare Treatment Results	25
Lesson 5: Analyze Gene Expression Data from Different Perspectives	29
Lesson 6: Perform a Survival Analysis.....	37
Lesson 7: Perform a Principal Component Analysis	41

Logging In to tranSMART Training

For this training, you will use a training server that is isolated from the “real-world” tranSMART environment. The login credentials and login address that you will use for this training apply to the training server only.

Login Credentials

Login credentials are as follows:

- ID: **publicuserxx**
where xx is a 1- or 2-digit number that the instructor will assign you.
For example, **publicuser5** or **publicuser21**.
- PWD: **training**

Credentials are case-sensitive.

Login Address

Please use the following address to log in to the training server:

<http://75.101.162.195/transmart>

Application and Data Differences

Due to periodic updates to the application and data on the training server, the figures and specific data references in this tutorial may be different than what you see on your screen.

Note: tranSMART includes features and access to data that are not available with the training version of tranSMART.

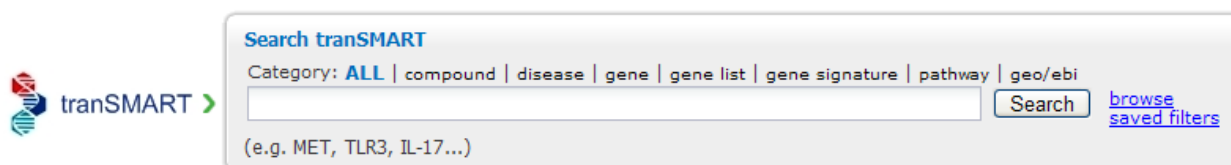
Remove Parts of a Search String

Lesson Goal: Remove parts of a search string, including individual genes in a pathway.

Scenario: You are interested in the relationship between melanoma and the gene EDNRB, and secondarily, in relationships between melanoma and the other genes in the melanogenesis pathway.

1. Start up tranSMART as described in [Logging In to tranSMART Training](#) on page 1.

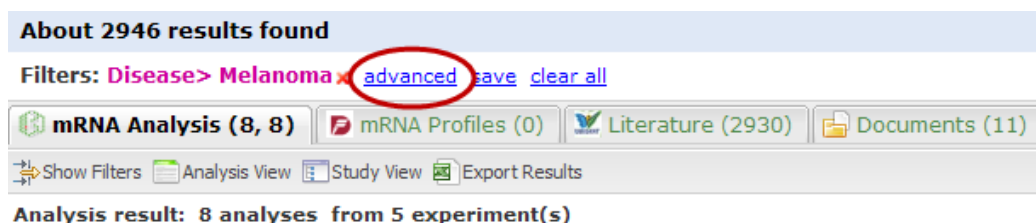
The tranSMART Search window appears:



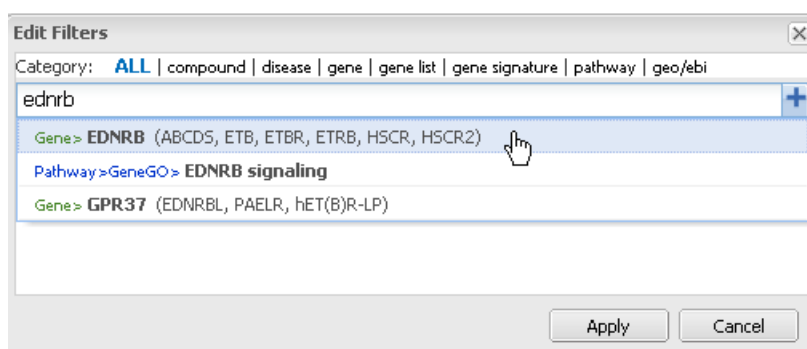
2. Type **melan** in the search field.
3. Click **Disease> Melanoma** in the dropdown list of search filters.

After the search result is returned, you want to add the gene EDNRB to the search string.

4. Click **advanced**:

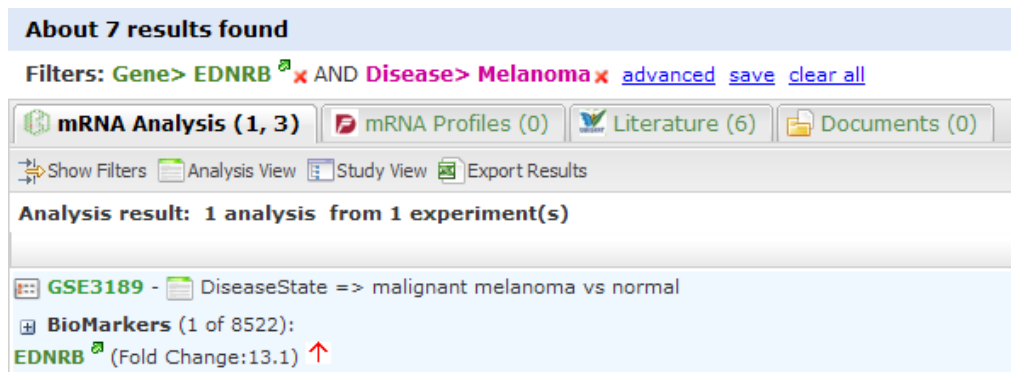


5. Type **EDNRB** in the search field of the Edit Filters dialog.
6. Click **Gene> EDNRB** in the dropdown list:



7. Click **Apply**.

As shown in the figure below, the search returns a single study, **GSE3189**, that matches the search criteria. The study has three matching analyses, but only one, **malignant melanoma vs normal**, that is considered statistically significant.



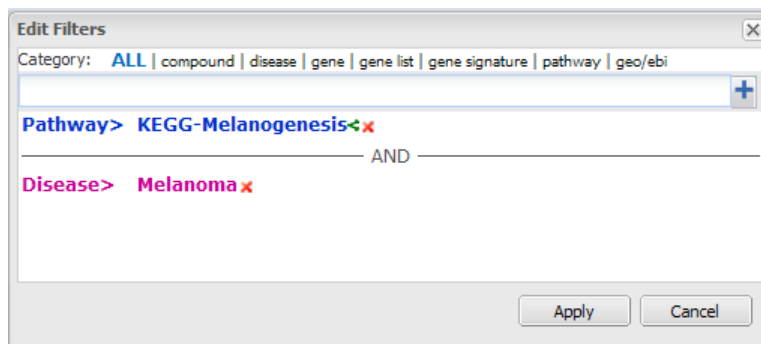
You now want to look for relationships between melanoma and the genes in the melanogenesis pathway. First, you want to remove EDNRB from the search string.

8. Click the red **X** after the name **EDNRB**:

A search begins immediately based on the modified search filter, and a result is returned.

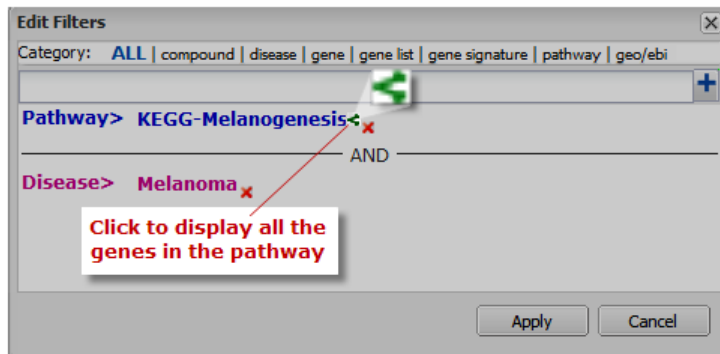
9. Click **advanced**.10. Type **melan** in the search field of the Edit Filters dialog.11. Click **Pathway>KEGG>Melanogenesis** in the dropdown list.

The Edit Filters dialog now looks as follows:



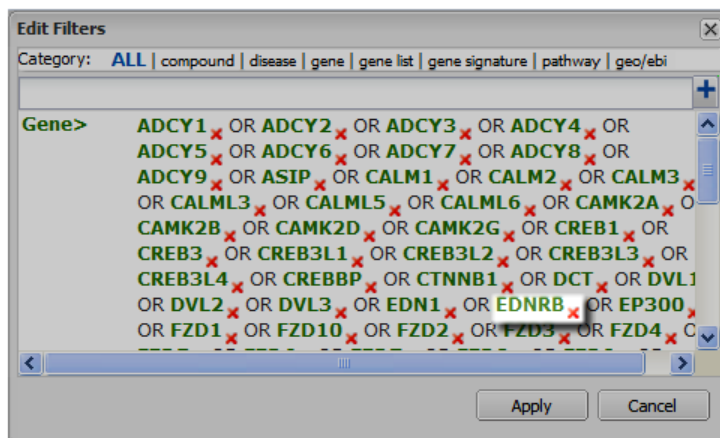
Since you've already searched for relationships between melanoma and EDNRB, you would like to conduct this new search without EDNRB as a factor. To do so, you must remove EDNRB from the melanogenesis pathway.

12. Click the green angle bracket immediately after the name of the melanogenesis pathway:



The name of the melanogenesis pathway is replaced by a list of the individual genes in the pathway.

13. Locate **EDNRB** in the alphabetized list, then click the red **X** after the gene name:



14. Click **Apply**.

tranSMART returns the new search result:

About 67 results found

Filters: Genes> ADCY1 OR ADCY2 OR ADCY3 OR ADCY4 OR ADCY5 OR ADCY6 OR ADCY7 OR ADCY8 OR ADCY9 OR ASIP OR CALM1 OR CALM2 OR CALM3 OR CALML3 OR CALML5 OR CALML6 OR CAMK2A OR CAMK2B OR CAMK2D OR CAMK2G OR CREB1 OR CREB3 OR CREB3L1 OR CREB3L2 OR CREB3L3 OR CREB3L4 OR CREBBP OR CTNNB1 OR DCT OR DVL1 OR DVL2 OR DVL3 OR EDN1 OR EP300 OR FZD1 OR FZD10 OR FZD2 OR FZD3 OR FZD4 OR FZD5 OR FZD6 OR FZD7 OR FZD8 OR FZD9 OR GNAI1 OR GNAI2 OR GNAI3 OR GNAO1 OR GNAQ OR GNAS OR GSK3B OR HRAS OR KIT OR KITLG OR KRAS OR LEF1 OR MAP2K1 OR MAP2K2 OR MAPK1 OR MAPK3 OR MC1R OR MITF OR NRAS OR PLCB1 OR PLCB2 OR PLCB3 OR PLCB4 OR POMC OR PRKACA OR PRKACB OR PRKACG OR PRKCA OR PRKCB OR PRKCG OR PRKX OR PRKY OR RAF1 OR TCF7 OR TCF7L1 OR TCF7L2 OR TYR OR TYRP1 OR WNT1 OR WNT10A OR WNT10B OR WNT11 OR WNT16 OR WNT2 OR WNT2B OR WNT3 OR WNT3A OR WNT4 OR WNT5A OR WNT5B OR WNT6 OR WNT7A OR WNT7B OR WNT8A OR WNT8B OR WNT9A OR WNT9B AND Disease> Melanoma [advanced](#) [save](#) [clear all](#)

mRNA Analysis (5, 6) mRNA Profiles (0) Literature (56) Documents (7)

Show Filters Analysis View Study View Export Results

Analysis result: 5 analyses from 3 experiment(s) [4 Significant TEA / 1 Insignificant TEA]
 Note, only significant TEA Analyses are displayed!

Create a Gene Signature

Lesson Goal: Use the tranSMART gene signature wizard to create a gene signature.

Scenario: You are interested in lung adenocarcinoma, and want to create a gene signature consisting of genes that were strongly up-regulated in an experiment involving lung adenocarcinoma patients.

This lesson involves two basic tasks:

1. Find the genes to include in the gene signature, then write the genes to a tab-delimited text file that can be imported into the gene signature.
2. Define the gene signature and import the file containing the genes.

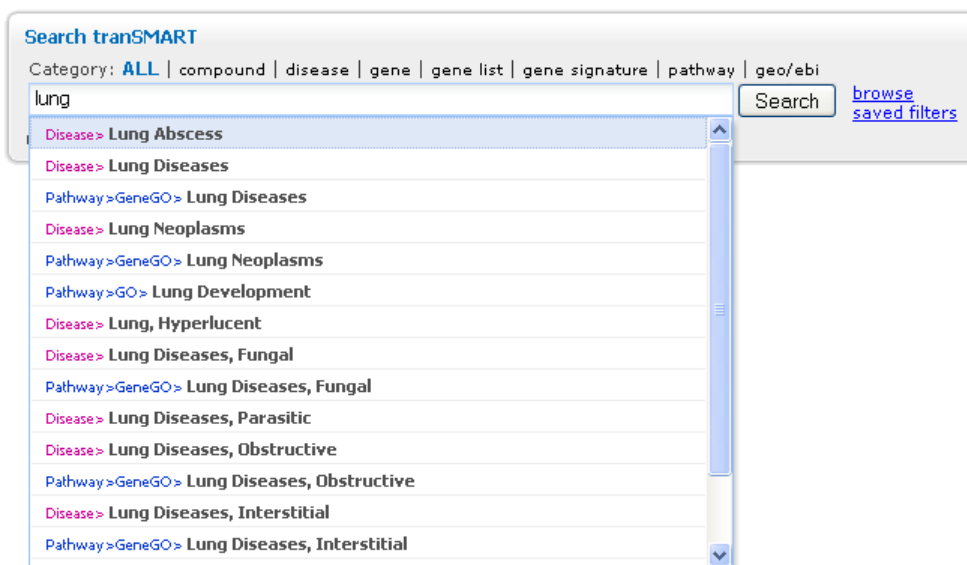
Task 1: Find the Genes for the Gene Signature and Write them to a File

1. Click the tranSMART logo to clear the results from the previous lesson:

Note: Your training application may have the following logo in the place of the tranSMART logo shown above:



2. Type **lung** in the Search field:

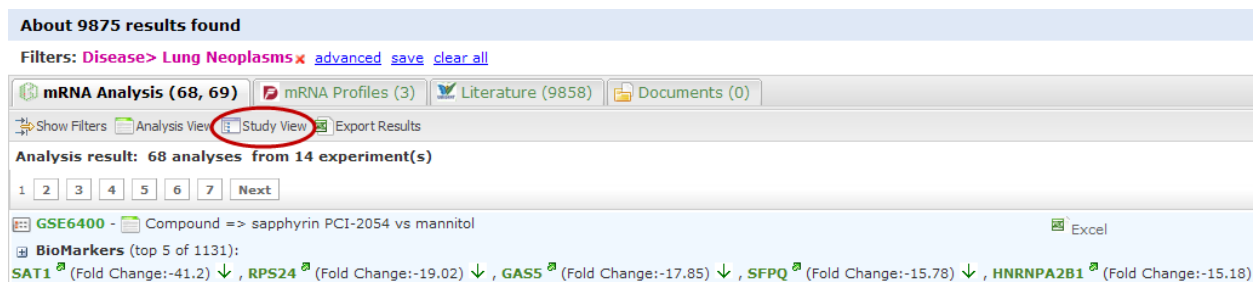


Lung adenocarcinoma is not listed in the dropdown list of search filters, but lung neoplasms is listed.

