





## Data Filtering and Differential Expression Analyses

For each parameter of interest in the study, the following comparison procedure is used.

- Signal quality: Genes are filtered to eliminate low-quality and less reliable measurements. The filtering is based on raw expression signal (after RMA or MAS5 processing for Affymetrix data). For each comparison (e.g., treated vs. control) the signal intensity must be above a certain threshold level of intensity in at least one of the two groups being compared. The threshold is calculated as the 20th percentile of all raw (probeset-level) measurements.
- Comparisons: Samples are grouped by curators, depending on the parameter of interest. The following guidelines apply:
  - In a time series study, all time points are compared against the zero time point.
  - In studies with one or more treatments, all treatments are compared against the control.
  - If the study involves investigating a disease vs. normal state, the disease samples are compared against normal samples. For studies investigating benign vs. primary vs. metastatic cancer, each cancer state is compared against benign, because benign is the closest to normal.
- Statistical tests: When comparing independent samples, a two-sample t test is performed on each gene, comparing the treated group to the control group (a one-sample t test is applied when there are no control samples). The Welch t test is used by default, when variances cannot be assumed to be equal. When comparing paired samples, a paired t test is performed.

- Gene filtering: Except for tissue atlas biosets where the full gene signature is presented to the user, genes in all other biosets are filtered using a p-value cutoff of 0.05, with no multiple-testing correction. The resulting set of genes is further filtered using a fold-change difference cutoff between the average intensity in test and control groups. Typically, the NextBio platform considers fold changes  $\geq 1.2$  based on microarray technology sensitivity, but the choice may also depend on the type of experiment, tissue source, and other factors. For example, changes in many cell lines, certain treatment types, or brain tissues can be a lot smaller than changes in blood cells of trauma patients. In this case, the cutoffs can be adjusted accordingly. Volcano plots can be used to evaluate fold-change cutoff thresholds using distribution of fold changes within a given comparison. For more “noisy” tissues (e.g., tumor, blood), cutoffs are more stringent relative to other types of samples. While the choice of fold-change cutoff is not an exact science, choices are limited to a small range that depends on the experimental context.

## Constructing a Bioset

The resulting final list of genes is used to construct the bioset, as shown in Figure 5. The bioset consists of the list of genes, along with associated p-value, fold change, and average expression level in each treatment group, for each gene. In addition, a short description of the bioset and a summary of the analysis performed can also be included in the file to be uploaded. All biosets in the NextBio library contain this information, which can be viewed on the bioset information page through the web interface.

For rank-based meta-analyses, the NextBio platform uses standard fields in the column headers to rank features in the bioset. If more than

Figure 5: Gene List Example

	A	B	C	D	E	F	G	H
1	Bioset summary = GSE7025 - PPARgamma ligands and platinum-based drugs in cancer							
2	Comparison = Differentially-expressed genes between samples in A549 cells treated with rosiglitazone vs. non-treated samples							
3	Analysis summary = CEL files were normalized using RMA (R-package). Statistical t-test assuming variances are not equal (Welsh t-test)							
4								
5								
6	Gene name	fold change	p-value	Custom field 1	Custom field 2	...		
7	200606_at	1.21	0.0107	259	213.4			
8	200762_at	1.29	0.0216	670.4	519			
9	200778_s_at	1.33	0.0291	1445.8	1084.8			
10	200789_at	1.53	0.00199	1145.7	749.5			
11	200798_x_at	1.25	0.0389	468.6	375.1			
12	200817_x_at	-1.24	0.0107	8177.9	10155.6			
13	200862_at	1.53	0.0236	329	214.6			
14	200878_at	1.48	0.0313	1112.9	753.2			
15	200903_s_at	-1.31	0.0474	1280.9	1677.6			
16	200906_s_at	-1.47	0.0178	386.9	570.1			
17	201015_s_at	-1.26	0.0424	214.2	270.6			
18	201029_s_at	-1.23	0.0495	1539.2	1894.5			
19	201060_x_at	1.39	0.0421	307.9	220.9			
20	201108_s_at	1.45	0.024	140.3	96.9			
21	201162_at	-1.22	0.0179	326.4	397.8			
22	201233_at	-1.21	0.0461	224.4	272.6			
23	201255_x_at	-1.24	0.031	393.4	486.5			
24	201369_s_at	-1.28	0.0219	60.1	77			
25	201398_s_at	-1.32	0.0186	1350.8	1785.7			
26	201432_at	1.57	0.00123	1263.1	806.4			
27	201461_s_at	1.38	0.0378	54.2	39.2			
28	201471_s_at	1.44	0.0267	5368.3	3740.8			
29	201480_s_at	-1.26	0.0369	425.6	536			
30	201578_at	1.75	0.0309	699.8	401			

Sample bioset from drug vs. control microarray samples analysis from A549 cell lines treated with rosiglitazone. In this example, the “fold change” column is selected for ranking during bioset import into the NextBio platform. Supplemental information columns can be included as custom columns with user-defined titles (currently, a maximum of five columns) to provide average expression level information in treated and control groups for interested users.

