

RPPA Analysis: Stuart Gallagher

Lucy Liu

Contents

Introduction	2
Experiment design	2
Batches	3
Data exploration	8
RFI for each protein	8
RFI for each sample	10
MDS plot of samples	12
Heatmap	13
Normalisation	17
Normalisation methods	17
Housekeeping protein	17
Median normalisation	17
Loading control	17
Invariable protein normalisation	18
How to evaluate normalisation methodologies	18
MA plot	18
RLE	18
Pairwise correlation coefficients (Spearman ρ)	18
RLE plots	19
Raw data	19
Median normalisation	19
Loading control normalisation	19
Global rank invariant	21
MA plots	22
Raw data	22
Median normalisation	27
Loading control normalisation	31
Global rank invariant	35
Pairwise correlation coefficient plots	39
Raw data	39
Median normalisation	39
Loading control normalisation	40
Global rank invariant	40
Normalised data	41
Average replicates then normalise	41
Normalise both replicates	41
References	42

Introduction

The aim of this analysis is to investigate the protein expression of ~52 melanoma cell lines. The data from this experiment will be used to explore a number of research questions:

1. A portion of cell lines (11 cell lines) have been treated with IBET151. The researcher would like to investigate how protein expression differs between treatment with IBET151 vs DMSO.
(Note: of the 11 cell lines treated with IBET151, there are only 3 that were also treated with DMSO - C106M, A06M & D14M2. All cell lines treated with DMSO were at time point 48h while the IBET151 treated cell lines were at 24h)
2. Is sensitivity of specific cell lines to a certain drug correlated to protein expression?
3. Is the mutation status/transcriptional signature of the melanoma cell lines correlated to its protein expression profile?

Experiment design

There were two biological replicates in this experiment. The samples are detailed below:

Lysate.ID	Cell.Line/Tissue.type	Treatment	Time.point
SJG-1	C078M	DMSO	48
SJG-2	D17M1	DMSO	48
SJG-3	D11M2	DMSO	48
SJG-4	D04M	DMSO	48
SJG-5	D23M	DMSO	48
SJG-6	A15M2	DMSO	48
SJG-7	D20M	DMSO	48
SJG-8	D35M1	DMSO	48
SJG-9	D38M2	DMSO	48
SJG-10	A11M	DMSO	48
SJG-11	A06M	DMSO	48
SJG-12	C044M	DMSO	48
SJG-13	C054M	DMSO	48
SJG-14	C055M	DMSO	48
SJG-15	C058M	DMSO	48
SJG-16	C062M2	DMSO	48
SJG-17	C074M	DMSO	48
SJG-18	C089M	DMSO	48
SJG-19	C027M	DMSO	48
SJG-20	C037M	DMSO	48
SJG-21	C057M	DMSO	48
SJG-22	C094M	DMSO	48
SJG-23	C013M	DMSO	48
SJG-24	C084M	DMSO	48
SJG-25	C086M	DMSO	48
SJG-26	C091M	DMSO	48
SJG-27	C100M	DMSO	48
SJG-28	C001M	DMSO	48
SJG-29	C002M	DMSO	48
SJG-30	C011M	DMSO	48
SJG-31	C017M	DMSO	48
SJG-32	C106M	DMSO	48
SJG-33	C006M	DMSO	48
SJG-34	C021M	DMSO	48

Lysate.ID	Cell.Line/Tissue.type	Treatment	Time.point
SJG-35	C065M	DMSO	48
SJG-36	C076M	DMSO	48
SJG-37	D22M1	DMSO	48
SJG-38	C092M	DMSO	48
SJG-39	A04M	DMSO	48
SJG-40	D36M	DMSO	48
SJG-41	D08M	DMSO	48
SJG-42	A02M1	DMSO	48
SJG-43	D05M1	DMSO	48
SJG-44	D14M2	DMSO	48
SJG-45	D41M	DMSO	48
SJG-46	C060M1	DMSO	48
SJG-47	C088M1	DMSO	48
SJG-48	D24M	DMSO	48
SJG-49	C077M1	DMSO	48
SJG-50	C016M1	DMSO	48
SJG-51	C022M1	DMSO	48
SJG-52	C025M1	DMSO	48
SJG-53	HDF	DMSO	48
SJG-54	HEMn-MP	DMSO	48
SJG-55	HEMn-LP	DMSO	48
SJG-56	CO57M	IBET151	24
SJG-57	CO02M	IBET151	24
SJG-58	C106M	IBET151	24
SJG-59	CO37M	IBET151	24
SJG-60	A06M	IBET151	24
SJG-61	CO92M	IBET151	24
SJG-62	D14M2	IBET151	24
SJG-63	CO65M	IBET151	24
SJG-64	CO25M	IBET151	24
SJG-65	CO44M	IBET151	24
SJG-66	CO76M	IBET151	24
SJG-67	Mel-XY3	DMSO	48
SJG-68	Mel-XY3_Plx-7day	DMSO	168hr
SJG-69	Mel-XY3_Sur	DMSO	48
SJG-70	Mel-XY3_Sur_post	DMSO	48

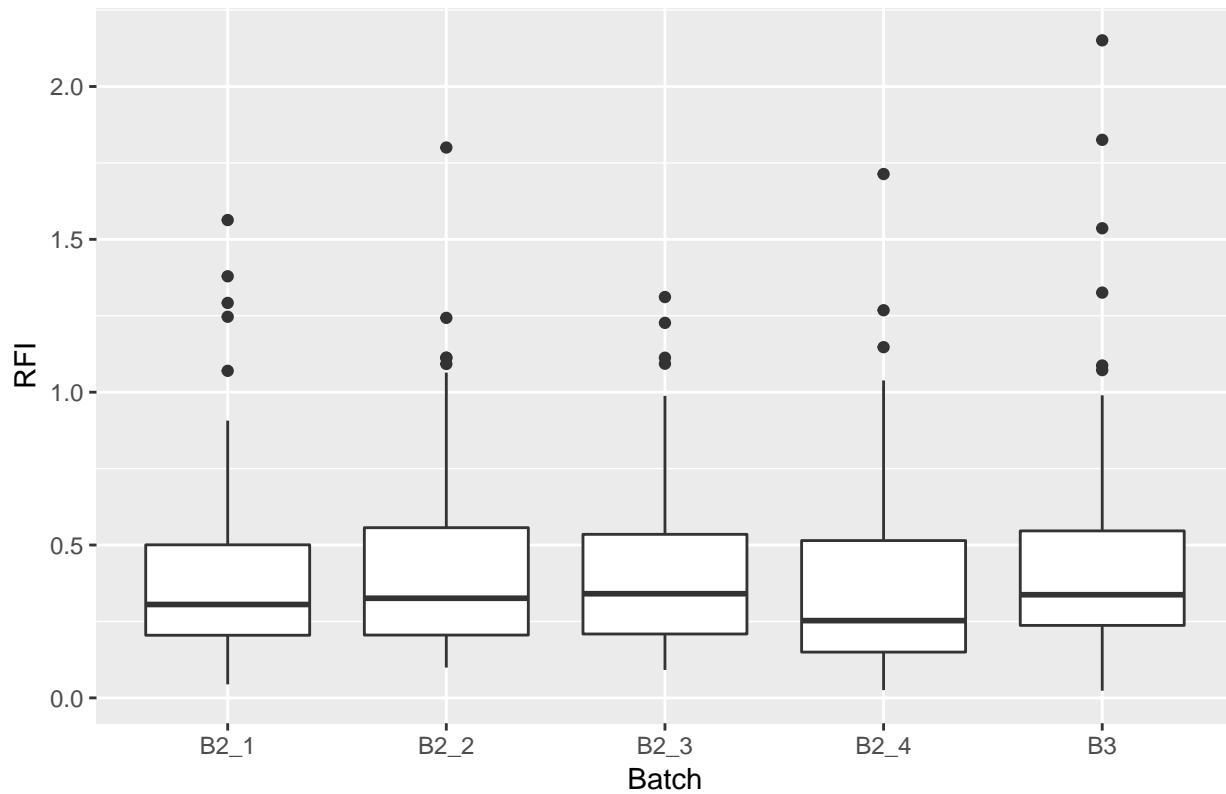
Batches

There were 2 main batches in the RPPA experiment. Note, in the pilot batch ('Batch 1'), 4 ABs were tested but were all repeated in batch 2. The 2 batches were thus:

- Batch 2: Full run of all ABs
- Batch 3: Samples and AB's that did not perform well in batch 2 are repeated in batch 3

First, we will compare the RFI values of the positive control AB (Prohibitin) in all batches. Note that Batch 2 was very large and Prohibitin was run 4 times, twice in the first run of ABs and twice in the second run of ABs.

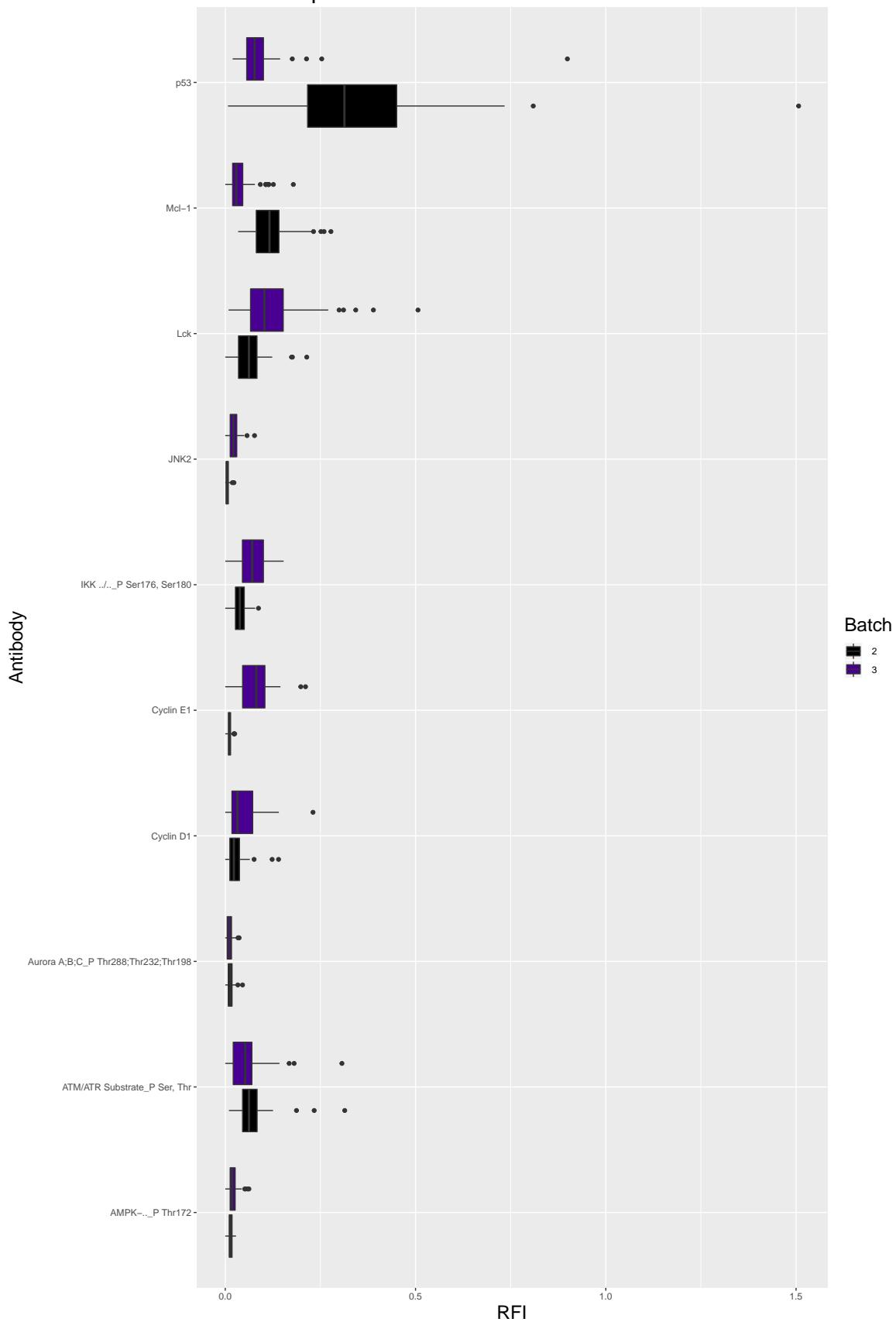
Prohibitin expression in each batch



The prohibitin expression across all batches are quite similar, which suggests that there were no significant batch effects.