3. Click **Disease > Lung Neoplasms**.

In a few seconds, the search result appears.

4. Click the **Study View** button in the **mRNA Analysis** tab:

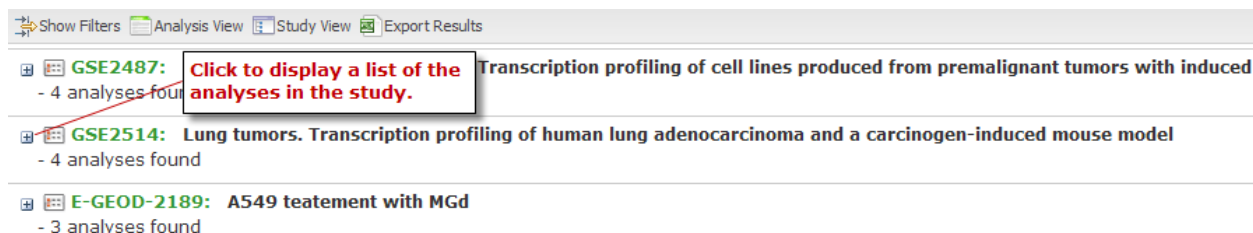


tranSMART displays a list of all the experiments related to lung neoplasms.

5. Scroll through the list of experiments until the experiment **GSE2514** appears.

You notice that this experiment focused specifically on lung adenocarcinoma.

6. Click the **+** icon (⊕) to the left of the experiment name:



A list of the analyses based on this experiment appears. The analysis **Lung adenocarcinoma vs normal** interests you.

7. Click the **Excel** button for the analysis **Lung adenocarcinoma vs normal** to export the analysis data to a Microsoft Excel file:

The screenshot shows the GSE2514 analysis interface with the following sections:

- GSE2514: Lung tumors. Transcription profiling of human lung adenocarcinoma and a carcinogen-induced mouse model** - 4 analyses found
- DiseaseStaging => late (42 weeks) vs normal** - Excel button
- BioMarkers (top 5 of 5556):**
 - Ereg (Fold Change:41.05) ↑, Meg3 (Fold Change:32.56) ↑, Dlk1 (Fold Change:25.92) ↑, Myl7 (Fold Change:-20.49) ↓, Klf2 (Fold Change:-18.91) ↓
- DiseaseStaging => early (24 to 26 weeks) vs normal** - Excel button
- BioMarkers (top 5 of 4928):**
 - Myl7 (Fold Change:-120.15) ↓, Meg3 (Fold Change:71.03) ↑, Ereg (Fold Change:67.54) ↑, Actc1 (Fold Change:-62.02) ↓, Tnnc1 (Fold Change:-44.16) ↓
- DiseaseState => lung adenocarcinoma vs normal** - **Excel** button (circled in red)
- BioMarkers (top 5 of 4163):**
 - COL11A1 (Fold Change:28.46) ↑, EEF1A2 (Fold Change:20.16) ↑, UBE2C (Fold Change:18.24) ↑, SLC6A4 (Fold Change:-18.08) ↓, SPP1 (Fold Change:17.75) ↑
- DiseaseStaging => late (42 weeks) vs early (24 to 26 weeks)** - Excel button

8. When the File Download dialog appears, click **Open**.

Excel starts up and displays the analysis data. The following figure shows the first few rows of data:

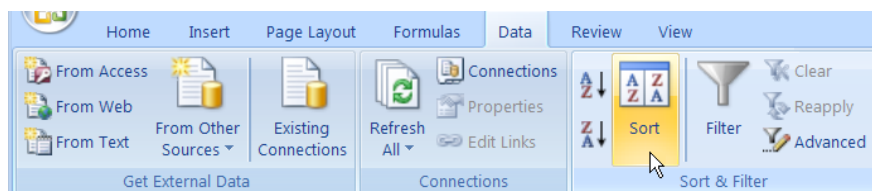
	A	B	C	D	E	F	G
1	Analysis	ProbeSet	Fold Change Ratio	p-Value	adjusted p-value	TEA p-Value	Gene
2	DiseaseState => lung adenocarcinoma vs norma	479_at	-1.45	0.0009	0.0051	0.2136	DAB2
3	DiseaseState => lung adenocarcinoma vs norma	36275_at	-2.84	0	0.00004	0.05363	SEMA6A
4	DiseaseState => lung adenocarcinoma vs norma	34830_at	1.18	0.0174	0.0558	0.26259	ECOP
5	DiseaseState => lung adenocarcinoma vs norma	1970_s_at	-5.31	0	0.00003	0.0011	FGFR2
6	DiseaseState => lung adenocarcinoma vs norma	33908_at	1.17	0.0313	0.0887	0.26451	CAPN1
7	DiseaseState => lung adenocarcinoma vs norma	32146_s_at	-2	0	0	0.13197	ADD1
8	DiseaseState => lung adenocarcinoma vs norma	33062_at	-1.55	0.0359	0.0987	0.1969	GSTA1

You want to select the most strongly up-regulated and down-regulated genes for your gene signature. To find them, you will sort the rows according to the **Fold Change Ratio** column.

Note: The instructions in this lesson for performing Excel operations are based on Microsoft Excel 2007. If you have a different version, some of the steps and graphics may be different for you.

9. Select all the data in all the columns, as follows:
 - a. Click the letter **A** above the leftmost column.
 - b. Press and hold down the **Shift** key, then click the letter above the rightmost column (letter **G** in the figure above).

10. Click the **Data** menu, then click **Sort**:

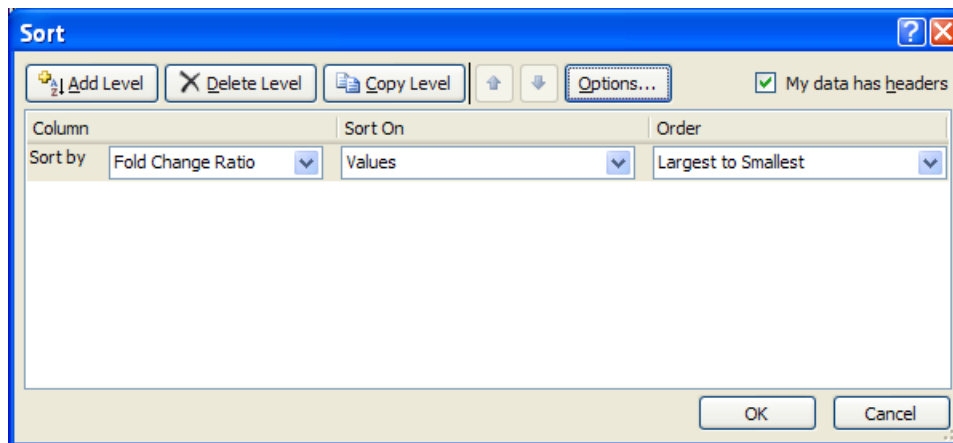


The Sort dialog appears.

11. Select **Fold Change Ratio** in the dropdown list under **Column**.
12. Select **Largest to Smallest** in the dropdown list under **Order**.

Note: If the choices in the dropdown are **A to Z** and **Z to A**, choose **Z to A**.

The Sort dialog now appears as follows:



13. Click **OK**.

Note: If the Sort Warning dialog appears, select **Sort anything that looks like a number, as a number**, then click **OK**.

The rows are now sorted from highest fold change value to lowest. You decide to create a gene signature based on the genes with a fold change value above absolute 10.


14. Select and copy all rows with a fold change value above 10, as follows:
 - a. In column A (the **Analysis** column), click the first cell under the column heading:

	A	B	C	D	E	F	G
1	Analysis	ProbeSet	Fold Change Ratio	p-Value	adjusted p-value	TEA p-Value	Gene
2	DiseaseState => lung adenocarcinoma vs norma	37892_at	28.46	0	0	0.00001	COL11A1
3	DiseaseState => lung adenocarcinoma vs norma	35174_i_at	20.16	0	0	0.00001	EEF1A2
4	DiseaseState => lung adenocarcinoma vs norma	1651_at	18.24	0	0	0.00001	UBE2C
5	DiseaseState => lung adenocarcinoma vs norma	34342_s_at	17.75	0	0	0.00001	SPP1
6	DiseaseState => lung adenocarcinoma vs norma	38582_at	16.83	0	0	0.00001	SPINK1
7	DiseaseState => lung adenocarcinoma vs norma	2092_s_at	15.01	0	0	0.00001	SPP1
8	DiseaseState => lung adenocarcinoma vs norma	37426_at	12.62	0	0.00001	0.00001	TOX3
9	DiseaseState => lung adenocarcinoma vs norma	1599_at	9.98	0	0.00001	0.00001	CDKN3
10	DiseaseState => lung adenocarcinoma vs norma	32154_at	8.2	0	0	0.00001	TFAP2A
11	DiseaseState => lung adenocarcinoma vs norma	38414_at	7.96	0	0	0.00001	CDC20
12	DiseaseState => lung adenocarcinoma vs norma	39677_at	7.82	0	0	0.00001	GIN51

- b. Press and hold down the **Shift** key, then click the cell in column G (the **Gene** column) in the last row with a fold change value above 10:

	A	B	C	D	E	F	G
1	Analysis	ProbeSet	Fold Change Ratio	p-Value	adjusted p-value	TEA p-Value	Gene
2	DiseaseState => lung adenocarcinoma vs norma	37892_at	28.46	0	0	0.00001	COL11A1
3	DiseaseState => lung adenocarcinoma vs norma	35174_i_at	20.16	0	0	0.00001	EEF1A2
4	DiseaseState => lung adenocarcinoma vs norma	1651_at	18.24	0	0	0.00001	UBE2C
5	DiseaseState => lung adenocarcinoma vs norma	34342_s_at	17.75	0	0	0.00001	SPP1
6	DiseaseState => lung adenocarcinoma vs norma	38582_at	16.83	0	0	0.00001	SPINK1
7	DiseaseState => lung adenocarcinoma vs norma	2092_s_at	15.01	0	0	0.00001	SPP1
8	DiseaseState => lung adenocarcinoma vs norma	37426_at	12.62	0	0.00001	0.00001	TOX3
9	DiseaseState => lung adenocarcinoma vs norma	1599_at	9.98	0	0.00001	0.00001	CDKN3
10	DiseaseState => lung adenocarcinoma vs norma	32154_at	8.2	0	0	0.00001	TFAP2A
11	DiseaseState => lung adenocarcinoma vs norma	38414_at	7.96	0	0	0.00001	CDC20
12	DiseaseState => lung adenocarcinoma vs norma	39677_at	7.82	0	0	0.00001	GINS1

- c. Press the **Ctrl + C** keys to copy the selected rows.

15. At the bottom of the Excel window, click the **Insert Worksheet** icon () to open a new worksheet:

	A	B	C	D	E	F	G
1	Analysis	ProbeSet	Fold Change Ratio	p-Value	adjusted p-value	TEA p-Value	Gene
2	DiseaseState => lung adenocarcinoma vs norma	37892_at	28.46	0	0	0.00001	COL11A1
3	DiseaseState => lung adenocarcinoma vs norma	35174_i_at	20.16	0	0	0.00001	EEF1A2
4	DiseaseState => lung adenocarcinoma vs norma	1651_at	18.24	0	0	0.00001	UBE2C
5	DiseaseState => lung adenocarcinoma vs norma	34342_s_at	17.75	0	0	0.00001	SPP1
6	DiseaseState => lung adenocarcinoma vs norma	38582_at	16.83	0	0	0.00001	SPINK1
7	DiseaseState => lung adenocarcinoma vs norma	2092_s_at	15.01	0	0	0.00001	SPP1
8	DiseaseState => lung adenocarcinoma vs norma	37426_at	12.62	0	0.00001	0.00001	TOX3
9	DiseaseState => lung adenocarcinoma vs norma	1599_at	9.98	0	0.00001	0.00001	CDKN3
10	DiseaseState => lung adenocarcinoma vs norma	32154_at	8.2	0	0	0.00001	TFAP2A
11	DiseaseState => lung adenocarcinoma vs norma	38414_at	7.96	0	0	0.00001	CDC20
12	DiseaseState => lung adenocarcinoma vs norma	39677_at	7.82	0	0	0.00001	GINS1
13	DiseaseState => lung adenocarcinoma vs norma	41104_at	7.56	0	0.00002	0.00001	CXCL13
14	DiseaseState => lung adenocarcinoma vs norma	35668_at	7.19	0	0	0.00002	RAMP1
15	DiseaseState => lung adenocarcinoma vs norma	32263_at	6.96	0	0	0.00003	CCNB2
16	DiseaseState => lung adenocarcinoma vs norma	37741_at	6.77	0	0	0.00004	PYCR1
17	DiseaseState => lung adenocarcinoma vs norma	35832_at	6.66	0	0.00003	0.00006	SULF1
18	DiseaseState => lung adenocarcinoma vs norma	40412_at	6.64	0	0	0.00006	PTTG1

sheet1

Select destination and press ENTER or choose Paste

Average: 4.609645714 Count: 49 Sum: 129.0

16. Press the **Ctrl + V** keys to paste the selected rows at the top of the new worksheet.
17. Return to the original worksheet.
18. Repeat Step 9 through Step 14, but this time sort from smallest fold change value to largest.
19. After copying the selected rows with **Ctrl + C**, return to the new worksheet.
20. Paste the selected rows (**Ctrl + V**) in the first empty row below the data you previously pasted.

The new worksheet now looks as follows:

	A	B	C	D	E	F	G
1	DiseaseSt	37892_at	28.46	0	0	0.00001	COL11A1
2	DiseaseSt	35174_i_at	20.16	0	0	0.00001	EEF1A2
3	DiseaseSt	1651_at	18.24	0	0	0.00001	UBE2C
4	DiseaseSt	34342_s_at	17.75	0	0	0.00001	SPP1
5	DiseaseSt	38582_at	16.83	0	0	0.00001	SPINK1
6	DiseaseSt	2092_s_at	15.01	0	0	0.00001	SPP1
7	DiseaseSt	37426_at	12.62	0	0.00001	0.00001	TOX3
8	DiseaseSt	34604_at	-18.08	0	0	0.00001	SLC6A4
9	DiseaseSt	773_at	-16.65	0	0	0.00001	MYH11
10	DiseaseSt	35868_at	-16.63	0	0	0.00001	AGER
11	DiseaseSt	38430_at	-13.92	0	0	0.00001	FABP4
12	DiseaseSt	34174_s_at	-13.22	0	0	0.00001	LPHN2
13	DiseaseSt	32527_at	-12.41	0	0	0.00001	C10orf116
14	DiseaseSt	39066_at	-11.81	0	0	0.00001	MFAP4
15	DiseaseSt	39577_at	-11.74	0	0	0.00001	SOSTDC1
16	DiseaseSt	37777_at	-10.87	0	0	0.00001	PTPRB
17	DiseaseSt	34637_f_at	-10.73	0	0	0.00001	ADH1A
18	DiseaseSt	34708_at	-10.61	0	0	0.00001	FCN3

In the next steps, you will organize the rows of data in the new spreadsheet into a format that can be imported into the gene signature, and then write the formatted data to a text file.

21. Select the column containing the fold change values (column C in the above figure – click the column letter to select the column).
22. Right-click the column of fold change values, then select **Cut**.
23. Select the column immediately to the right of the column containing the gene names.
24. Right-click in the column you just selected, then select **Paste**.

The columns of data should now appear as shown below:

	A	B	C	D	E	F	G	H
1	DiseaseSt	37892_at		0	0	0.00001	COL11A1	28.46
2	DiseaseSt	35174_i_at		0	0	0.00001	EEF1A2	20.16
3	DiseaseSt	1651_at		0	0	0.00001	UBE2C	18.24
4	DiseaseSt	34342_s_at		0	0	0.00001	SPP1	17.75
5	DiseaseSt	38582_at		0	0	0.00001	SPINK1	16.83
6	DiseaseSt	2092_s_at		0	0	0.00001	SPP1	15.01
7	DiseaseSt	37426_at		0	0.00001	0.00001	TOX3	12.62
8	DiseaseSt	34604_at		0	0	0.00001	SLC6A4	-18.08
9	DiseaseSt	773_at		0	0	0.00001	MYH11	-16.65
10	DiseaseSt	35868_at		0	0	0.00001	AGER	-16.63
11	DiseaseSt	38430_at		0	0	0.00001	FABP4	-13.92
12	DiseaseSt	34174_s_at		0	0	0.00001	LPHN2	-13.22
13	DiseaseSt	32527_at		0	0	0.00001	C10orf116	-12.41
14	DiseaseSt	39066_at		0	0	0.00001	MFAP4	-11.81
15	DiseaseSt	39577_at		0	0	0.00001	SOSTDC1	-11.74
16	DiseaseSt	37777_at		0	0	0.00001	PTPRB	-10.87
17	DiseaseSt	34637_f_at		0	0	0.00001	ADH1A	-10.73
18	DiseaseSt	34708_at		0	0	0.00001	FCN3	-10.61

You will now delete all columns except for the **Gene** and **Fold Change Ratio** columns (in the figure above, columns G and H, respectively).

25. Click column A to select it.

Ensure that only column A is selected (highlighted).

26. Press and hold down the **Shift** key, then click the column (F in the above figure) immediately to the left of the **Gene** column.

27. Right-click anywhere in the selected area, then select **Delete**.

The remaining data appears as follows:

	A	B
1	COL11A1	28.46
2	EEF1A2	20.16
3	UBE2C	18.24
4	SPP1	17.75
5	SPINK1	16.83
6	SPP1	15.01
7	TOX3	12.62
8	SLC6A4	-18.08
9	MYH11	-16.65
10	AGER	-16.63
11	FABP4	-13.92
12	LPHN2	-13.22
13	C10orf116	-12.41
14	MFAP4	-11.81
15	SOSTDC1	-11.74
16	PTPRB	-10.87
17	ADH1A	-10.73
18	FCN3	-10.61

28. Click the Office button (), then click **Save As > Other Formats**.

The Save As dialog appears.

29. In the **Save in** field at the top of the dialog, select the root directory on the C:\ drive – for example:

Local Disk (C:)

30. In the **File name** field, type **lung adenocarcinoma vs normal**.

31. In the **Save as type** field, select **Text (Tab delimited) (*.txt)**.

32. Click **Save**.

A warning dialog appears, advising you that you can only save the active worksheet.

33. Click **OK** to acknowledge the warning.

Another warning dialog appears, informing you that some features will be lost when saving to a text file.

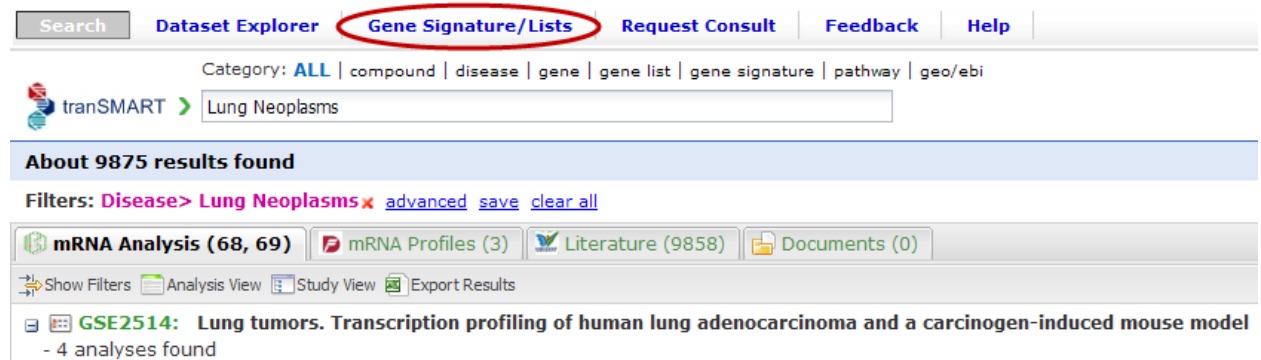
34. Click **Yes** to acknowledge the warning.

35. Close Excel.

36. Click **No** if prompted to save the file. It has already been saved.

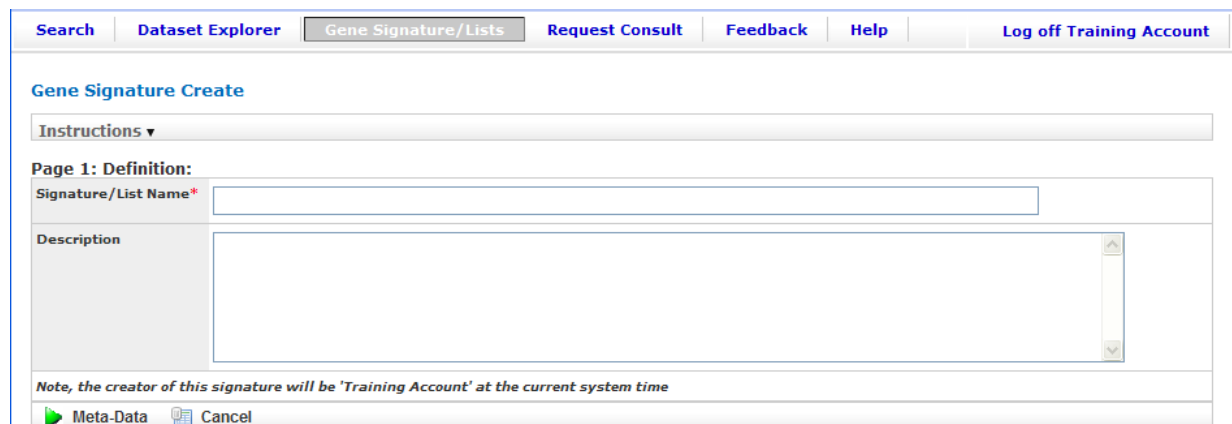
Task 2: Define the Gene Signature and Import the Gene File

1. Click the tranSMART **Gene Signature/Lists** tab to open the Gene Signature tool:

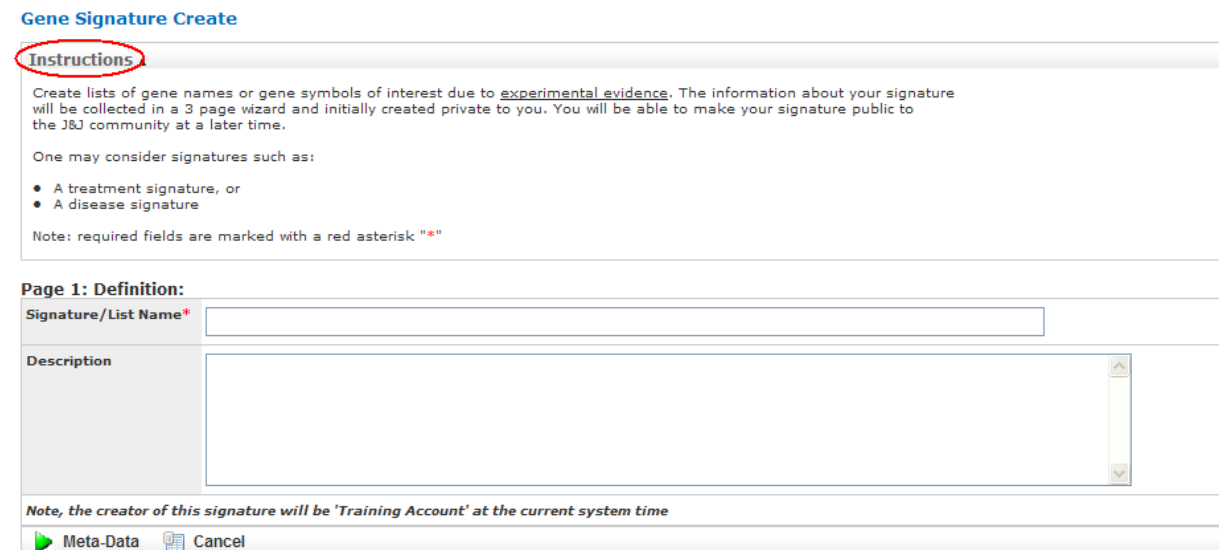


2. Click the **New Signature** button.

The first page of the gene signature wizard appears:



3. Click **Instructions** to read the instructions for creating a gene signature:



4. In **Signature/List Name**, type **<Your Training ID> Gene Signature**.

For example, if your training ID is **publicuser9**, type:

PublicUser9 Gene Signature

5. In **Description**, type the following text:

Genes from lung adenocarcinoma experiment with fold change value above absolute 10.

6. Click **Meta-Data** to proceed to the next wizard page.
7. In **Source of list**, select **Experiment**.
8. In **Owner of data**, select **GEO**.
9. Scrolling down to **PMIDs**, type **16314486**.
10. In **Species**, select **Human**.
11. In **Technology Platform**, select **Affymetrix - HG_U95Av2 [GPL8300]**.
12. In **Tissue Type**, select **Lung**.

The wizard page now appears as follows:

Gene Signature Create

Instructions ▼

Page 2: Meta-Data:

Source of list	Experiment ▼	
Owner of data	GEO ▼	
Stimulus	i.e. LPS, polyIC, etc:	<input type="text"/>
	Dose, units, and time:	<input type="text"/>
Treatment	Drug treatment used in assay:	<input type="text"/>
	Dose, units, and time:	<input type="text"/>
	OR Enter:	
	J&J Compound:	select compound ▼
	Protocol Number:	<input type="text"/>
PMIDs (comma separated)	16314486	
Species*	Human ▼	
Technology Platform*	Affymetrix - HG_U95Av2 [GPL8300] ▼	
Tissue Type	Lung ▼	
Experiment Type	select experiment type ▼	
	If applicable, ATCC designation: <input type="text"/>	

Definition Next Cancel

13. Click **Next** to proceed to the final wizard page.

14. In **P-value Cutoff**, select .05:

Gene Signature Create

Instructions ▾

Page 3: Analysis Meta-Data:

Analysis Performed By:

Normalization Method:

Analysis Info:

Category:

Method:

Multiple Testing Correction Employed? ☐ Yes ☐ No

P-value Cutoff*

File Upload Information (tab delimited text only, no .xls Excel files): [See Samples](#)

File Information*

File schema:

Fold change metric:

Upload File* (tab delimited text files only)

You are now ready to specify the format of the gene file you created in Task 1, and then upload the file into the gene signature.

15. In **File Information**:

- ☐ Select **Gene Symbol <tab> Metric Indicator** in **File schema**.
- ☐ Select **actual fold change** in **Fold change metric**.

16. Click the **Browse** button to the right of the **Upload File** field.

17. In the Choose File dialog, navigate to the **C:** directory and select the file **lung adenocarcinoma vs normal.txt**.