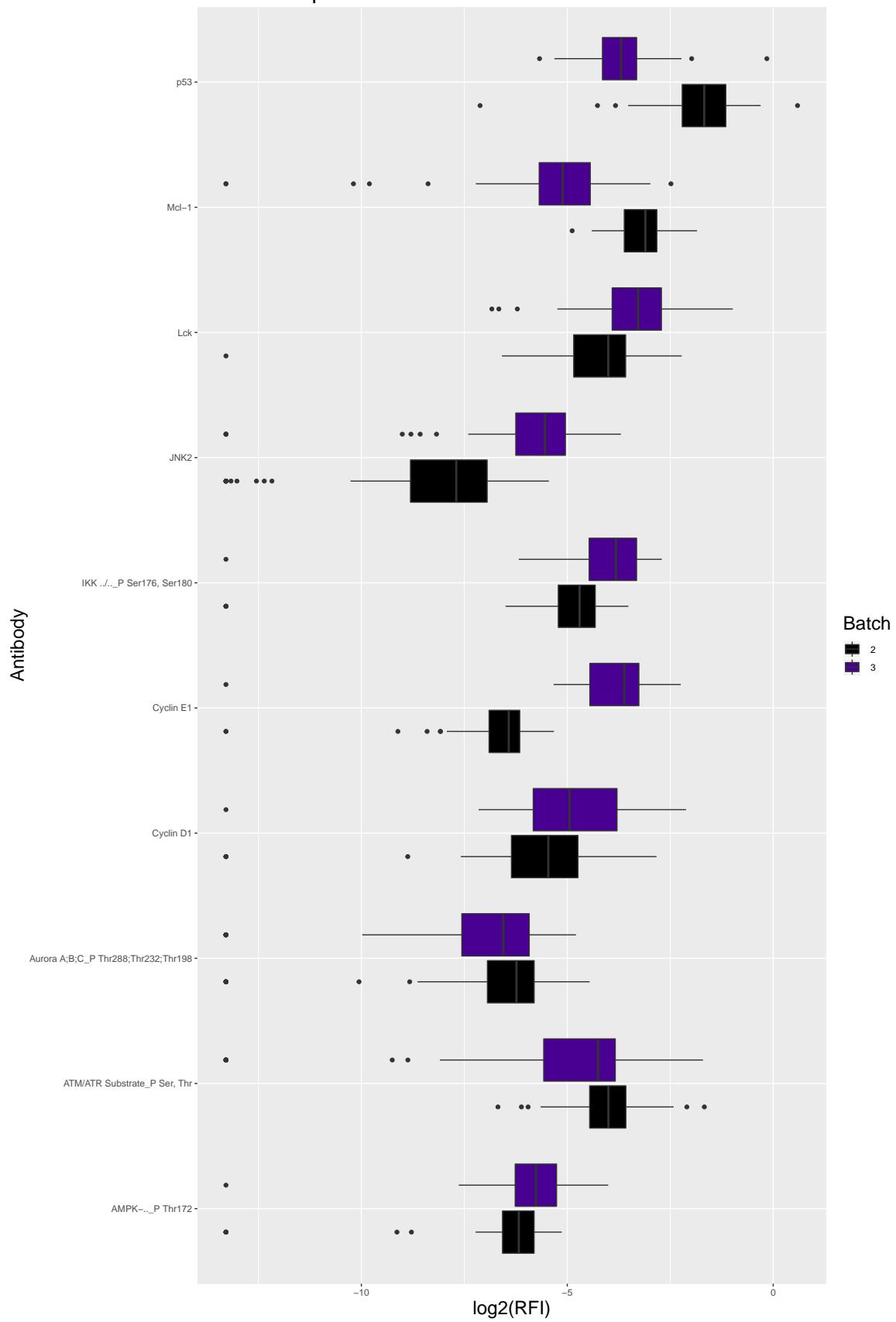
Most ABs were tested solely within batch 2. This makes it difficult to investigate batch effects as we do not know if the difference in RFI values between batches is due to between array variation or differences in protein expression. There were 10 ABs tested on both batch 2 and 3 (e.g. run in batch 2 for replicate 1 samples (i.e. SJG1-SJG70) and in batch 2 for replicate 2 samples (i.e. SJG71-SJG140)). We compare the RFI values of these ABs between the two batches.

Comparison of RFI between batches for each AB



Log RFI values and replot to better visualise ABs with lower RFI values:

Comparison of RFI between batches for each AB

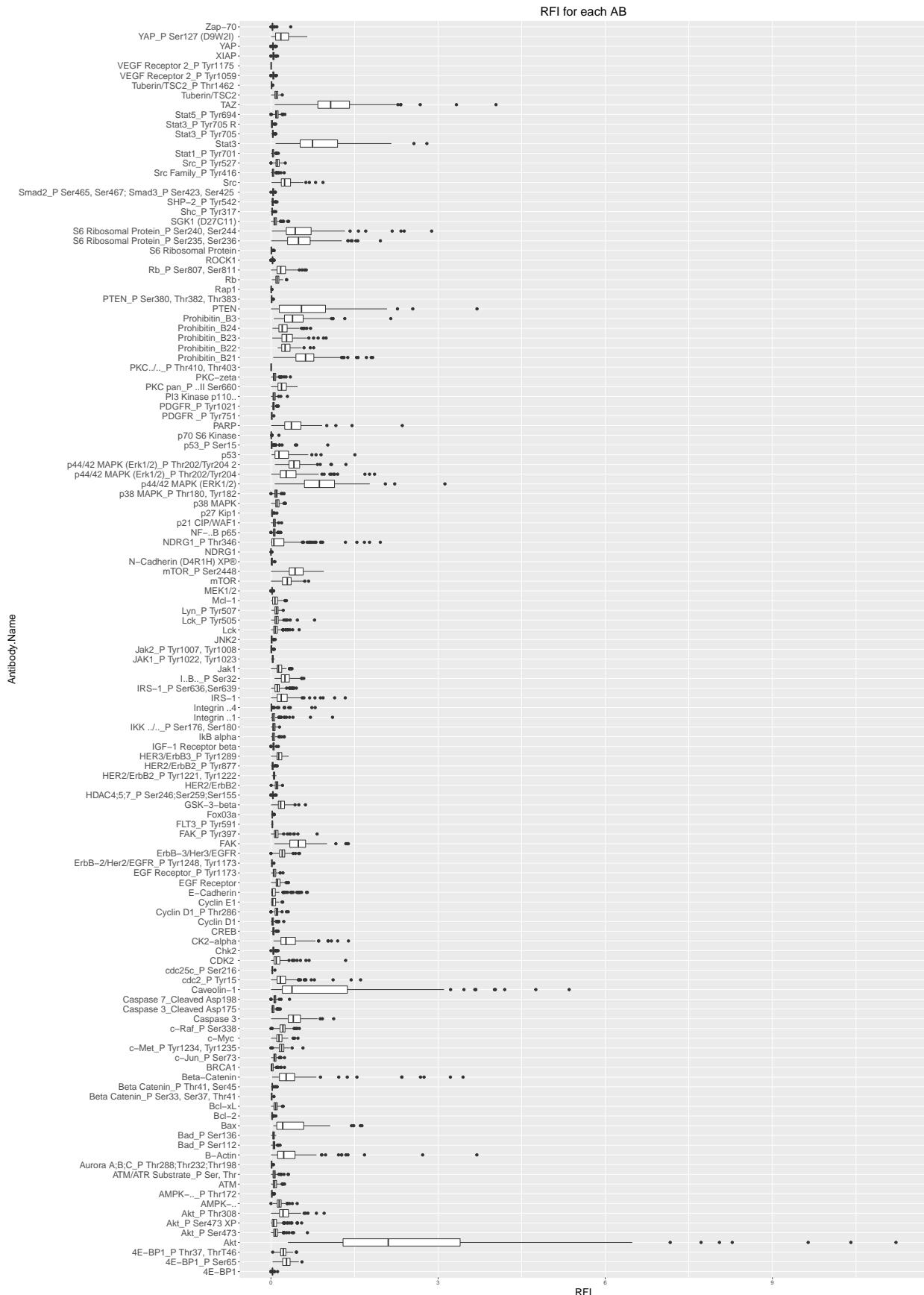


Overall, the concordance between batches is not great. However, as mentioned above, it would be difficult to correct for batch effects as the majority of ABs were tested solely in one batch only (thus we do not know if the difference in RFI values between batches is due to between array variation or differences in protein expression). Further, as the vast majority of ABs were tested in batch 2 and the data will be normalised overall for proteins effects, batch correction will not be further considered.

Data exploration

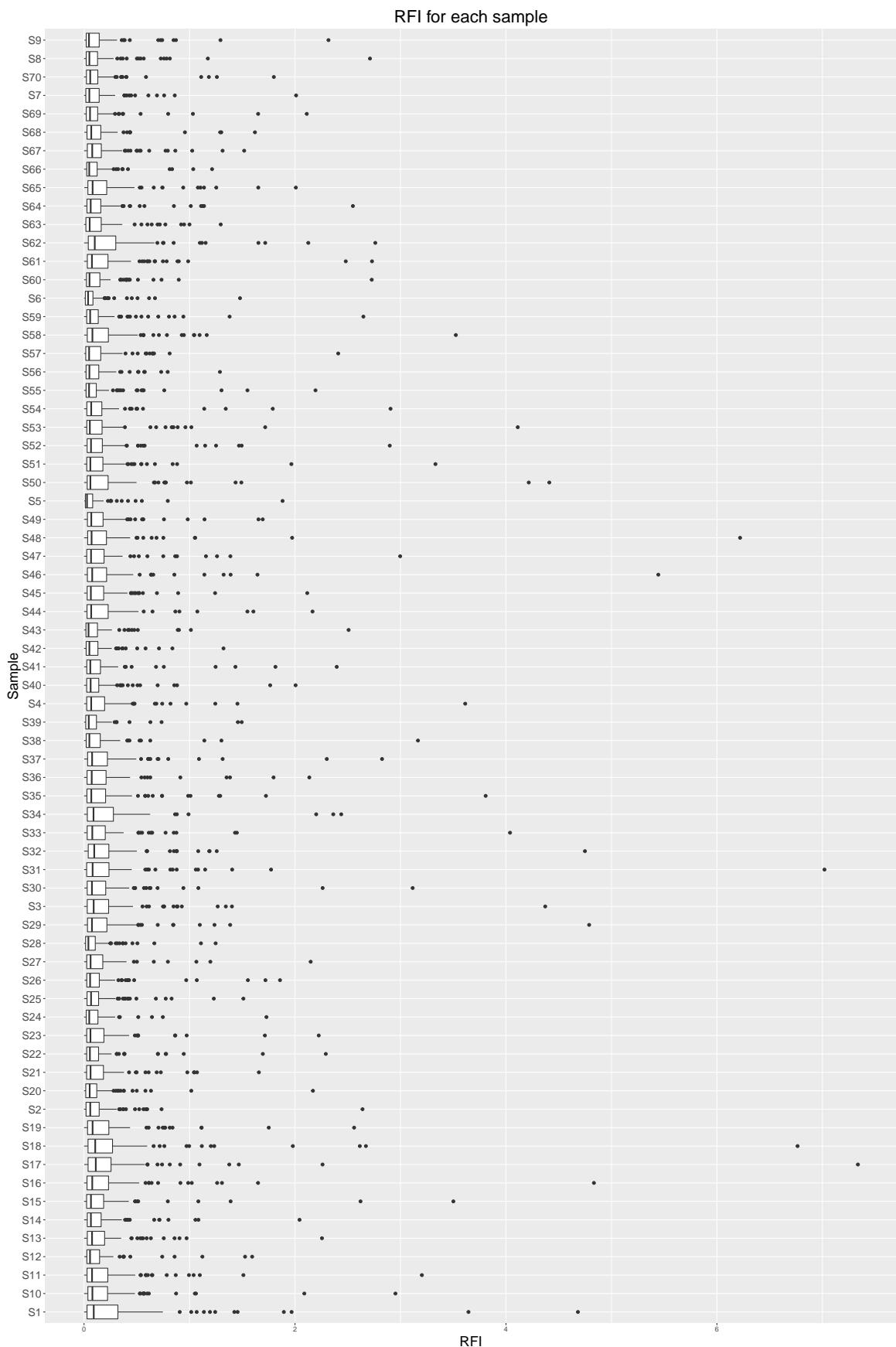
RFI for each protein

A boxplot of the raw RFI values of all samples (both replicates) are plotted for each AB.



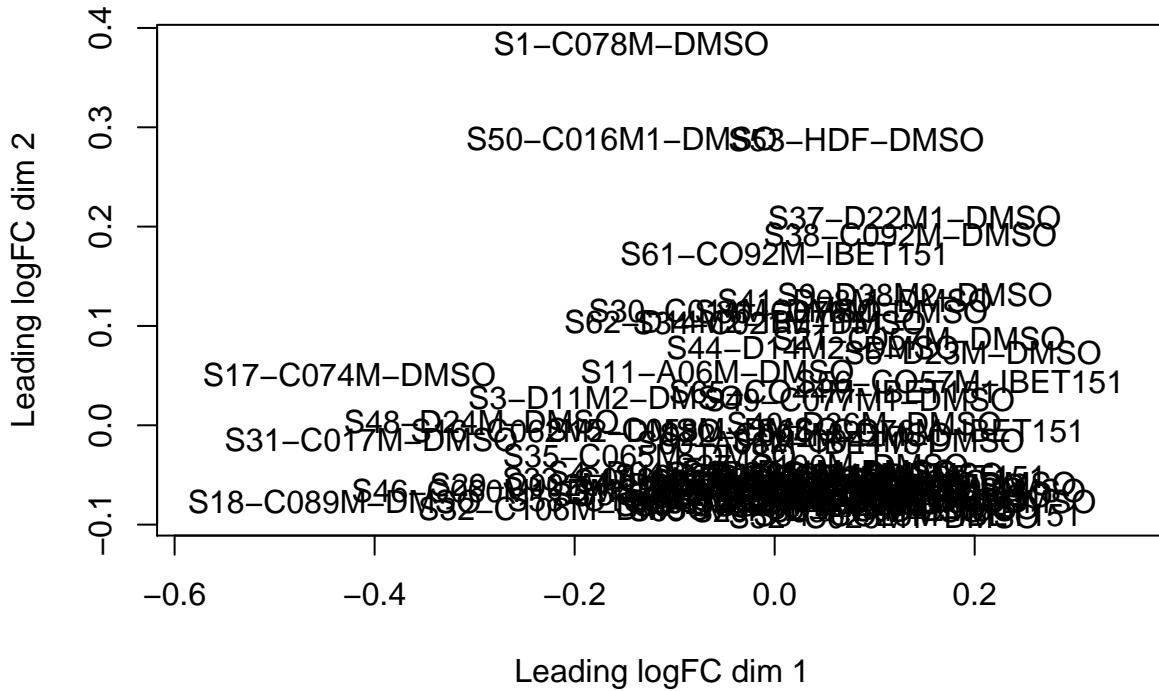
RFI for each sample

Average RFI from the two replicates are plotted for each sample.

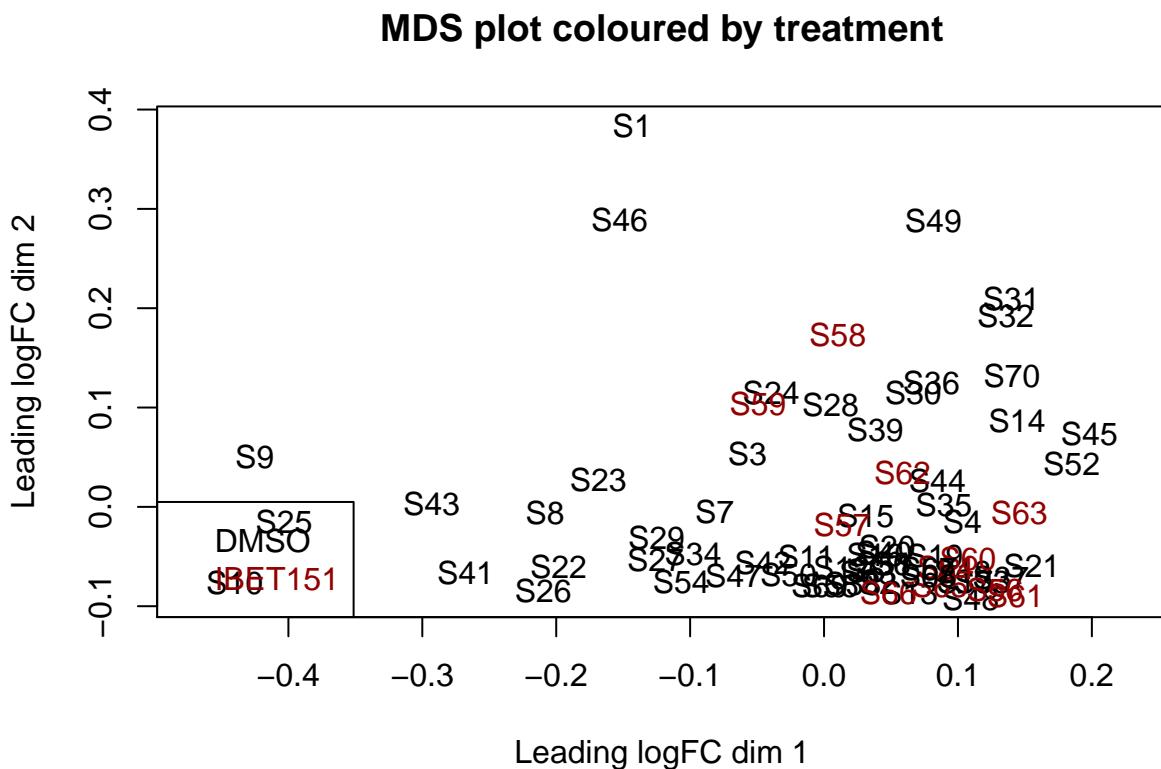


MDS plot of samples

A multidimensional scaling (MDS) plot graphically represents relationships between objects. The distance between two samples approximates their similarity or dissimilarity. The MDS plot is useful for identifying if there are any distinct groups within our samples.



Re-plot MDS using sample names S1-S70 so labels are more legible, and colour sample names by treatment (DMSO or IBET151).



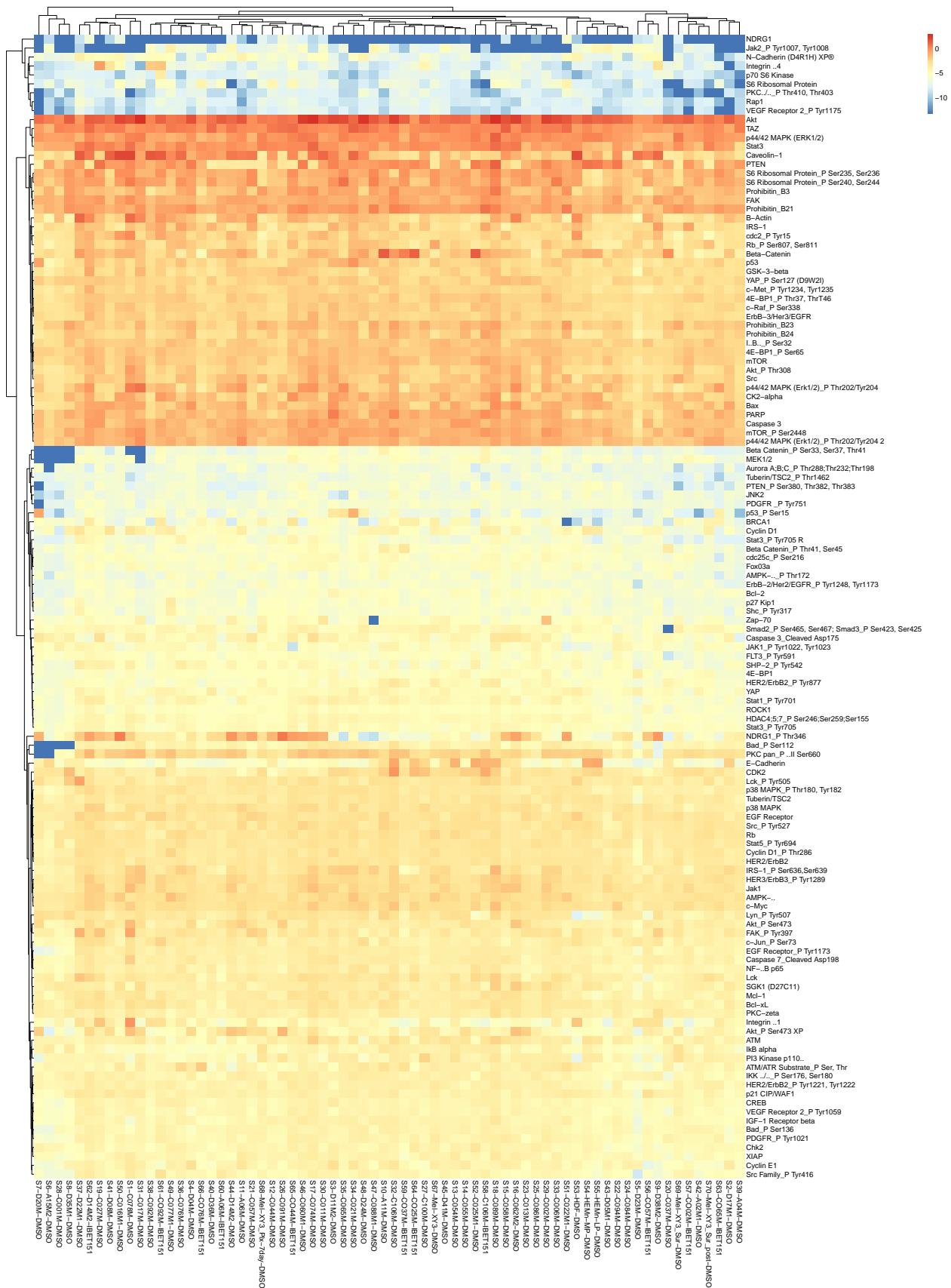
There does not appear to be distinct groups of cell lines within the experiment. The DMSO and IBET151 treated cell lines also do not appear to separate. This is not unexpected as most of the cell lines in the experiment are different (there are only three cell lines that feature twice in the experiment). This may be the reason that the IBET151 treated cell lines do not cluster separately from the DMSO treated cell lines.

Heatmap

Heatmap and dendrogram of raw RFI values.



To aid in better discriminating the range of RFI values, a heatmap of log 2 RFI values is plotted as well.



Normalisation

There are three primary factors that need to be normalised for in an RPPA experiment (Liu et al. 2014, Wachter et al. (2015)):

1. Spatial bias: Differences in intensity caused by location of the lysate spot on the slide (e.g. rim effects). The Zepto system already accounts for this using the BSA control spots, thus spatial normalisation will not be further considered.
2. Total amount of protein (loading) of different samples on the slides: Although the total proteins in the lysate are gauged before they are printed, this is confounded by lipids and other biological materials in the samples. Thus, the total protein measurement is only a rough estimate.
3. Non linearity of variances: A MA (differences versus means) plot of the differences (M) in intensity of ABs between two samples varies across the spectrum of A , the mean intensity of the two samples. This is especially the case at the upper range and sometimes the lower range of A . This can be seen best when comparing a pair of technical replicates. As there would not be any differential expression, you would expect their intensities to be equal. Some imbalance often occurs and this imbalance may vary depending on the average intensity (A). Variation in the imbalance When normalised well all the points of a MA plot should align with a horizontal line and be evenly distributed about this line.

Normalisation methods

Housekeeping protein

This normalisation methodology is based on the assumption that the levels of housekeeping (HK) proteins such as β -Actin is uniform across samples and experiment conditions. Thus any differences in level of these housekeeping proteins is due to differing amounts of protein loading. Our previous experience suggests that HK proteins that are effective in western blots are not in RPPA due to the increased sensitivity of RPPA.

Median normalisation

This method assumes that all measured proteins reflect the total protein amount of one sample. Thus the median AB of a sample estimates sample loading. The median value of all AB signals for a sample is used to normalise the raw intensity values for each AB. This is one of the ‘simplest’ methods (involving the least steps) however is biased when the number of ABs is <100 (Liu et al. 2014).

Loading control

This approach is utilised by MD Anderson.

Protein effects on intensity are accounted for by dividing all the raw linear intensity by the median for each AB (across all samples). This is the ‘median centered ratio’. Sample effects are accounted for by taking the median of the median centered ratios for each sample (across all ABs). This becomes the correction factor for that sample. Raw intensity values are divided by this correction factor to obtain the normalised intensity.

The steps are outlined below:

1. Determine median RFI for each AB (across all samples)
2. Divide each RFI by the median within each AB to get the ‘median-centred ratio’.
3. Calculate the median median-centred ratio for each sample (across all ABs). This is the correction factor for each sample.
4. Divide each median-centred ratio by the correction factor for each sample.

Invariable protein normalisation

This methodology was proposed by Liu et al. (Liu et al. 2014) and was originally suggested to normalise microarray data (Pelz et al. 2008).

The aim of this methodology is to determine the most uniformly expressed proteins and use these as an effective ‘housekeeping’ protein to normalise to.

The steps are outlined below:

1. Rank the intensity of ABs for each sample so you have ranked ABs for each sample from the highest expressing AB to lowest expressing AB.
2. Calculate the variance of the ranks for each AB. Remove the AB with the highest rank variance.
3. Re-rank the intensity of ABs.
4. Repeat the steps 2 and 3 until the number of remaining markers reaches a predetermined number (100 kept in Liu et al. study).
5. Trim intensities for each AB e.g. the highest and lowest 25% of values are removed from the data set.
6. Average the remaining values of every AB across all samples and use this as a virtual reference sample.
7. Normalise each sample to the virtual reference sample by lowess smoothing using a MA plot. The normalised values are generated using the residuals of the fit.