18. Click **Open**.

The path and file name of the **lung adenocarcinoma vs normal.txt** file now appears in the **Upload File** field:

File Upload Information (tab delimited text only, no .xls Excel files): [See Samples](#)

File Information*

File schema:


Fold change metric:

Upload File* (tab delimited text files only)

19. Click **Save** to save the new gene signature.

The new signature now appears in the **My Signatures** section of the gene signature list:

Gene Signature List

My Signatures (1) ▲										
Name	Author	Date Created	Species	Tech Platform	Tissue Type	Public	Gene List	# Genes	# Up-Regulated	# Down-Regulated
 PublicUser Gene Signature	Public Training Account	2010-04-08	Human	GPL8300	Lung	No	No	18	7	11
										-- Select Action -- ▼

20. Click the **Select Action** dropdown at the right of the gene signature entry to see a list of the actions you can perform on the gene signature:

# Up-Regulated	# Down-Regulated	
7	11	-- Select Action -- ▼
		<div> -- Select Action -- Clone Delete Edit Edit Items Excel Download Make Public </div>

Notice that one of the actions is **Make Public**. If a gene signature is private, only you or an administrator can view it and use it as a filter in a tranSMART Search operation (as you will do in the next lesson). If a gene signature is public, anyone can view it and use it as a search filter.

Note: You will use your new gene signature in the next lesson.

Search for Studies Using a New Gene Signature as a Filter

Lesson Goal: Use a newly created gene signature to generate hypotheses in tranSMART.

Scenario: You want to find studies where the differentially regulated genes overlap with the genes contained in your new gene signature. This will generate a set of hypotheses about diseases or treatments that may have similar genes dysregulated, and that can help you develop a further set of experiments.

1. Click the tranSMART **Search** tab to display the Search window.

Search Dataset Explorer Gene Signature/Lists Request Consult Feedback Help Log off Public Training Account

GeneSignature 'PublicUser Gene Signature' was created on: Thu Apr 08 07:51:04 EDT 2010

New Signature

Gene Signature List

My Signatures (1) ▲

Name	Author	Date Created	Species	Tech Platform	Tissue Type	Public	Gene List	# Genes	# Up-Regulated	# Down-Regulated	
PublicUser Gene Signature	Public Training Account	2010-04-08	Human	GPL8300	Lung	No	No	18	7	11	-- Select Action -- ▼

2. Type **public** in the search field:

Search tranSMART

Category: ALL | compound | disease | gene | gene list | gene signature | pathway | geo/ebi

public

Search browse saved filters

Gene Signature>Internal> PublicUser Gene Signature

Gene List>Internal> PublicUser Gene Signature

Note: The following steps use **publicuser** as the login ID that is used as part of the gene signature name. Substitute your login ID for **publicuser**.

3. Click **Gene Signature>Internal> PublicUser Gene Signature** in the dropdown list.

The search results include studies involving genes that matched the ones in your gene signature:

About 528 results found

Filters: **Genesig> PublicUser Gene Signature** [advanced](#) [save](#) [clear all](#)

mRNA Analysis (513, 953) mRNA Profiles (83) Literature (90) Documents (0)

Show Filters Analysis View Study View Export Results

Analysis result: 513 analyses from 231 experiment(s) [432 Significant TEA / 81 Insignificant TEA]
Note, only significant TEA Analyses are displayed!

1 2 3 4 5 6 7 8 9 10 .. 44 Next

[**co-regulated** ∞] Excel

GSE15245 - DiseaseState => Definite MS vs CIS

BioMarkers (25 signature/pathway genes matched):

You are only interested in studies related to lung neoplasms, so you want to filter the results further.

4. Click **advanced**:

About 528 results found

Filters: **Genesig> PublicUser Gene Signature** **advanced** [save](#) [clear all](#)

mRNA Analysis (513, 953) mRNA Profiles (83) Literature (90) Documents (0)

Show Filters Analysis View Study View Export Results

Analysis result: 513 analyses from 231 experiment(s) [432 Significant TEA / 81 Insignificant TEA]
Note, only significant TEA Analyses are displayed!

The Edit Filters dialog appears:

Edit Filters

Category: **ALL** | compound | disease | gene | gene list | gene signature | pathway | geo/ebi

Gene Signature> **PublicUser Gene Signature**

Apply Cancel

5. Type **lung neo** in the search field.
6. Click **Disease> Lung Neoplasms** in the dropdown list.

7. Click **Apply**.

The following figure shows a portion of the mRNA Analysis results. The analyses that are returned involve both lung neoplasms and one or more genes in your gene signature:

About 14 results found

Filters: **Genesig> PublicUser Gene Signature** **AND Disease> Lung Neoplasms** [advanced](#) [save](#) [clear all](#)

mRNA Analysis (20, 28) mRNA Profiles (3) Literature (0) Documents (0)

Show Filters Analysis View Study View Export Results

Analysis result: 20 analyses from 11 experiment(s) [17 Significant TEA / 3 Insignificant TEA]
Note, only significant TEA Analyses are displayed!

1 2 Next

co-regulated

GSE2514 - DiseaseState => lung adenocarcinoma vs normal

BioMarkers (24 signature/pathway genes matched):

co-regulated

GSE7670 - DiseaseState => Tumor vs normal

BioMarkers (21 signature/pathway genes matched):

co-regulated

GSE3268 - DiseaseState => lung cancer vs normal

BioMarkers (11 signature/pathway genes matched):

anti-regulated **21.831**

GSE10245 - DiseaseState => squamous cell carcinoma vs adenocarcinoma

BioMarkers (5 signature/pathway genes matched):

Now you want to browse through the analyses to see how the genes that match those in your gene signature behaved during the experiments.

In the figure above, the analysis **lung adenocarcinoma vs normal** in experiment **GSE2514** is the first in the list. This is the analysis from which you derived the genes for your gene signature during the previous exercise.

8. Click the **+** icon () to the left of the label **BioMarkers** for the analysis **lung adenocarcinoma vs normal**. A partial list of the matching genes is shown below:

co-regulated					
GSE2514 - DiseaseState => lung adenocarcinoma vs normal					
BioMarkers (24 signature/pathway genes matched):					
AGER	ProbeSet: 35868_at	Gene: AGER	Fold Change: -16.63	p-Value: 0.00	TEA p-Value: 0.00001
EEF1A2	ProbeSet: 35174_i_at	Gene: EEF1A2	Fold Change: 20.16	p-Value: 0.00	TEA p-Value: 0.00001
C10orf116	ProbeSet: 32527_at	Gene: C10orf116	Fold Change: -12.41	p-Value: 0.00	TEA p-Value: 0.00001
SOSTDC1	ProbeSet: 39577_at	Gene: SOSTDC1	Fold Change: -11.74	p-Value: 0.00	TEA p-Value: 0.00001
ADH1A	ProbeSet: 34637_f_at	Gene: ADH1A	Fold Change: -10.73	p-Value: 0.00	TEA p-Value: 0.00001
SPINK1	ProbeSet: 38582_at	Gene: SPINK1	Fold Change: 16.83	p-Value: 0.00	TEA p-Value: 0.00001
FABP4	ProbeSet: 38430_at	Gene: FABP4	Fold Change: -13.92	p-Value: 0.00	TEA p-Value: 0.00001

Notice that the list of genes in the analysis contains the same genes as in your gene signature, the same gene expression values (for values above absolute 10), and the same probe sets that produced the expression results in your gene signature.

You now want to browse through the list of genes in the other analyses. You are interested in finding those genes whose expressions were produced by the same probe set as in your gene signature.

If you compared the genes and probe sets for each of the analyses to the ones in your gene signature, you would find that the most strongly up-regulated gene in your gene signature, COL11A1, is the only gene that is associated with the same probe set (37892_at) in both your gene signature and in any of the analyses.

You decide to view a profile of COL11A1 to learn more about it.

9. Click the **mRNA Profiles** tab.
10. Select **COL11A1** from the **Gene** dropdown.
11. Select **Lung Neoplasms** from the **Disease** dropdown.
12. Select **37892_at** from the **Probe Set** dropdown.

This is the probe set that produced the expression results for COL11A1 in your pathway and in some of the analyses.

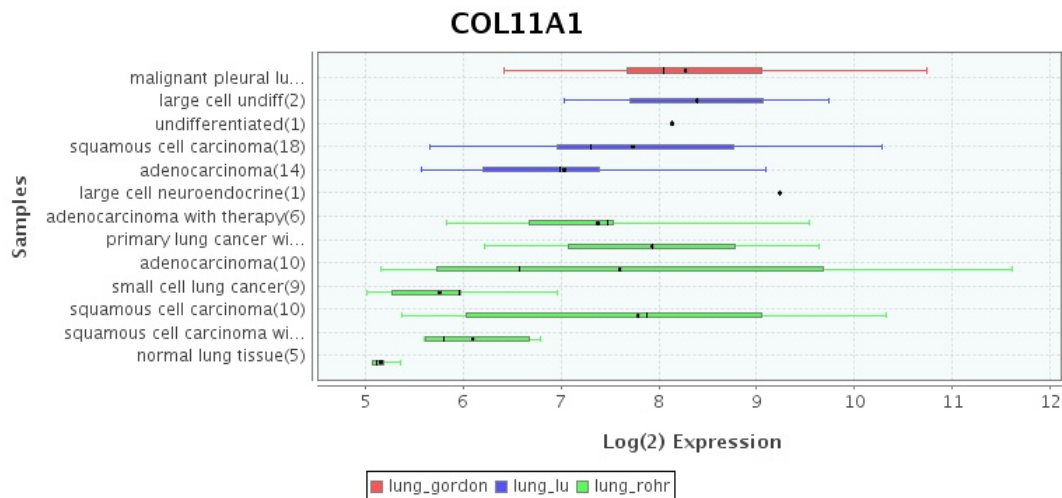
The following chart and the datasets related to the COL11A1 gene are displayed for your further study:

Filter (Note: search found 3 studies)

Gene: COL11A1

Disease: Lung Neoplasms

Probe Set: 37892_at



Print Chart

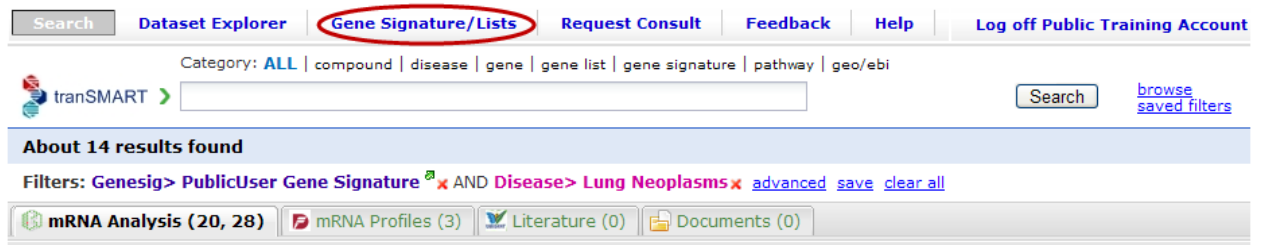
Information on individual datasets

Dataset	No. Samples	Experiment
malignant pleural lung mesothelioma	33	lung_gordon : Translation of microarray data into clinically relevant cancer diagnostic tests
large cell undiff	2	lung_lu : A gene expression signature predicts survival of patients with stage I non-small
undifferentiated	1	lung_lu : A gene expression signature predicts survival of patients with stage I non-small
squamous cell carcinoma	18	lung_lu : A gene expression signature predicts survival of patients with stage I non-small
adenocarcinoma	14	lung_lu : A gene expression signature predicts survival of patients with stage I non-small
squamous cell carcinoma	18	lung_lu : A gene expression signature predicts survival of patients with stage I non-small
large cell neuroendocrine	1	lung_lu : A gene expression signature predicts survival of patients with stage I non-small
adenocarcinoma with therapy	6	lung_rohr : Genetic programming and gene expression profiling for molecular discriminator
primary lung cancer with therapy	2	lung_rohr : Genetic programming and gene expression profiling for molecular discriminator
adenocarcinoma	10	lung_rohr : Genetic programming and gene expression profiling for molecular discriminator
small cell lung cancer	9	lung_rohr : Genetic programming and gene expression profiling for molecular discriminator
squamous cell carcinoma	10	lung_rohr : Genetic programming and gene expression profiling for molecular discriminator
squamous cell carcinoma with therapy	5	lung_rohr : Genetic programming and gene expression profiling for molecular discriminator
normal lung tissue	5	lung_rohr : Genetic programming and gene expression profiling for molecular discriminator

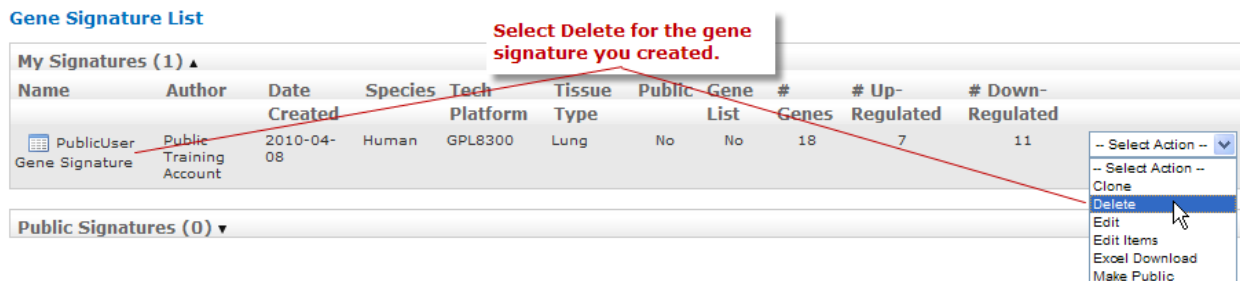
Data courtesy of **Dana-Farber Cancer Institute GeneChip Oncology Database**.

When finished, you want to delete the gene signature you created.