How to evaluate normalisation methodologies

To be able to evaluate each normalisation methodology, one must be able to compare how effectively each method normalises for total protein loading and non-linearity of variances. Normalisation of non-linearity of variances is effectively assessed using the MA plot. Normalisation of total protein loading is more difficult to assess and must be done indirectly. Two techniques have been suggested and are discussed below.

MA plot

Normalisation of non-linearity of variances is effectively assessed using a MA plot. A MA plot of replicate samples show the difference in RFI values between the two samples against the average RFI value of the two samples, for each AB.

RLE

RLE (Relative log expression) plots are effective for the assessment for unwanted variation (Gandolfo and Speed 2018). It plots log expression relative to the median of that AB for each sample. Assuming the expression levels for the majority of proteins are stable across cell lines, the boxplots in a RLE plot should be roughly centered on 0 and would be roughly the same size (height).

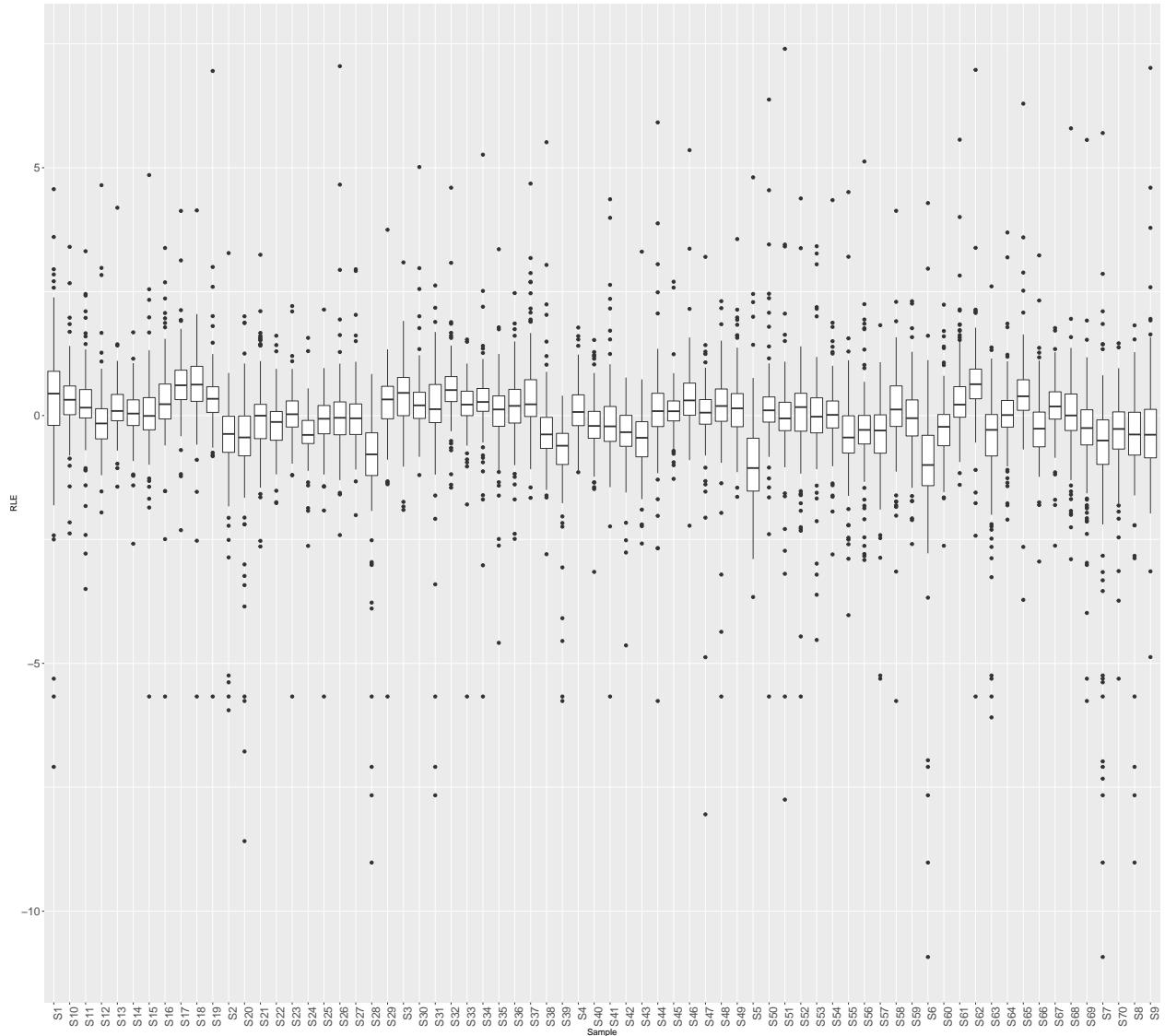
Pairwise correlation coefficients (Spearman ρ)

This plot is used to evaluate the between protein correlations. For each pair of proteins (e.g. for 10 proteins, there would be 45 different possible pairs) the Spearman ρ is calculated across all samples. This is then sorted and plotted. We expect the positive and negative correlation coefficients to be largely equal. If there is an abnormally large number of protein pairs with positive correlation, it may be due to a sample loading effect where some proteins are high in all samples, thus are ranked similarly highly for those samples. Samples with low loading result in low levels of all proteins, and are ranked similarly low for those samples. There would thus be high correlation between a larger proportion of proteins. We expect protein expression to be inherently different due to expression, phosphorylation levels and/or AB affinity.

RLE plots

Raw data

Plot an RLE graph of the raw data:



There is variation in boxplot position and height and shows taht the data needs to be normalised.

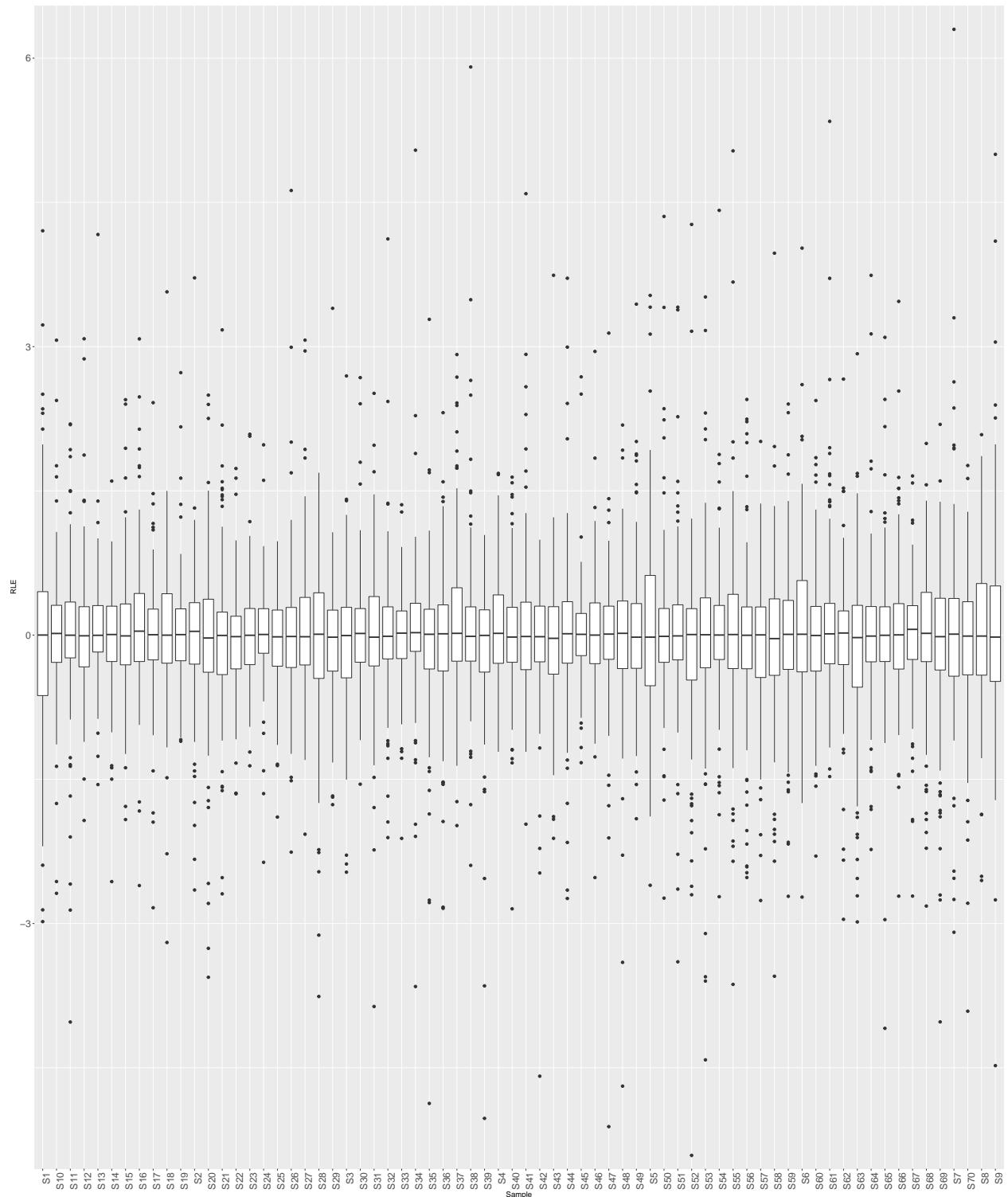
Median normalisation

Median normalisation would not change the above plot as the median value of the median normalised expression values (for each AB) is just 1.

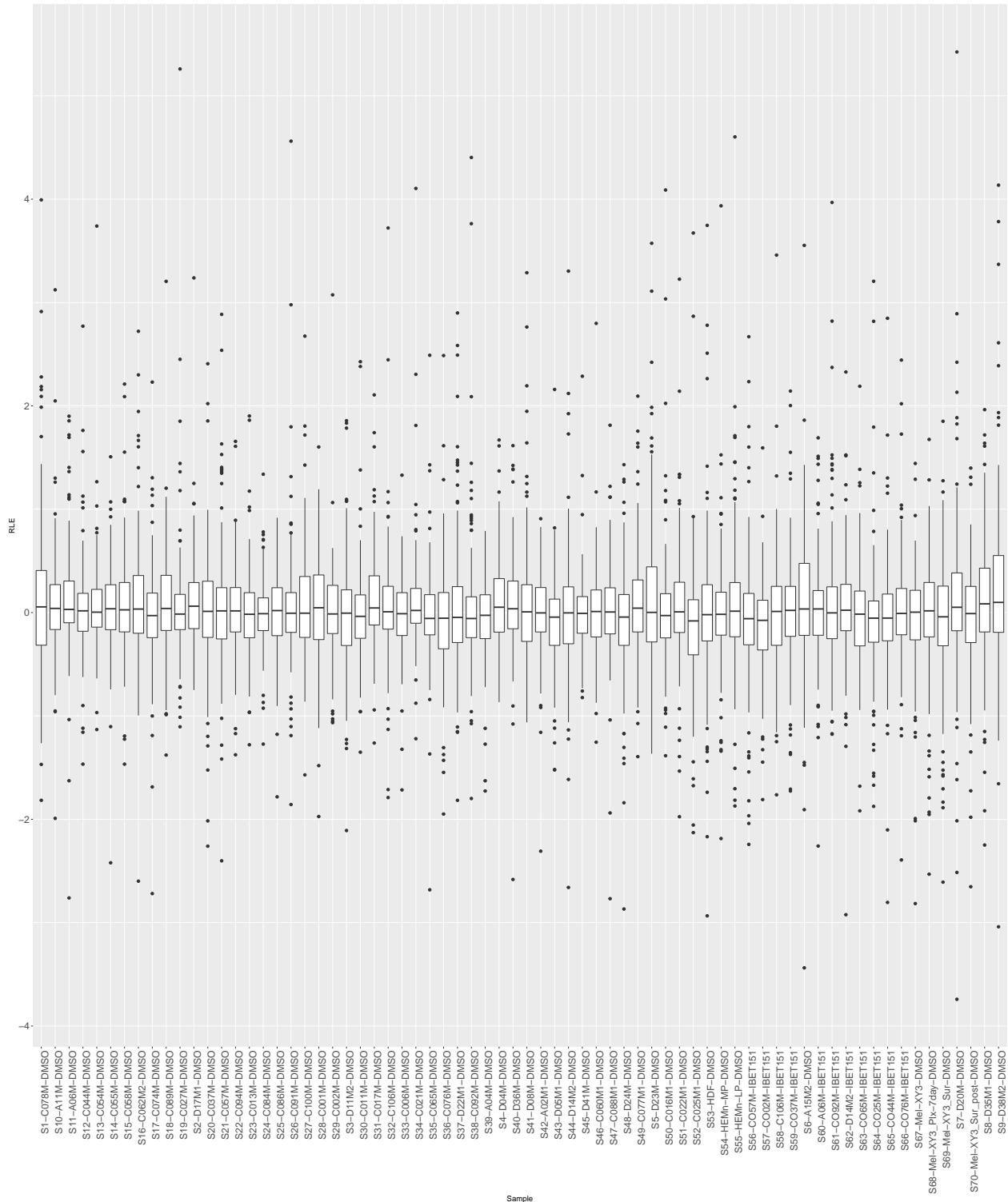
Loading control normalisation

Loading control normalisation causes the median value of each sample to be 0, by virtue of the normalisation approach itself. Nonetheless, the height of each boxplot can be used to assess quality of normalisation. Note

that for many samples there are 1 or 2 very low RLE values which are excluded from the graph to allow easier comparison.



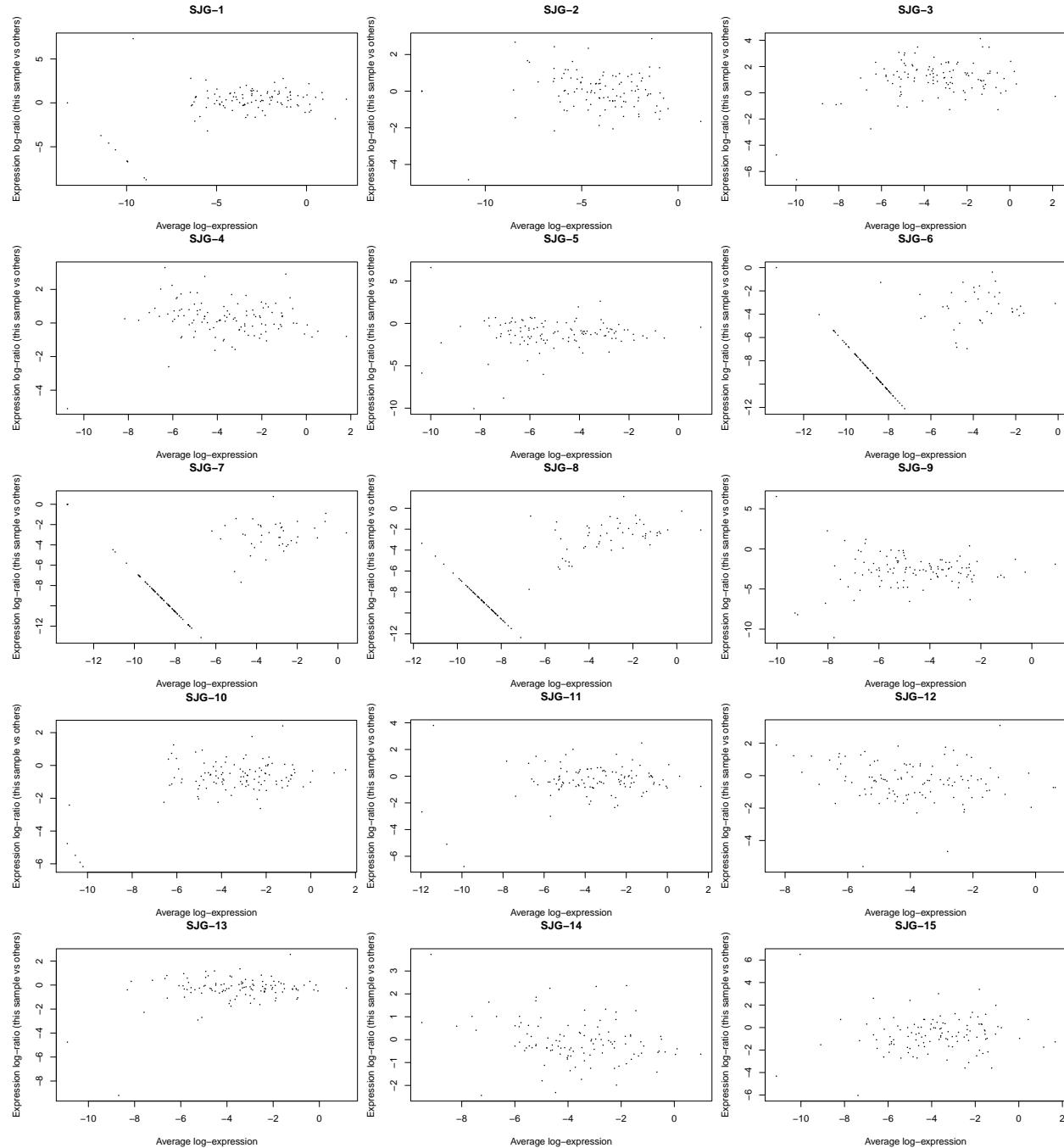
Global rank invariant

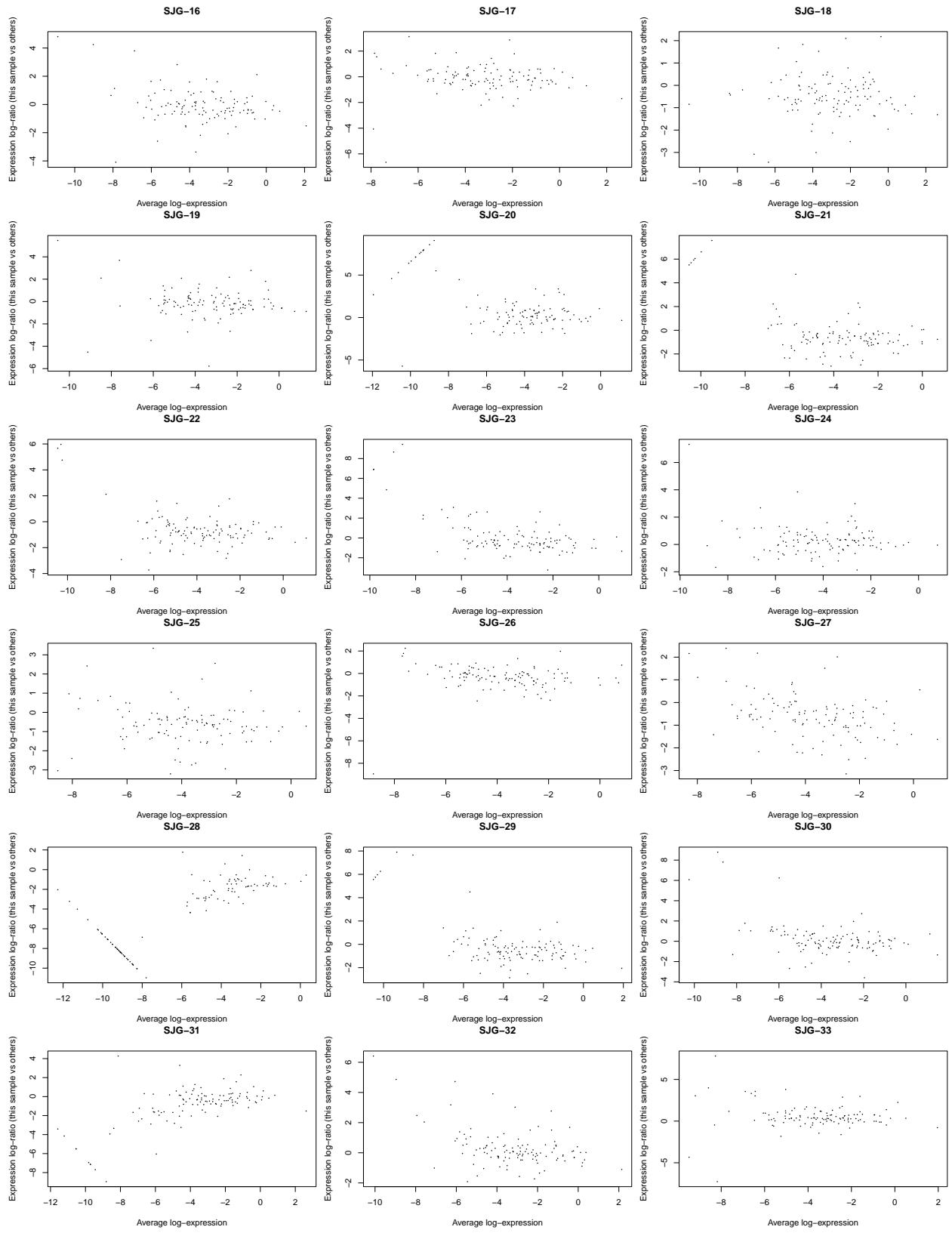


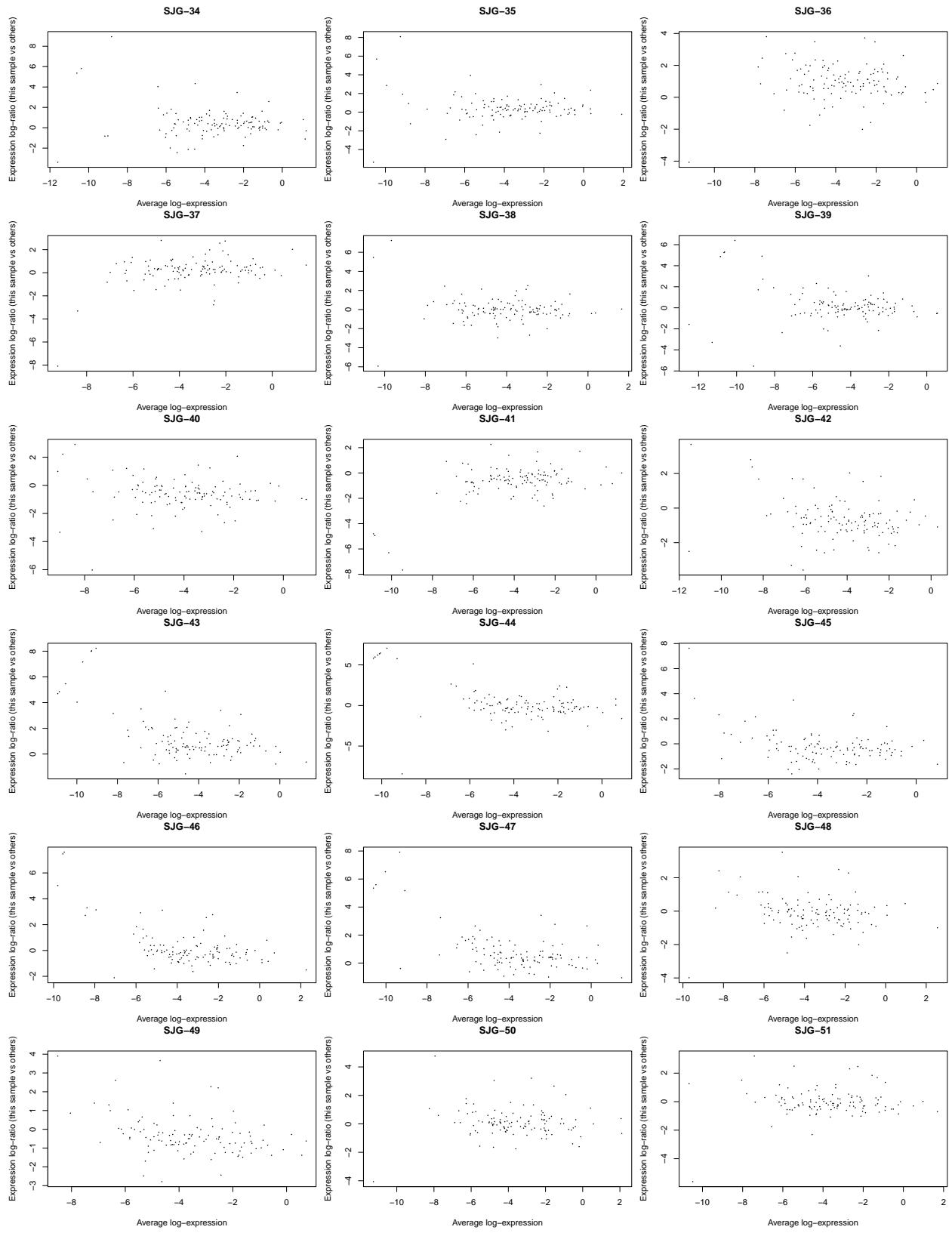
MA plots

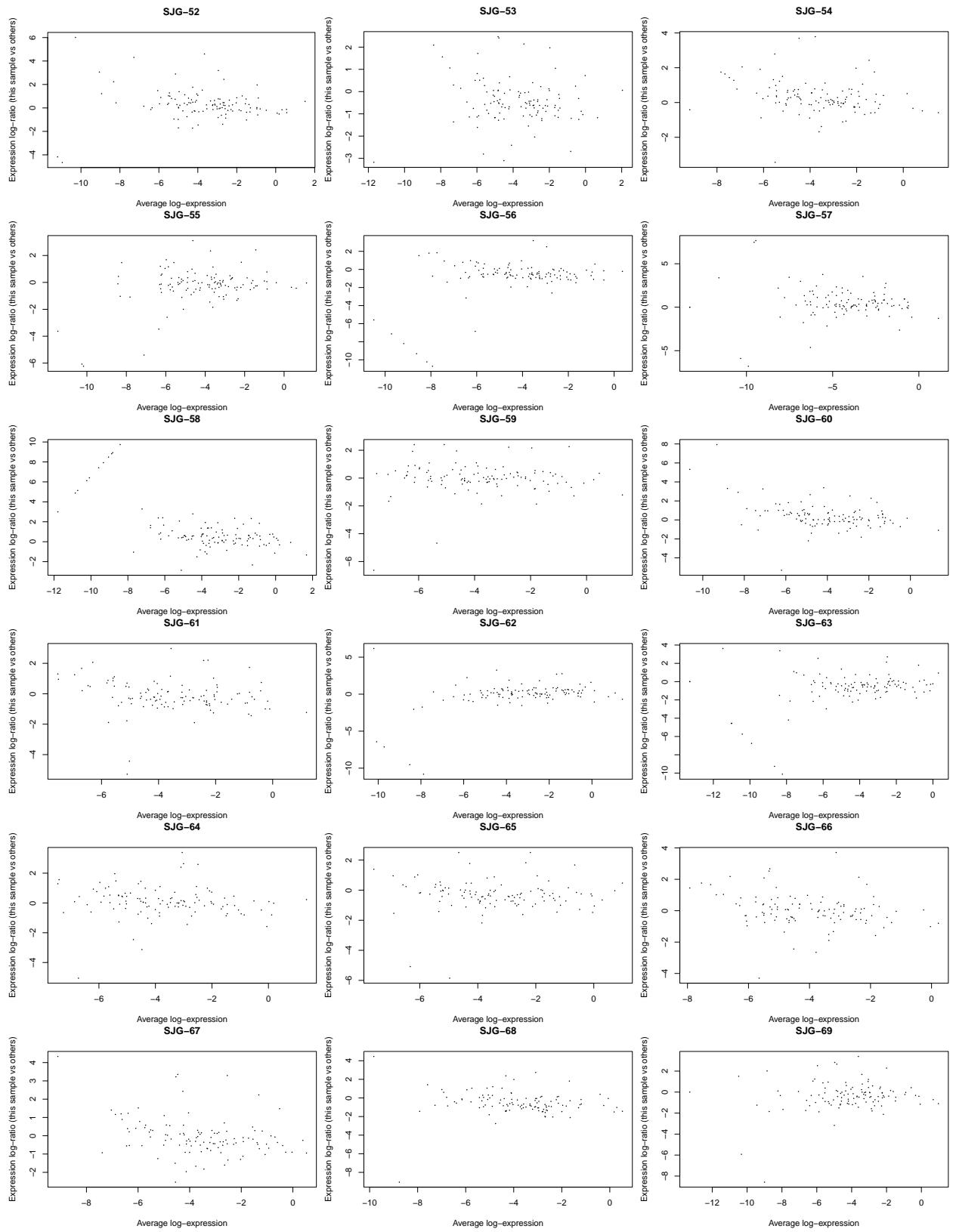
Each point in an MA plot represents an AB. On the y axis is the difference in RFI between two replicates and on the x axis is the mean RFI of the two replicates. For well normalised data, you expect the points to lie along 0 and the variation of points around 0 to be similar along A.