13. Click **Gene Signature/Lists**:



14. Select **Delete** in the **Select Action** dropdown for your gene signature:



15. Click **OK** to confirm the deletion.

Use a Heat Map to Compare Treatment Results

Lesson Goal: Use the tranSMART Dataset Explorer to create a heat map, and save the heat map data to a file.

Scenario: You are analyzing the results of a study on rheumatoid arthritis. You want to see a visualization of gene expression data for the gene REL in two cohorts: those who responded to anti-Tnf therapy and those who did not.

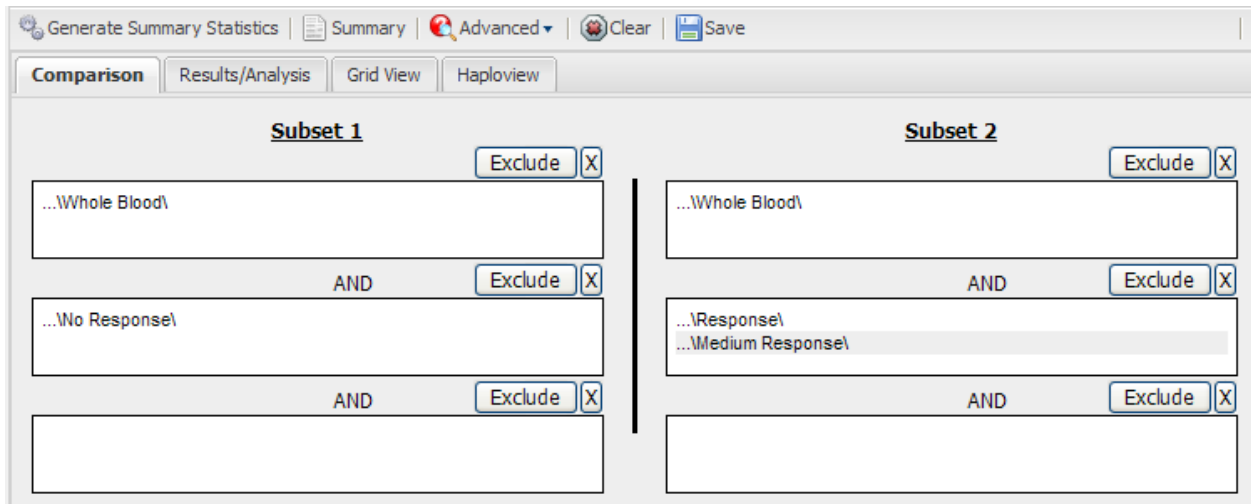
Note: For a heat map to be generated in Dataset Explorer, at least one of the subsets must contain the following elements:

- One or more test subjects.
- A platform of biomarkers (for example, a gene pathway or RBM antigens).

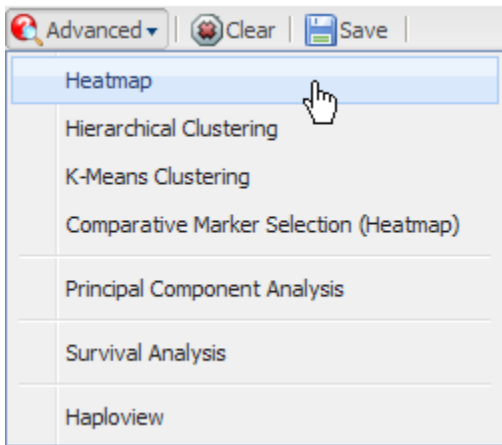
Note: The heat map illustrated in this section was generated with Internet Explorer version 8. If you are using a different browser or a different version of Internet Explorer, you might see slight differences between the illustration and the heat map displayed on your screen.

1. Click the tranSMART **Dataset Explorer** tab to display the Dataset Explorer window.
2. In the left pane of Dataset Explorer, click the **Navigate Terms** tab.
3. In **Public Studies**, open the study **Bienkowska_RheumatoidArthritis_GSE15258**.
4. Open the following nested nodes in the following order:
 - a. Biomarker Data
 - b. Affymetrix GeneChip Human Genome U133 Plus 20 Array
5. Drag **Whole Blood** into subset definition boxes in Subset 1 and Subset 2.
6. Open the following nested nodes in the following order:
 - a. Clinical Data
 - b. Response To Anti Tnf Therapy
7. Drag **No Response** into an empty box in Subset 1.
8. Drag **Response** into an empty box in Subset 2.
9. Drag **Medium Response** into the same box where you placed **Response**.

The subset definition boxes look as follows:



10. Click the **Advanced** tab, then click **Heatmap**:



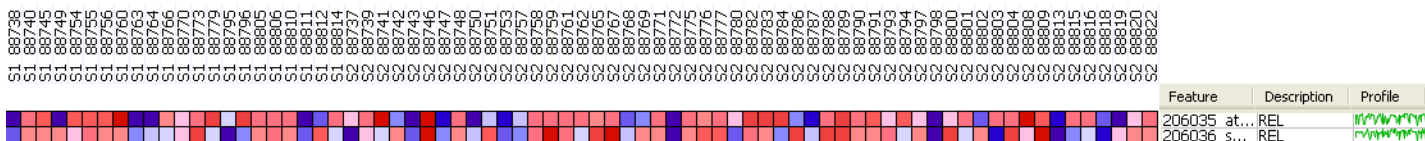
The Compare Subsets-Pathway Selection dialog appears.

11. Type **rel** in the **Select a Gene/Pathway** field.

12. Click **Gene> REL** in the dropdown list.

13. Click **Run Workflow**.

In a few seconds, the heat map appears in a new browser window:



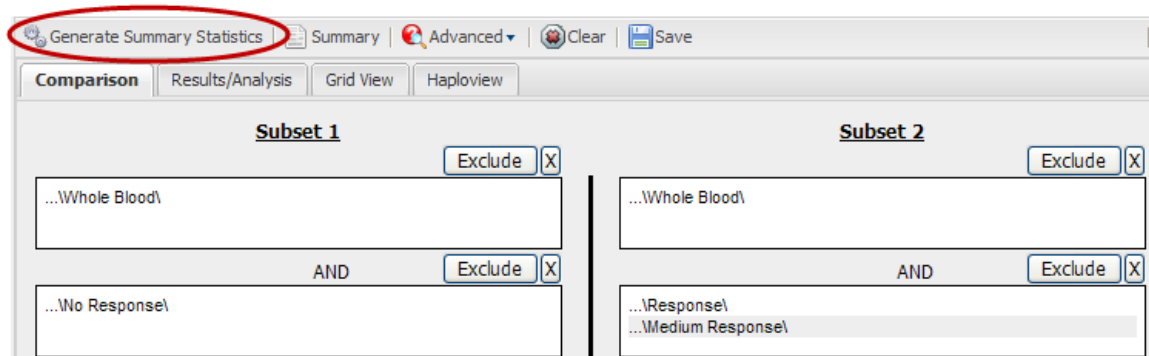
Notice that:

- The column headings represent the subjects in the study.
 - The prefix S1_ represents Subset 1 subjects. Prefix S2_ represents Subset 2 subjects. The numbers following the prefixes are the IDs of the subjects.
 - REL expression data is represented by the colored cells – up-regulation is expressed in shades of red. Down-regulation is expressed in shades of blue.
 - In this example, two probe sets were used, yielding two sets of REL data.
14. When finished comparing the biomarker metrics, close the browser window containing the heat map.

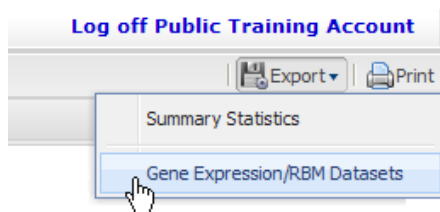
You decide to save the REL expression data to a Microsoft Excel spreadsheet.

Note: With private studies, you must have permission from the study owner to save study data to a file.

15. Click **Generate Summary Statistics**:



16. Click **Export**, then click **Gene Expression/RBM Datasets**:



17. Click **Open** in the File Download dialog.

REL data observed in several of the subjects in Subset 1 appear below:

	A	B	C	D	E	F	G	H	I	J	K	L
1	NAME	Descriptio	S1_88738	S1_88740	S1_88745	S1_88749	S1_88754	S1_88755	S1_88756	S1_88760	S1_88763	S1_88764
2	206035_at	REL	-2.47133	0.13106	0.27076	-2.35473	0.24986	0.27869	0.23062	0.80504	-2.3824	-2.47734
3	206036_s	REL	-1.54764	0.04338	0.10132	0.04961	-0.25269	0.34217	0.1117	0.02272	-1.38417	-0.66737

18. Close the Excel file without saving it.

Normally you would save the file for future reference.

Analyze Gene Expression Data from Different Perspectives

Lesson Goal: Create different heat map visualizations of study data.

Scenario: You want to analyze gene expression data for the gene IL6R, collected in a study of patients diagnosed with multiple myeloma. You are particularly interested in the data collected from the study's proliferation group.

In the previous lesson, you created a standard heat map. A standard heat map organizes its data points according to the numeric order of the IDs of the subjects in the subset(s).

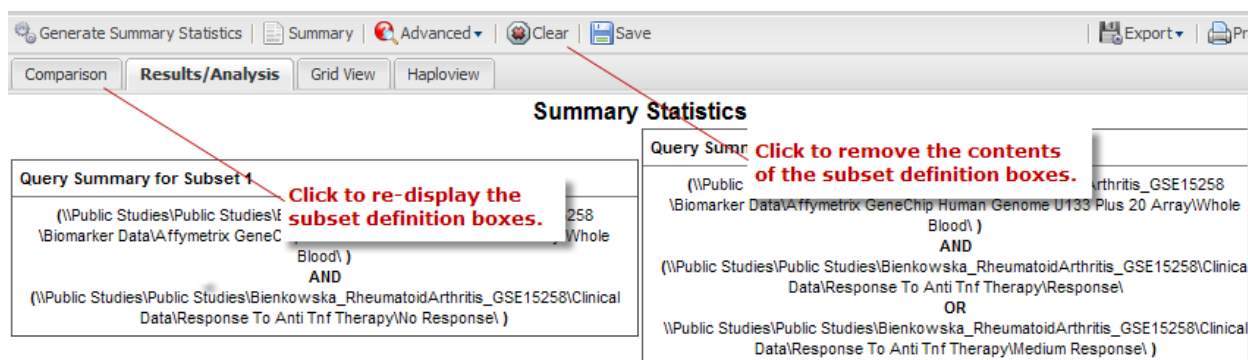
You can also organize a heat map's data points by their gene expression values. The following types of heat maps are organized by gene expression values:

- Class discovery (hierarchical clustering) heat map – A visualization of patterns of related data points in gene expression and RBM data.
- Class discovery (k-means clustering) heat map – A visualization of groupings of the most closely related data points, based on the number of groupings you specify.
- Differential Analysis/Marker Selection heat map – A visualization of differentially expressed genes in distinct phenotypes.

You will generate all these types of heat maps in this lesson.

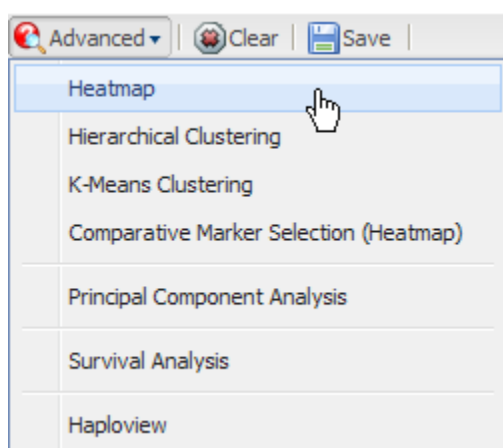
Note: The heat maps illustrated in this section were generated with Internet Explorer version 8. If you are using a different browser or a different version of Internet Explorer, you might see slight differences between the illustrations and the heat maps displayed on your screen.

1. Click the Dataset Explorer **Comparison** tab, then click the **Clear** button. These actions clear the subset definitions and results from the previous lesson.



2. In **Public Studies**, open the study **Shaughnessy_MultipleMyeloma_GSE2658**.

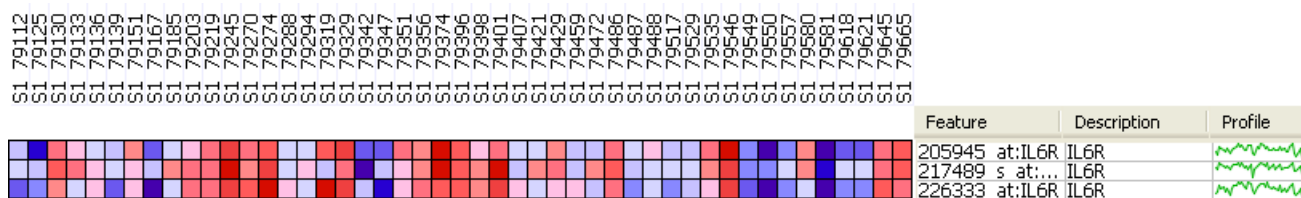
3. Open the following nested nodes in the following order:
 - a. Biomarker Data
 - b. Affymetrix GeneChip Human Genome U133 Plus 20 Array
4. Drag **Bone Marrow** into a subset definition box in Subset 1.
5. Open the following nested nodes in the following order:
 - a. Published Conclusions
 - b. Disease Subtype Classification
 - c. RNA
6. Drag **Proliferation group** into an empty box in Subset 1.
7. Click the **Advanced** tab, then click **Heatmap** to generate a standard heat map:



The Compare Subsets-Pathway Selection dialog appears.

8. Type **il6r** in the **Select a Gene/Pathway** field.
9. Click **Gene> IL6R** in the dropdown list.
10. Click **Run Workflow**.

In a few seconds, the standard heat map appears in a new browser window:



Leave the heat map window open, so you can compare it with the other heat maps you will generate. Shades of color are more easily distinguished on your screen than on the black-and-white printed page.