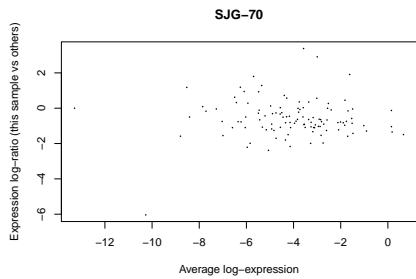
Raw data





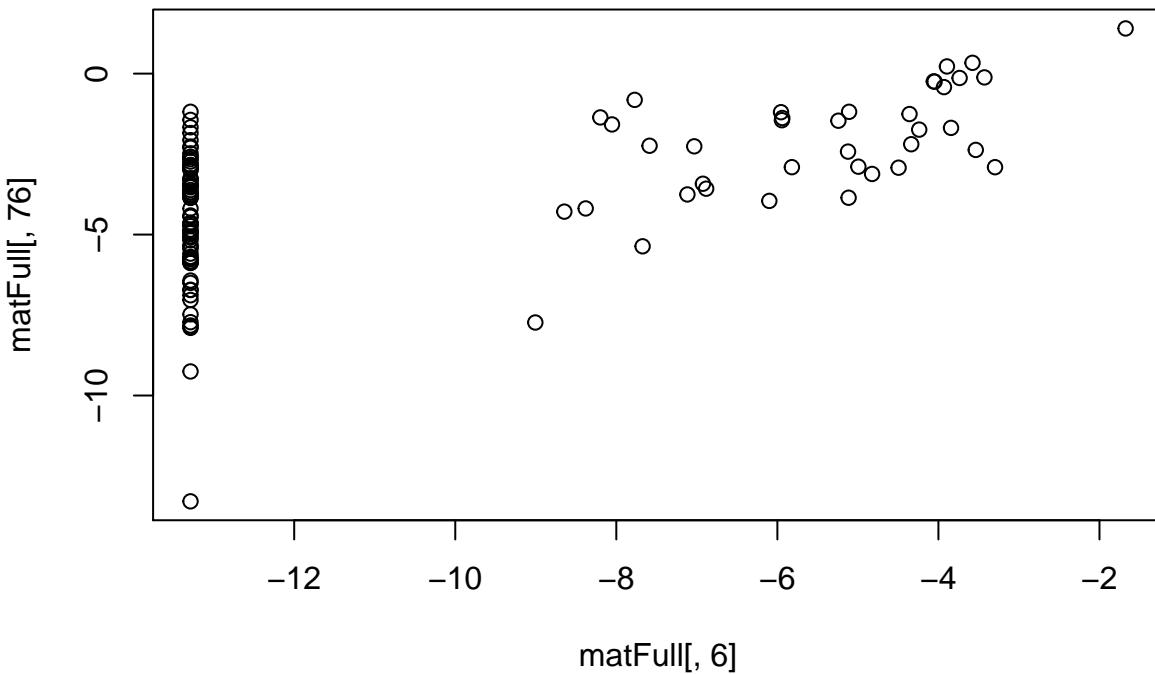






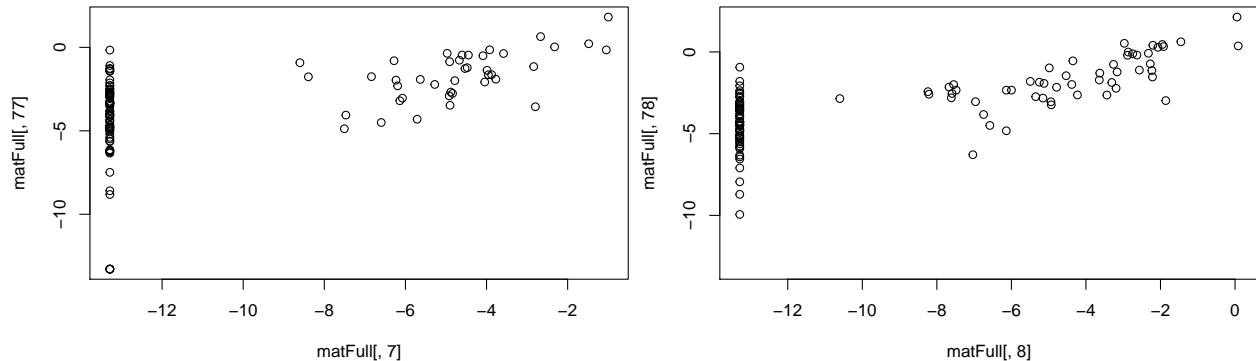
A number of sample pairs have unusual MA plot patterns and are further investigated here:

The RFI values for sample 6 and 76 are plotted:



Can see from the above plot that sample 6 has a large number of uniformly low AB expression values.

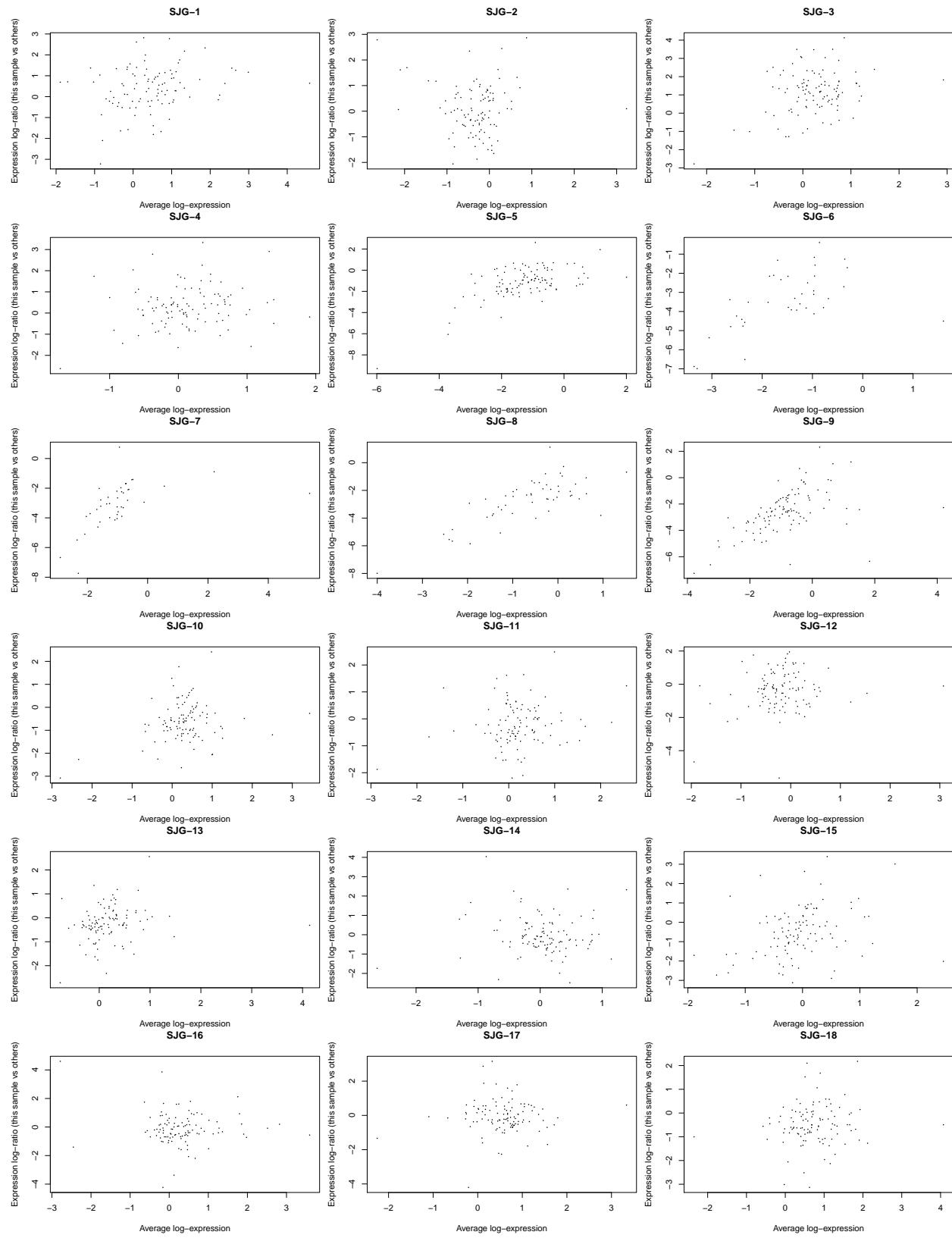
A similar pattern is seen for samples 7 & 77 and 8 & 78:

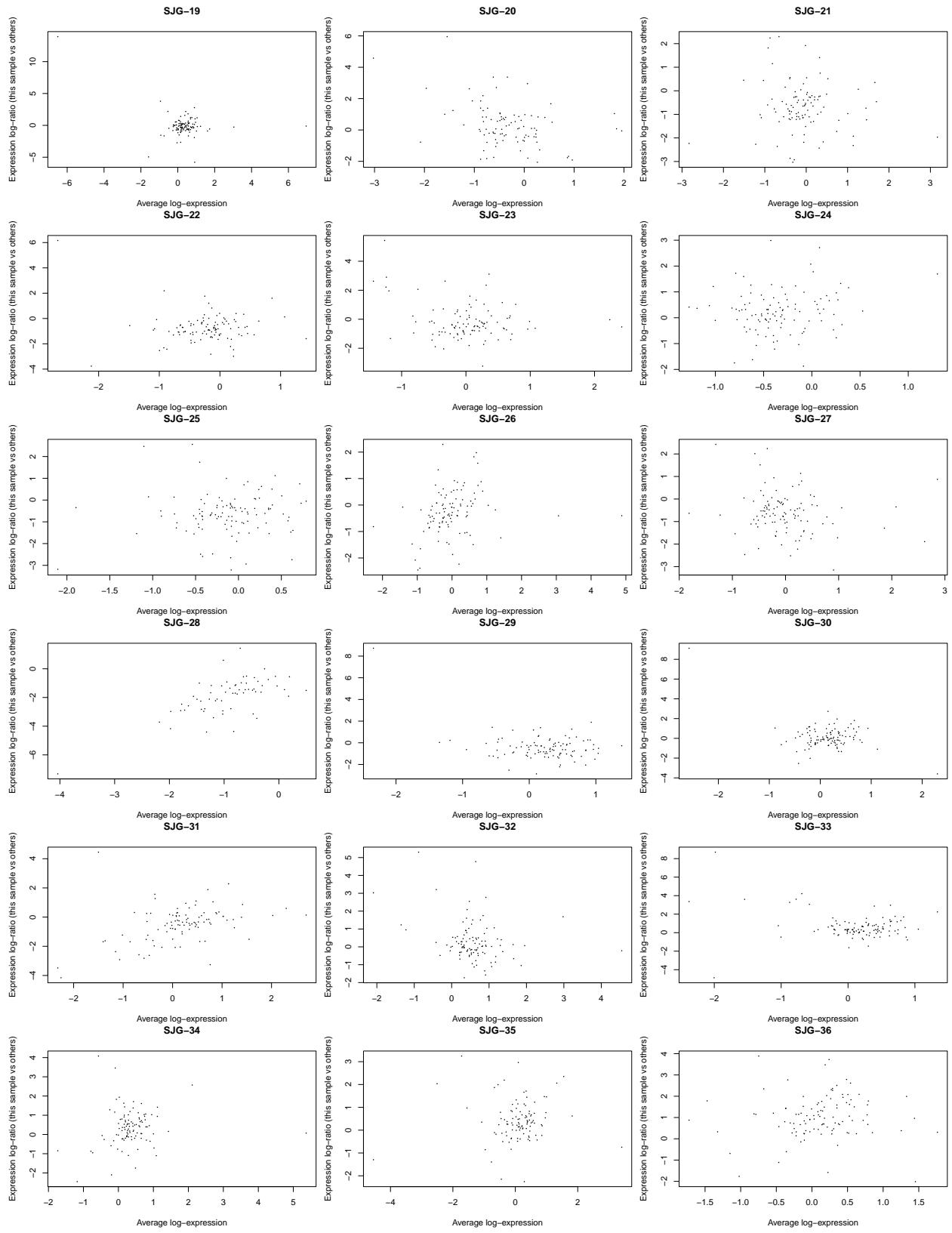


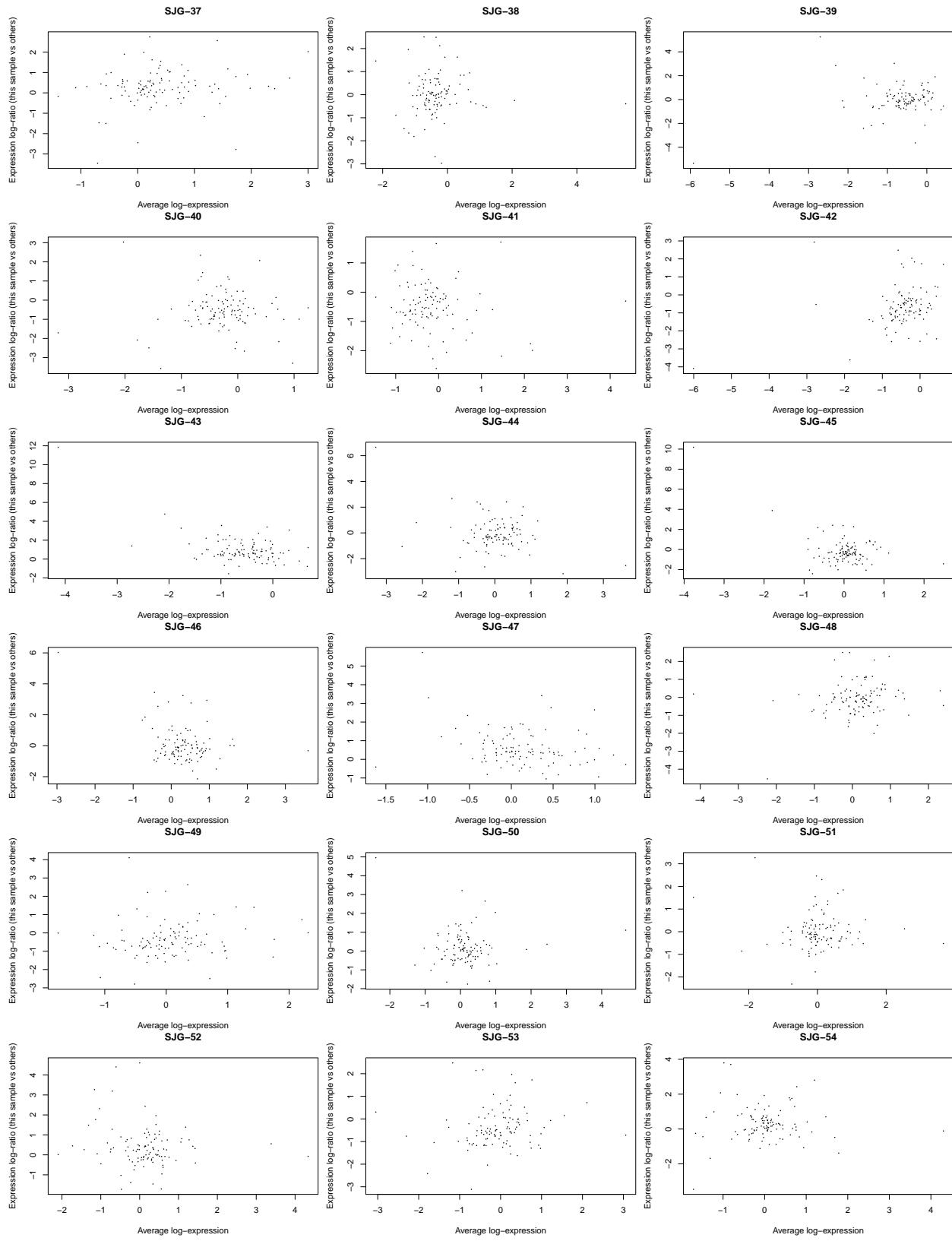
For sample's 20 & 90, 28 & 98 and 58 & 128, there is a similar cause to their unusual MA plots, except it is the second replicate that has the uniformly low RFI values.

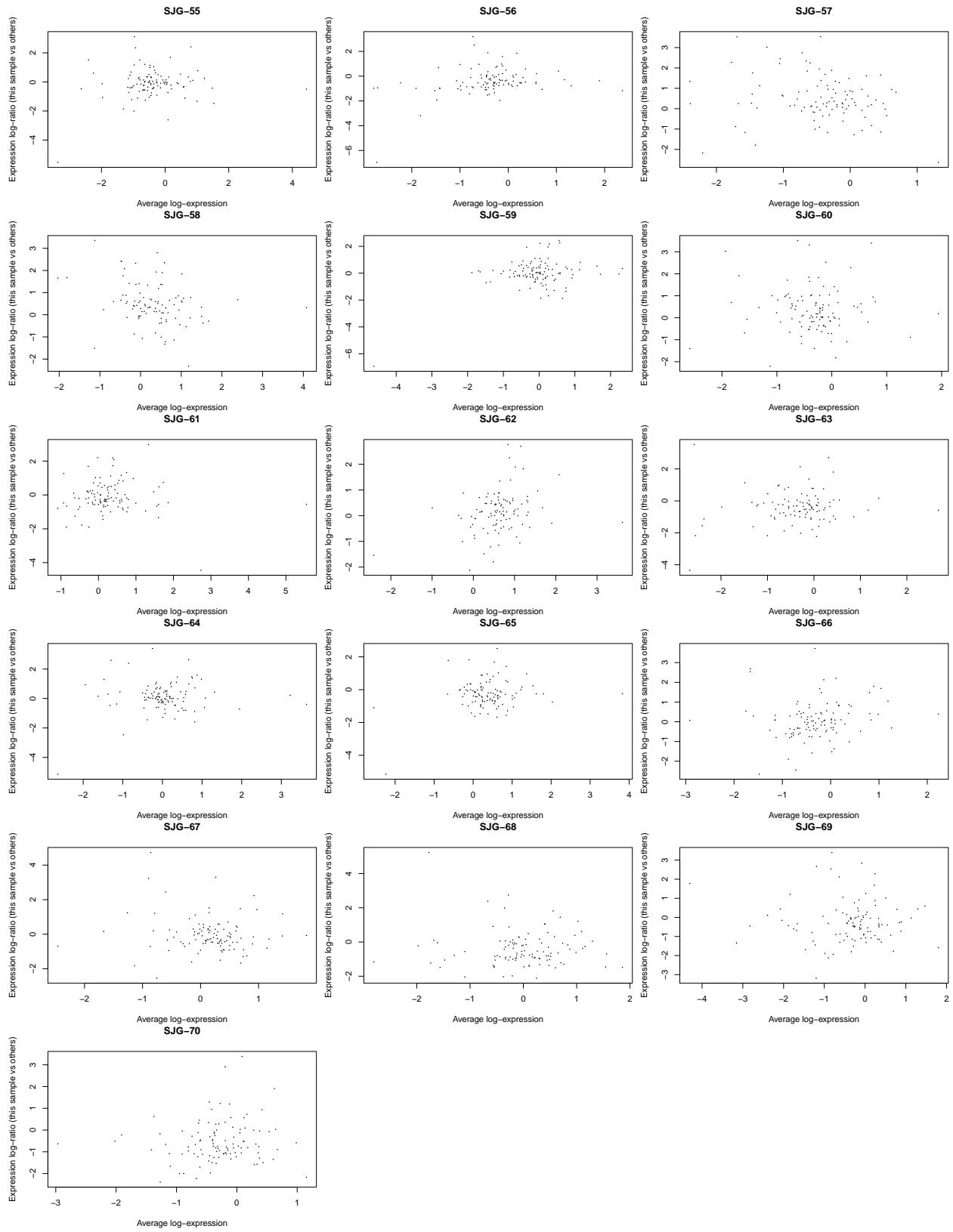
These plots suggest that it may be appropriate to remove an entire replicate and/or entire AB's that have a large proportion of low/NA RFI values. This is explored further in the "Filtering" document.