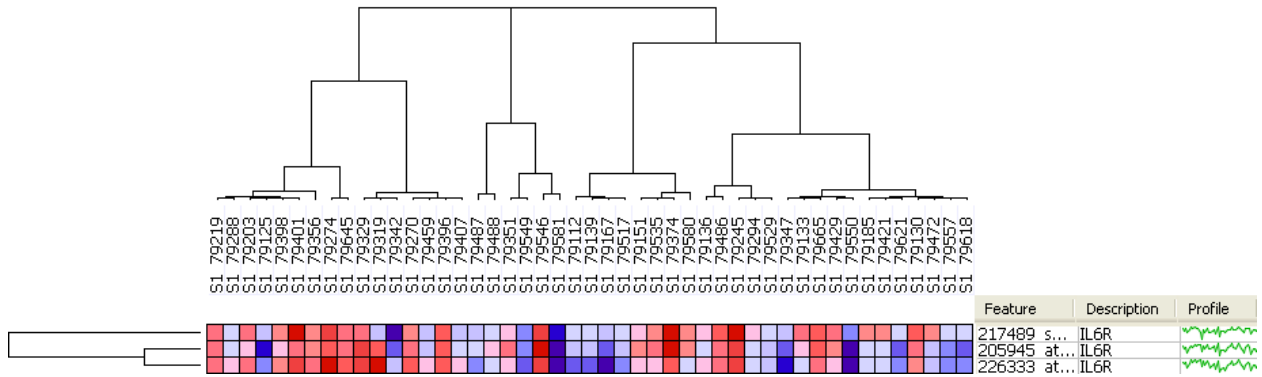
Now you want to organize the same data in hierarchical patterns of related data points.

11. Click the **Advanced** tab, then click **Hierarchical Clustering**.

The Compare Subsets-Pathway Selection dialog appears, with IL6R already in the **Select a Gene/Pathway** field.

12. Click **Run Workflow**.

The hierarchical clustering heat map appears in a new browser window:



Now you want to organize the same data in three clusters of the most related data points:

13. Click the **Advanced** tab, then click **K-Means Clustering**.

The Compare Subsets-Pathway Selection dialog again appears with IL6R in the **Select a Gene/Pathway** field, but it now contains a new field, **Select the number of Clusters**:

Compare Subsets-Pathway Selection

SUBSET 1	SUBSET 2
Platform: MRNA	Platform:
GPL Platform: Affymetrix GeneChip Human Ger	GPL Platform:
Sample: Bone Marrow	Sample:
Tissue Type:	Tissue Type:
Timepoint:	Timepoint:
Select a Gene/Pathway: IL6R	
Select the number of Clusters: 2	
<div>Run Workflow</div> <div>Cancel</div>	

14. Type **3** in the **Select the number of Clusters** field, overwriting the default value of **2**.

15. Click **Run Workflow**.

The k-means clustering heat map appears in a new browser window:



Finally, you want to generate a heat map of differentially expressed genes in the proliferation group as compared with another set of phenotypes that you define.

16. Return to the node Biomarker Data > Affymetrix GeneChip Human Genome U133 Plus 20 Array, and drag **Bone Marrow** into a subset definition box in Subset 2.

17. Return to the node Published Conclusions > Disease Subtype Classification > RNA, and drag **Cyclin D1 deregulation group** into an empty box in Subset 2.

18. Drag **Cyclin D2 deregulation group** into the same box.

The subset definition boxes now look as follows:

19. Click the **Advanced** tab, then click **Comparative Marker Selection (Heatmap)**.

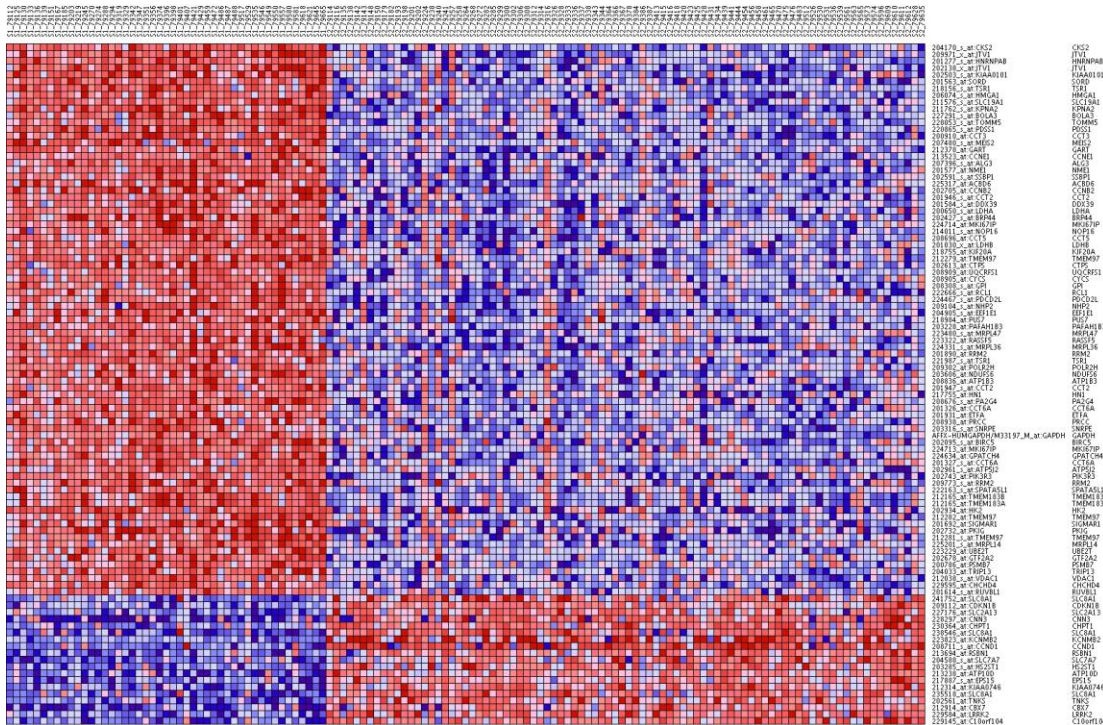
This time, when the Compare Subsets-Pathway Selection appears, the **Select a Gene/Pathway** field is not on it. With this type of heat map, tranSMART searches for *all differentially expressed genes* between the two subsets.

20. Click **Run Workflow**.

Note: Due to the large amount of data being searched, the heat map may take several minutes to appear. Meanwhile, explore the interactive features of heat maps. See the section [Interactive Heat Maps](#) on page 33.

Note: If you are using Internet Explorer 8 and a window labeled **Upregulated features** appears, click **View > Heatmap** to display the heat map.

A miniaturized version of the heat map appears below:



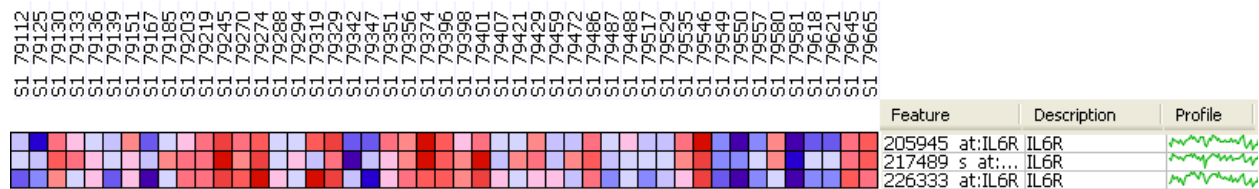
21. Keep the standard heat map open (the one you generated in Step 10), but close the other heat map windows. You will use the standard heat map in the next section.

Interactive Heat Maps

Dataset Explorer heat maps generated with Internet Explorer versions below version 8 are static. With Internet Explorer 8, you can generate interactive heat maps.

Continue with this section if you are using Internet Explorer 8.

Return to the standard heat map you created in Step 2 through Step 10 of the previous section. The heat map looks as follows:



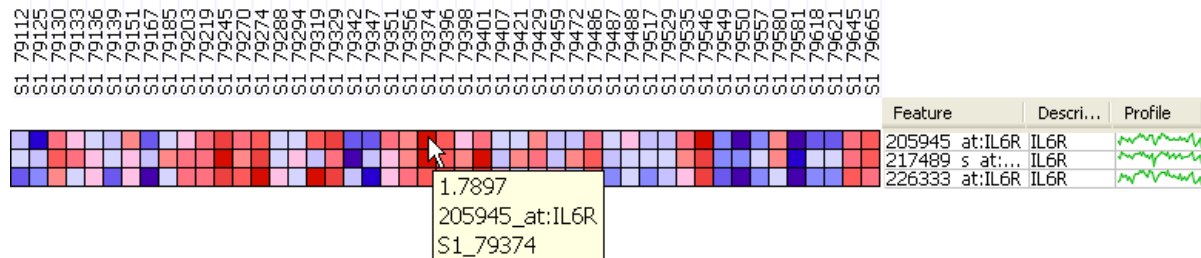
The following sections illustrate some of the features and visualizations available with an interactive heat map.

Note: A standard heat map is used in the following examples, but the same features apply to all interactive heat maps you generate with Dataset Explorer.

View a Particular Data Point

To view a particular data point value:

- Hover the mouse pointer over the cell representing the data point of interest:



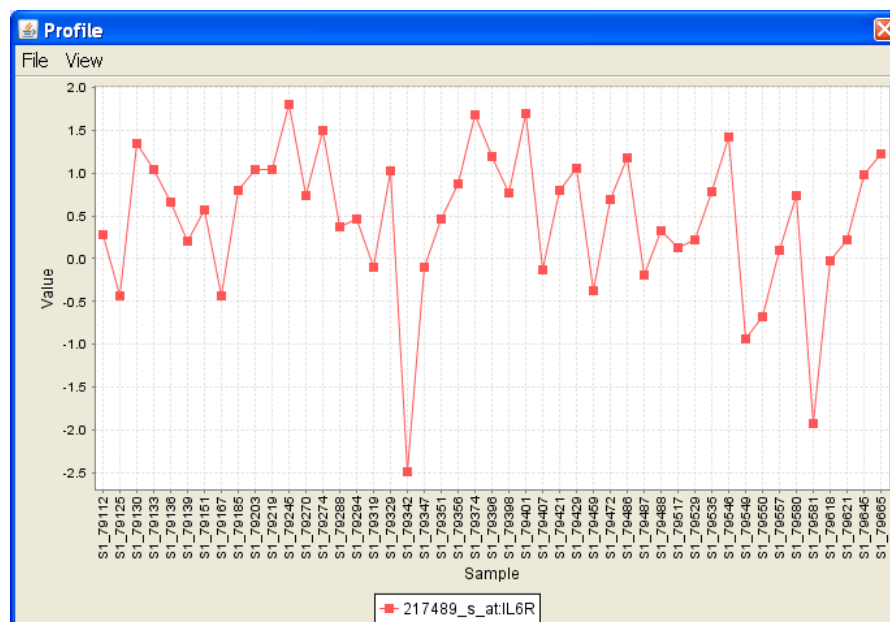
View a Profile of All Data Points for a Probe Set

To view a profile of all data points for a given probe set:

1. In the **Feature** column to the right of the heat map, click probe set **217489_s_at** to select it.
2. Click the green line in the **Profile** column for the selected probe set:



The following profile chart appears:

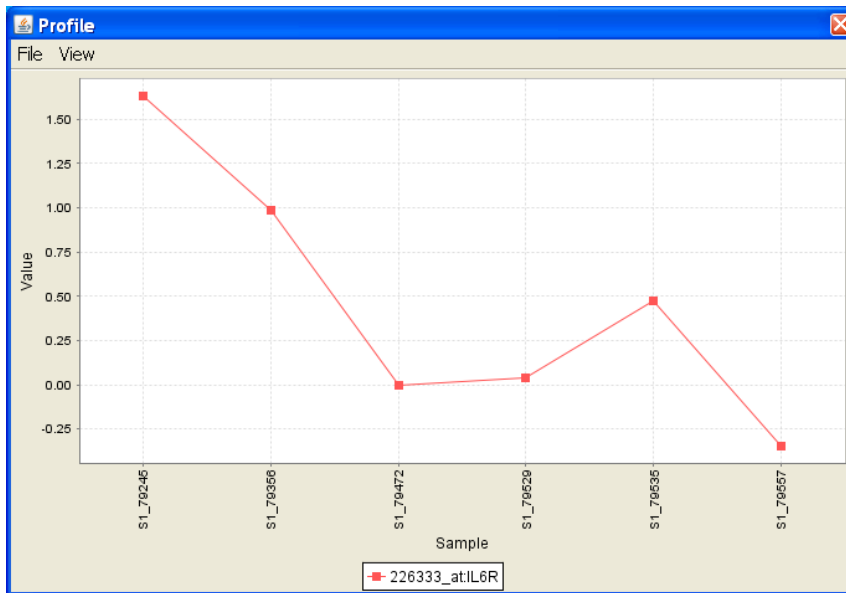


3. When finished analyzing the chart, close the chart window.

View a Profile of Selected Data Points for a Probe Set

To view a profile of selected data points for a given probe set:

1. Hold down the **Ctrl** key, then click the ID of each subject whose data point you want to include in the profile:
2. Click probe set **226333_at** to select it.
3. Click the green line in the **Profile** column for the selected probe set:

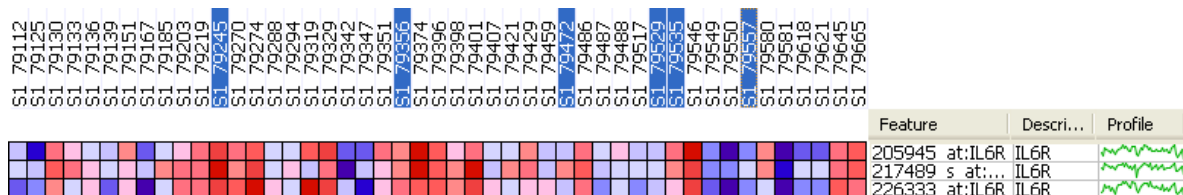


4. When finished analyzing the chart, close the chart window.

View a Profile of Selected Data Points for Multiple Probe Sets

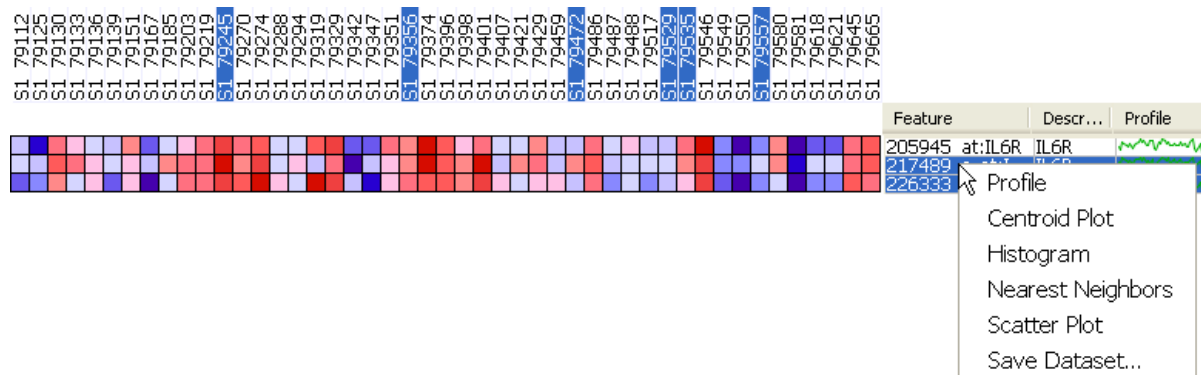
To view a profile of selected data points for multiple probe sets:

1. Hold down the **Ctrl** key, then click the ID of each subject whose data point you want to include in the profile:



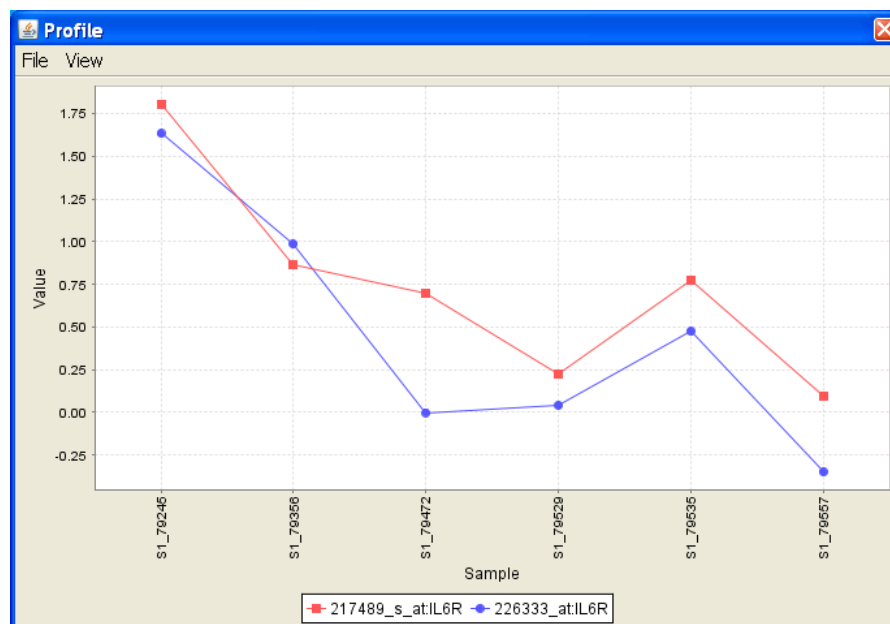
2. Click probe set **217489_s_at** to select it.
3. Hold down the **Ctrl** key, then click probe set **226333_at** to select it:

4. Right-click either of the selected probe sets:



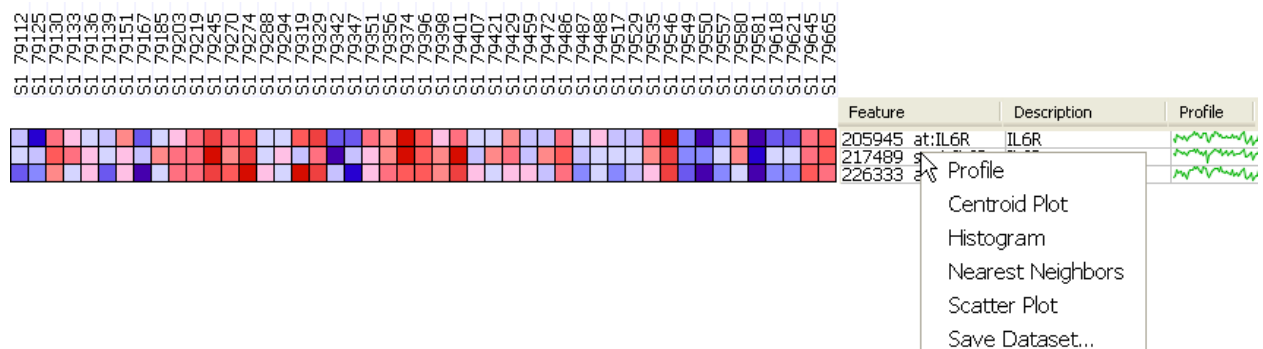
5. Click **Profile**.

The selected data points for each probe set are profiled:



6. When finished analyzing the chart, close the chart window.

7. In the heat map window, right-click any probe set to view the other operations you can perform:



Perform a Survival Analysis

Lesson Goal: Generate a survival analysis visualization and statistics.

Scenario: You hypothesize that breast cancer patients with positive estrogen receptors have a better overall survival rate than breast cancer patients with negative estrogen receptors. You want to test that hypothesis against data from a study in Dataset Explorer.

To generate a survival analysis in Dataset Explorer, you must introduce the following criteria into the two groups you are comparing:

- The observed survival times of the individuals in each group.
- The specific event (death) being tracked for the individuals in the study, and optionally, any censoring factors that occurred before the event took place.

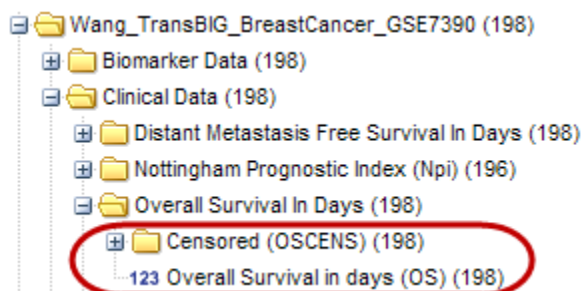
A censoring factor might be the withdrawal of an individual from the study, or the conclusion of the study before the event occurred for a given individual.

- At least one variable that distinguishes the two groups – in this case, positive vs negative estrogen receptors.

To generate the survival analysis:

1. Click the Dataset Explorer **Comparison** tab, then click the **Clear** button to remove the subset definitions from the previous lesson.
2. In **Public Studies**, open the study **Wang_TransBIG_BreastCancer_GSE7390**.
3. Open the following nested nodes in the following order:
 - a. Clinical Data
 - b. Overall Survival in Days

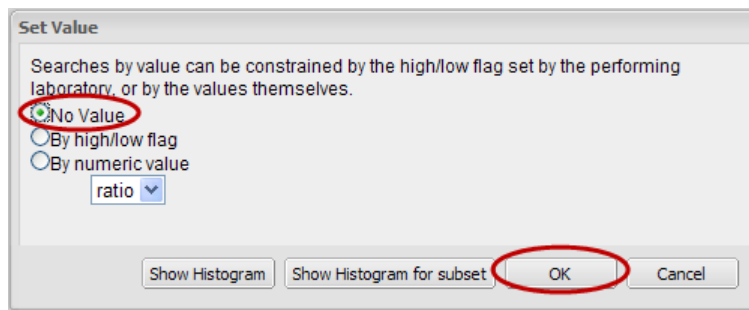
The folder **Overall Survival in Days** contains two items. In the following steps, you will drag them both into subset boxes.



4. Drag the time-to-event dataset **Overall Survival in days (OS)** into Subset 1.

Don't drag the folder into the subset box – just the dataset included in the circle above.

- In the Set Value dialog, select **No Value**, then click **OK**:



Specifying a value limits the values in the dataset. By not specifying a value, the entire time-to-event dataset is used.

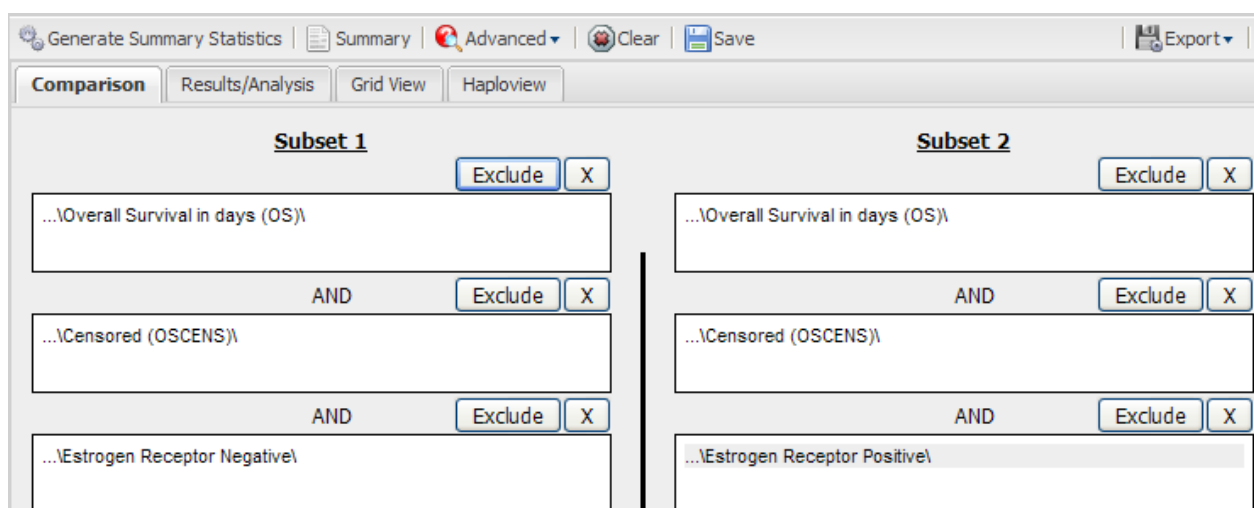
- Repeat step 4 and step 5 for Subset 2.
- In the same **Clinical Data** node, drag the **Censored (OSCENS)** folder into empty boxes in Subset 1 and Subset 2.

The contents of the **Censored (OSCENS)** folder are **No** and **Yes**. This concept introduces the **Event** and **Censored** datasets into the analysis.

Now you will introduce the variable whose effect on survivability you want to test.

- Open the following nested nodes in the following order:
 - Subjects
 - Medical History
 - Estrogen Receptor Status
- Drag **Estrogen Receptor Negative** into an empty box in Subset 1.
- Drag **Estrogen Receptor Positive** into an empty box in Subset 2.

The subset boxes are now defined as follows:



11. Click the **Advanced** tab, then click **Survival Analysis**.

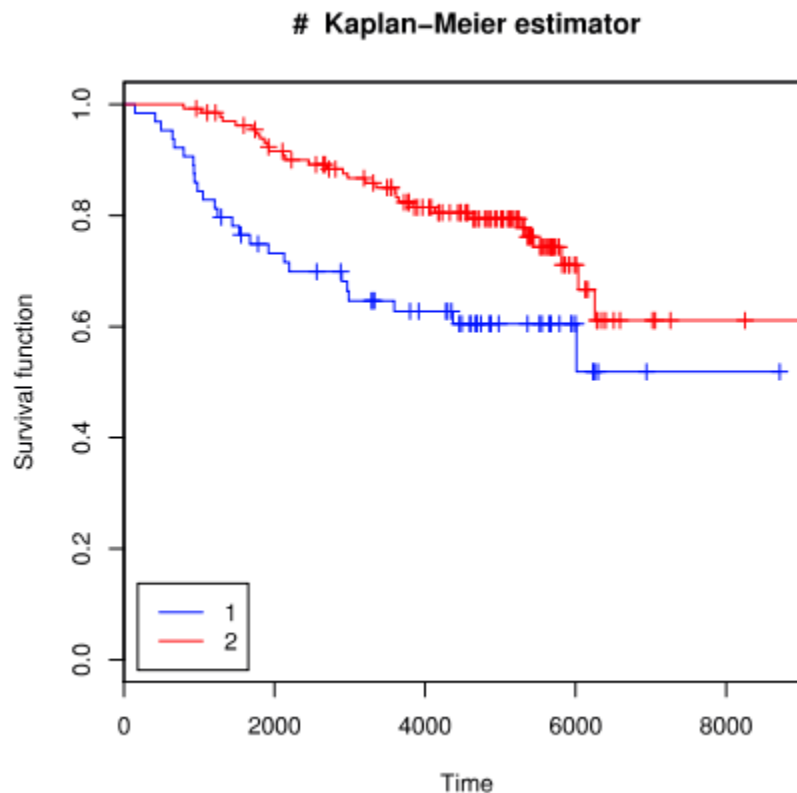
In a few seconds, the Survival Analysis window appears. It has four sections:

- At the top, a summary of the subset definitions.
- Next are two tables that contain the survival analysis statistics.

In the first table, the Hazard Ratio and Relative Risk statistics compare Subset 2 results against Subset 1 results. In this example, Subset 2 patients (those with positive estrogen receptors) have a lower hazard ratio and lower relative risk than those in Subset 1 (those with negative estrogen receptors).

Number of Subjects	198
Hazard Ratio (95% CI)	0.476 (0.281 - 0.806)
Relative Risk (p Value)	-0.743 (0.0058)

- The final section of the Survival Analysis window contains Kaplan-Meier curves of the survival times of each group:



In the figure, the x-axis represents survival time in days, and the y-axis represents the percentage of subjects who were still alive at a given point in time during the study.

Note the hashmarks in the plot lines. These represent censored data – for example, subjects who dropped out of the study before the event (death) occurred.

12. When finished viewing the analysis, close the Survival Analysis window.

Perform a Principal Component Analysis

Lesson Goal: Generate a principal component analysis (PCA) visualization and statistics.

Scenario: You are interested in a study on the effect of strenuous exercise on neutrophils in the 12 healthy male subjects. You want to see if exercise causes significant change on the gene expression profiles of these 12 subjects, before and after exercise.

It is expected that for one subject, one set of genes will show some changes, and another set of genes will show changes for another subject. Direct comparisons cannot answer the question if the group's gene profiles change after exercise. PCA analyzes the multiple dimensional data and finds several linear combinations of the original dimensions to represent the most variance in the data. Each linear combination of the original dimension is named a "Principal Component." PCA can be used to provide multiple snapshots of the data, and show if the two groups of data have different distributions on the graph.

1. Click the Dataset Explorer **Comparison** tab, then click the **Clear** button to remove the subset definitions from the previous lesson.
2. In **Public Studies**, open the study **Radom-Azik_Exercise_GSE8668**.
3. Open the following nested nodes in the following order:
 - a. Biomarker Data
 - b. Affymetrix GeneChip Human Genome U133 Plus 20 Array
 - c. Neutrophils
4. Drag **Before exercise** into an empty box in Subset 1.
5. Drag **After exercise** into an empty box in Subset 2.
6. Click the **Advanced** tab, then click **Principal Component Analysis**.

The Compare Subsets-Pathway Selection dialog appears, with its fields already filled out with the default values for the analysis:

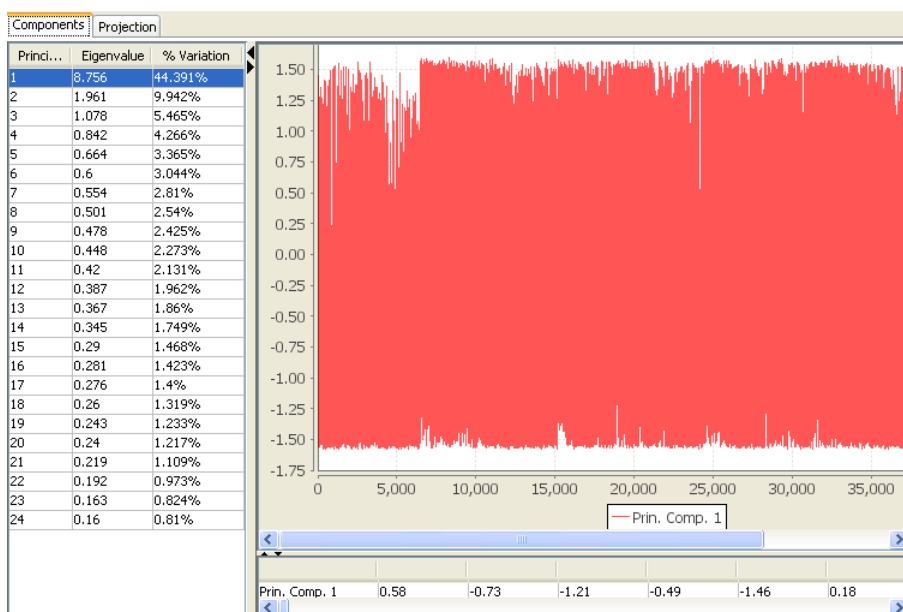
SUBSET 1	SUBSET 2
Platform: MRNA	Platform: MRNA
GPL Platform: Affymetrix GeneChip Human Ger	GPL Platform: Affymetrix GeneChip Human Ger
Sample: Neutrophils	Sample: Neutrophils
Tissue Type:	Tissue Type:
Timepoint: Before exercise	Timepoint: After exercise

Run Workflow Cancel

7. Click **Run Workflow**.

Due to the large amount of gene expression data being processed, it may take a minute or two for the PCA Viewer window to appear. When the window does appear, you see the **Components** tab by default. Take a minute to familiarize yourself with its contents:

- The left pane contains a table listing the principal components. The components are listed in order of variability of data along this vector (the combination of original dimensions), from the greatest variability of these measurements within a component to the least.
- The right pane contains a visualization of the eigenvectors for each principal component selected in the table on the left.
- Below the visualization is a table of the gene expression measurements represented in the visualization for each selected component.

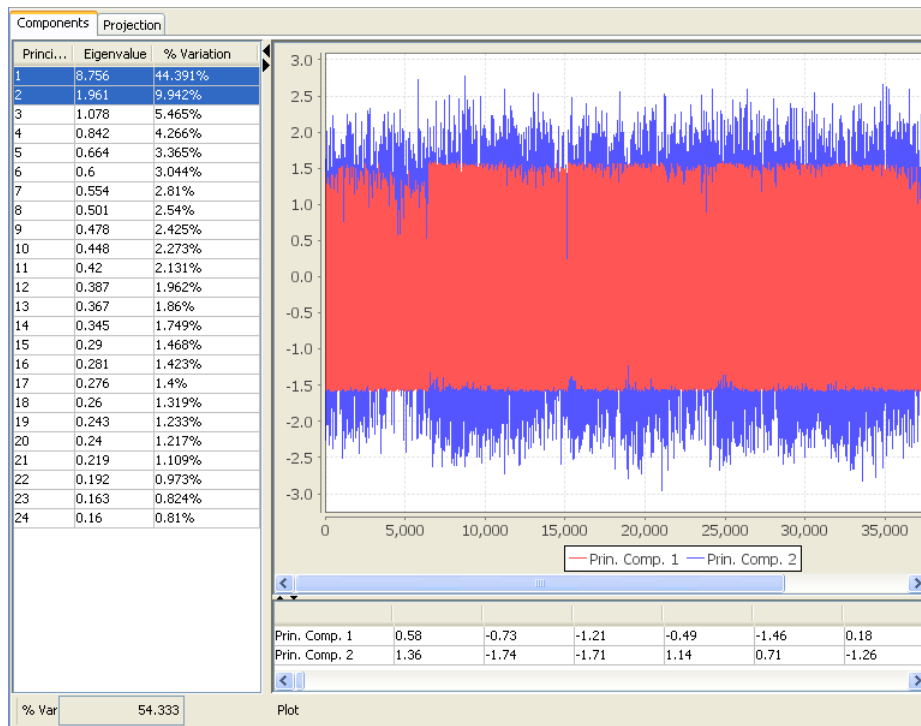


By default, the visualization shows the values for the first principal component only. You want to view a visualization of the first and second principal components.

8. While holding down the **Ctrl** key, click the second principal component in the table on the left.

9. Click the **Plot** button at the bottom of the window.

The visualization presents the selected components in different colors, and the table below the visualization now contains data from both components:



You can generate a visualization of any combination of principal components that you select in the table on the left.

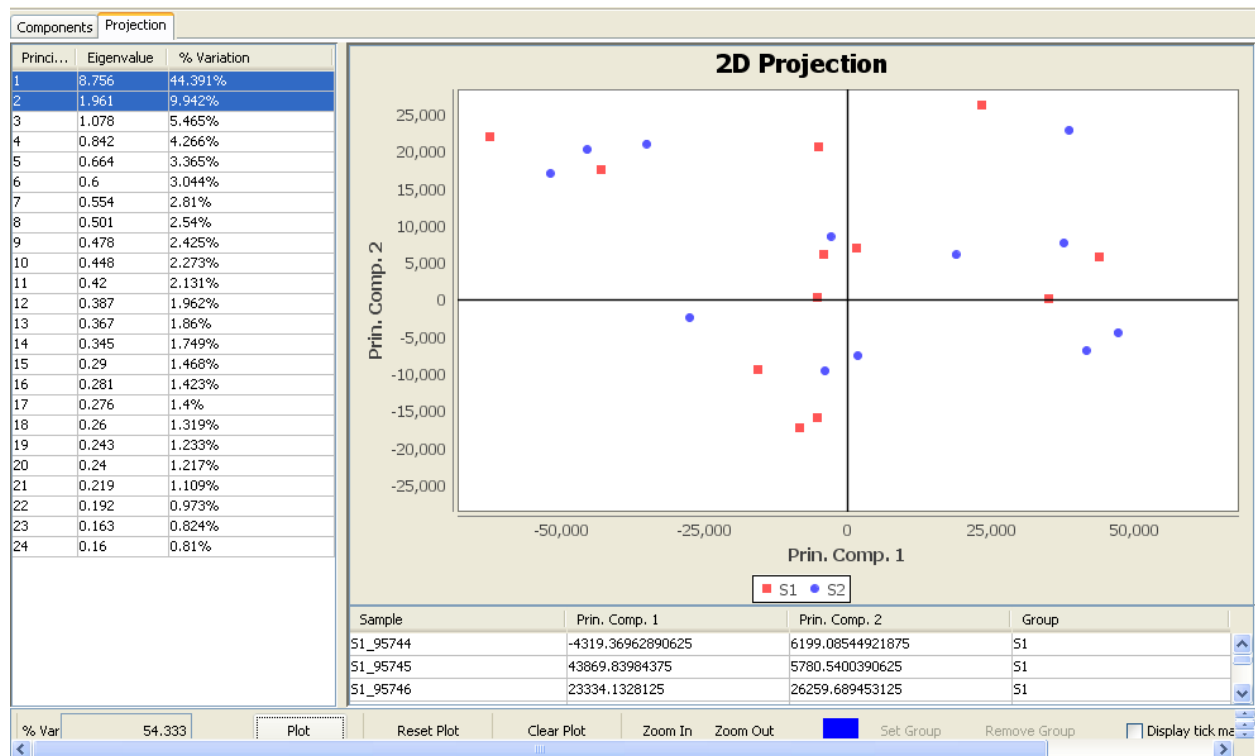
Now you want to see the distribution of the 12 subjects' data, before and after exercise, on these new dimensions. If the before and after data groups occupy completely different sectors in the 2D or 3D view, it indicates that these two groups are significantly different.

10. Click the **Projection** tab:

Component:	Projection	
Princi...	Eigenvalue	% Variation
1	8.756	44.391%
2	1.961	9.942%
3	1.078	5.465%
4	0.842	4.266%

11. Select the first two principal components, then click the **Plot** button at the bottom of the PCA Viewer window.

A two-dimensional projection like the following appears:



Note: You can also generate a three-dimensional projection by selecting three principal components. However, three-dimensional projections require that you have Java 3D installed.

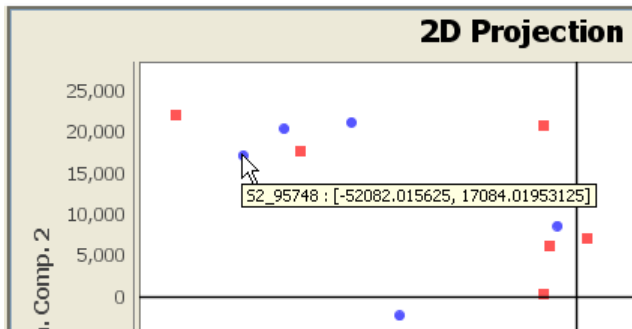
In examining the two-dimensional projection and data, you observe that:

- The two new dimensions of this 2D view represent 44.4% and 9.9% variability of the data. If there is significant change caused by exercise, there should be separation of two groups of subject data (blue and red data points) in this view.
- Even distribution of red data points in this view indicates no outlier among the 12 gene expression profiles before the exercise. Similarly, blue data points indicate no outlier for those after the exercise.
- The evenly mixed red and blue data points indicate there is no significant change in gene expression profiles for these 12 patients before and after exercise.

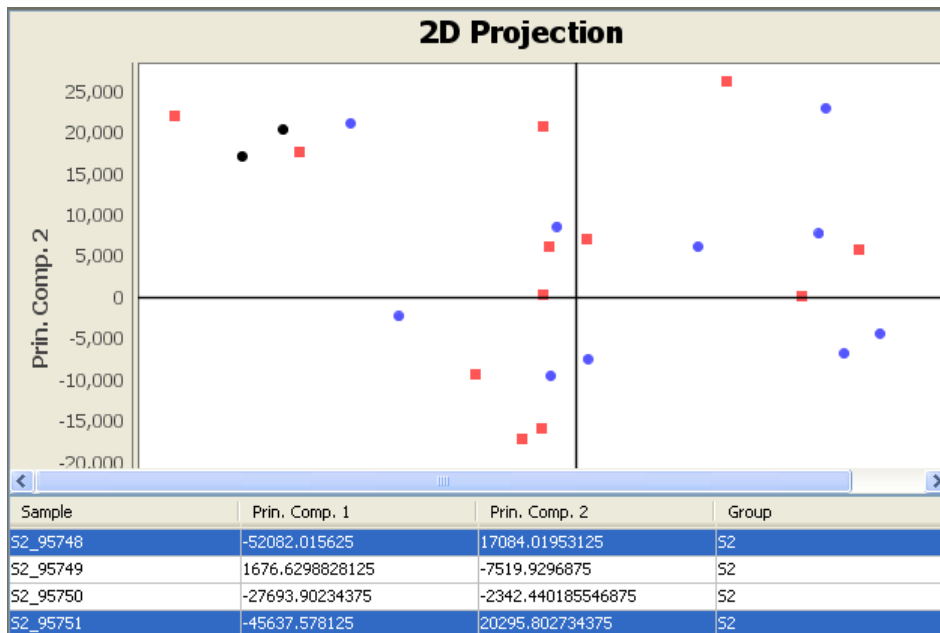
The observations above are based on the capability and limitation of PCA method. Other analysis methods may yield different observations and conclusions.

12. Familiarize yourself with the interactive features of the visualization – for example:

- Hover the mouse pointer over a data point to display its data:



- Click on one or more data points (selected data points become black) to highlight the corresponding sample's data in the table below the visualization:



- Conversely, select one or more samples in the table to locate the corresponding data points in the two-dimensional projection (the data points become black).

13. When finished viewing the projection, close the PCA Viewer window.