Median normalisation

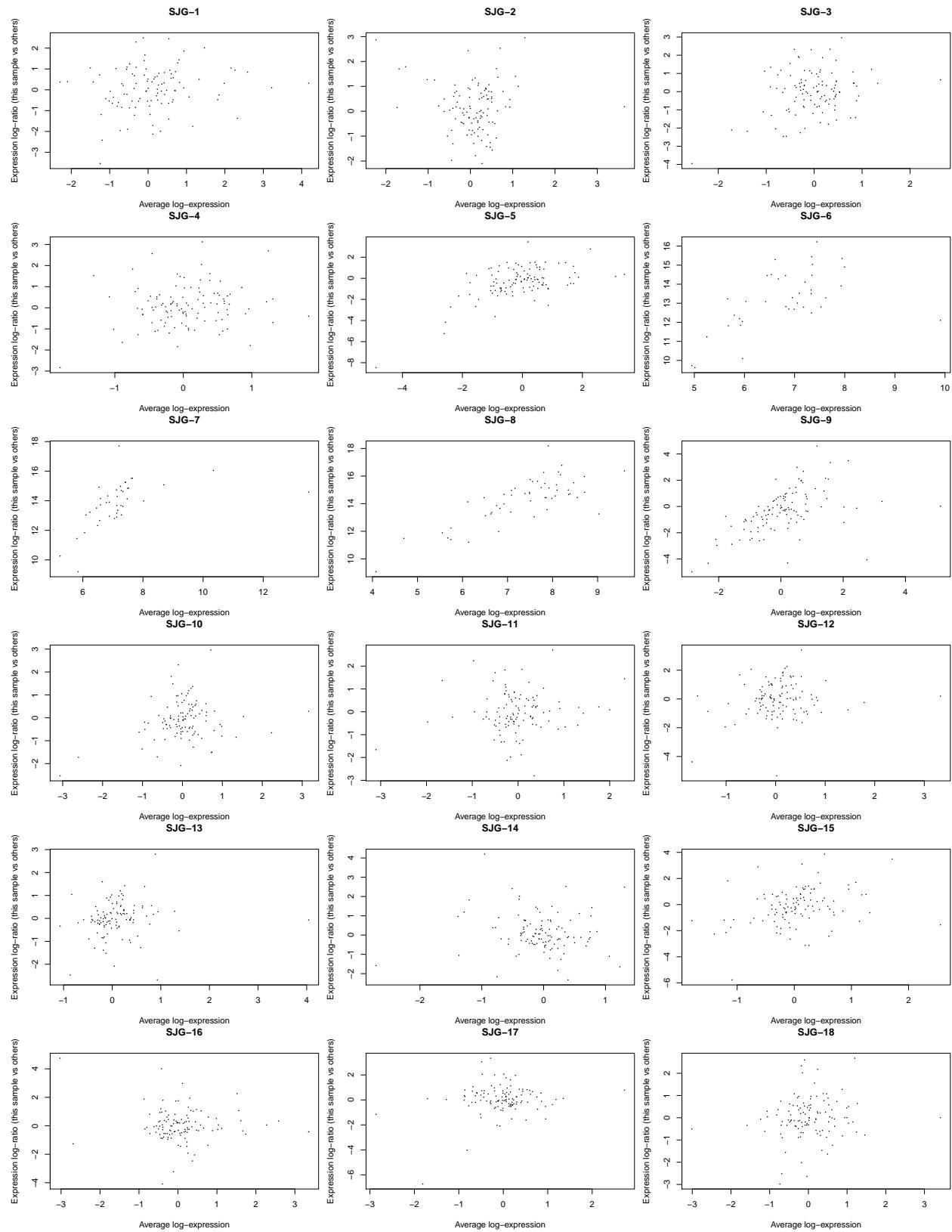


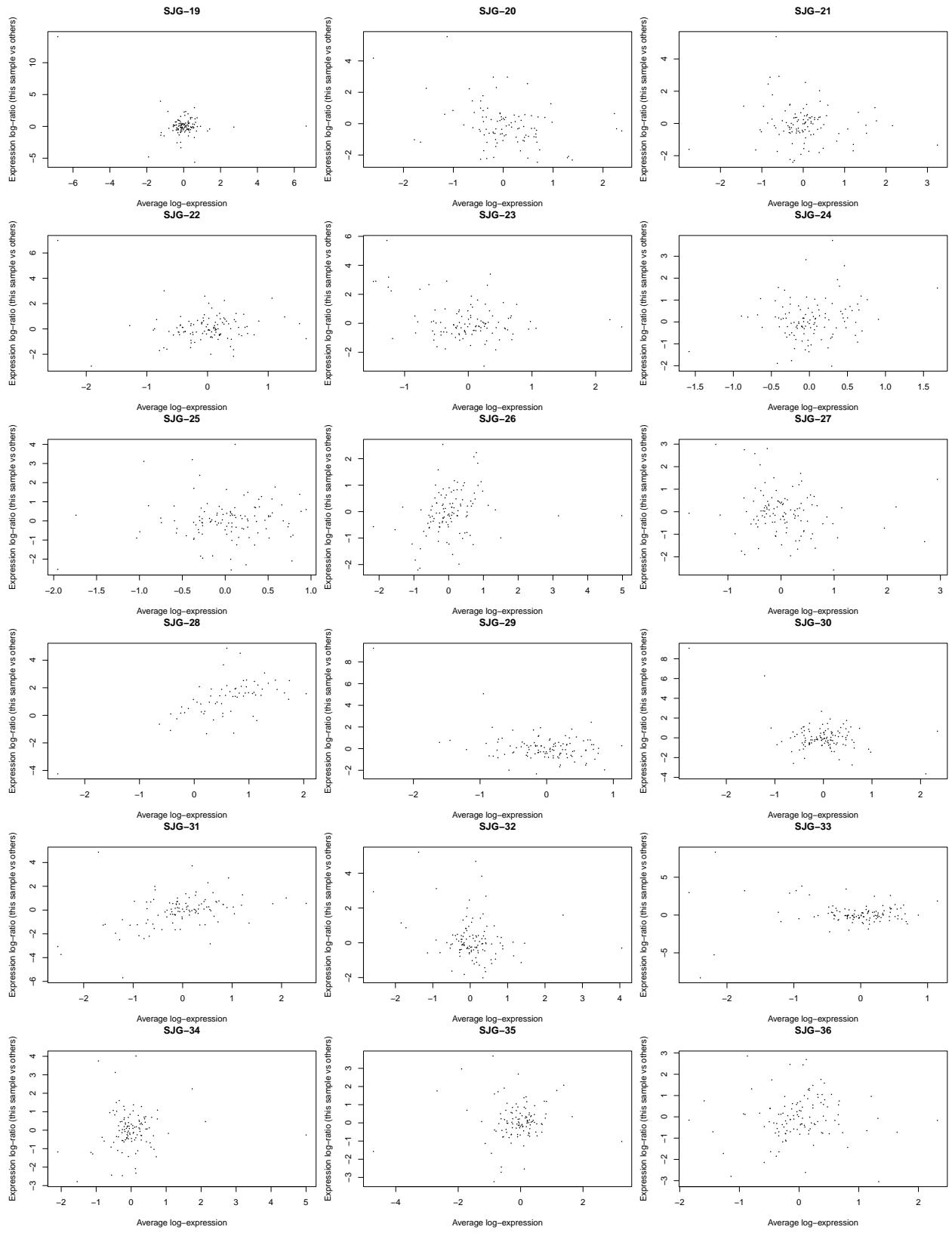


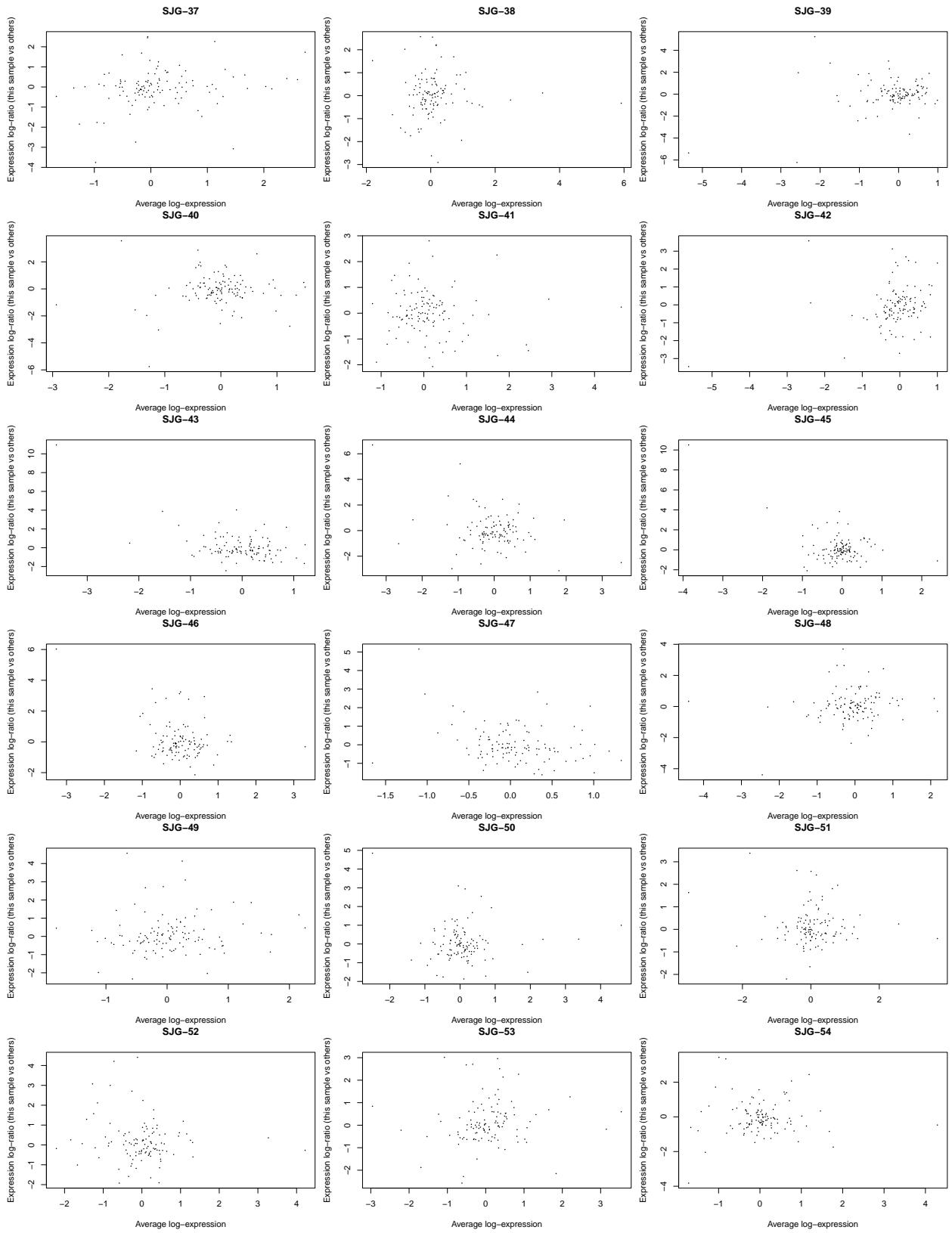


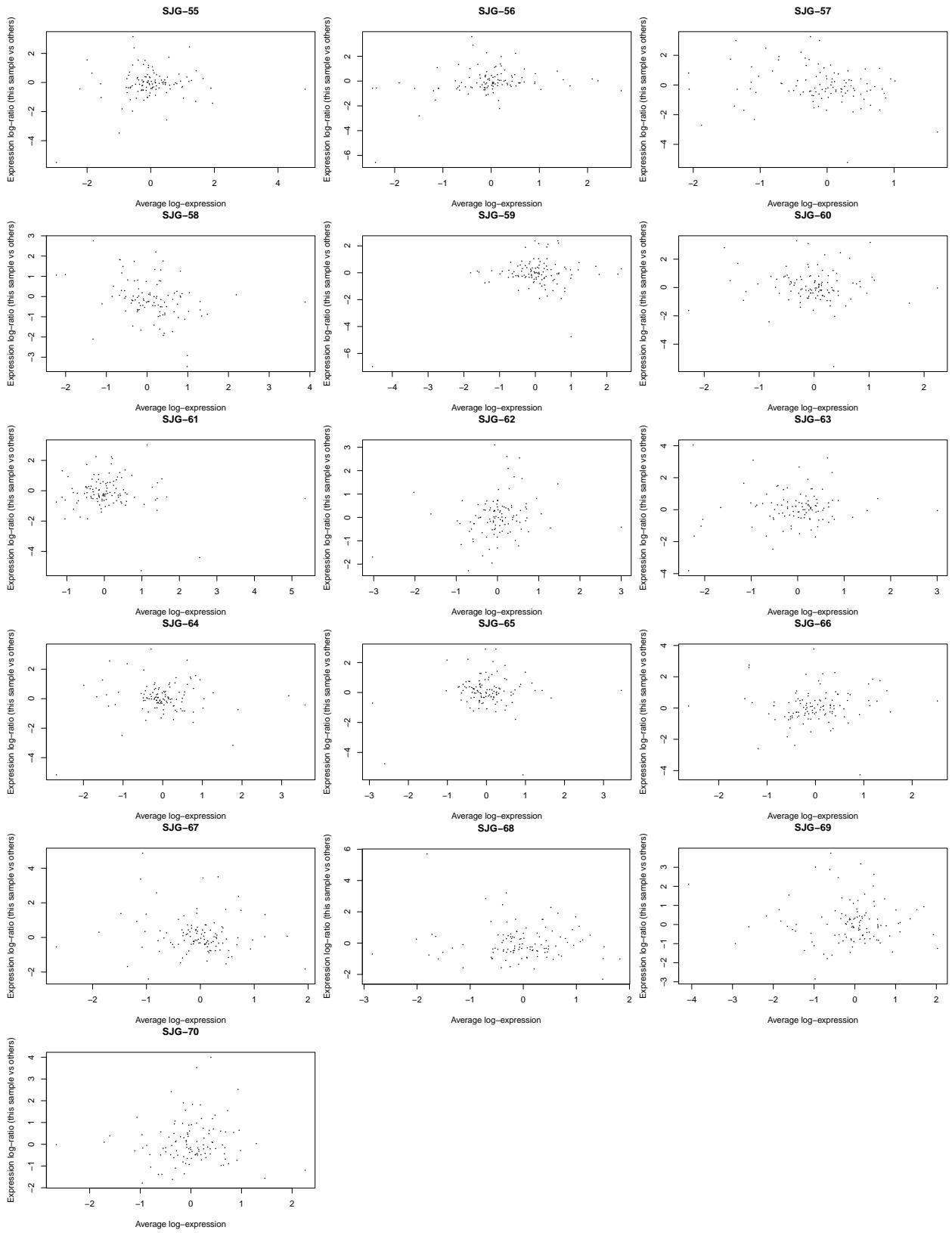


Loading control normalisation

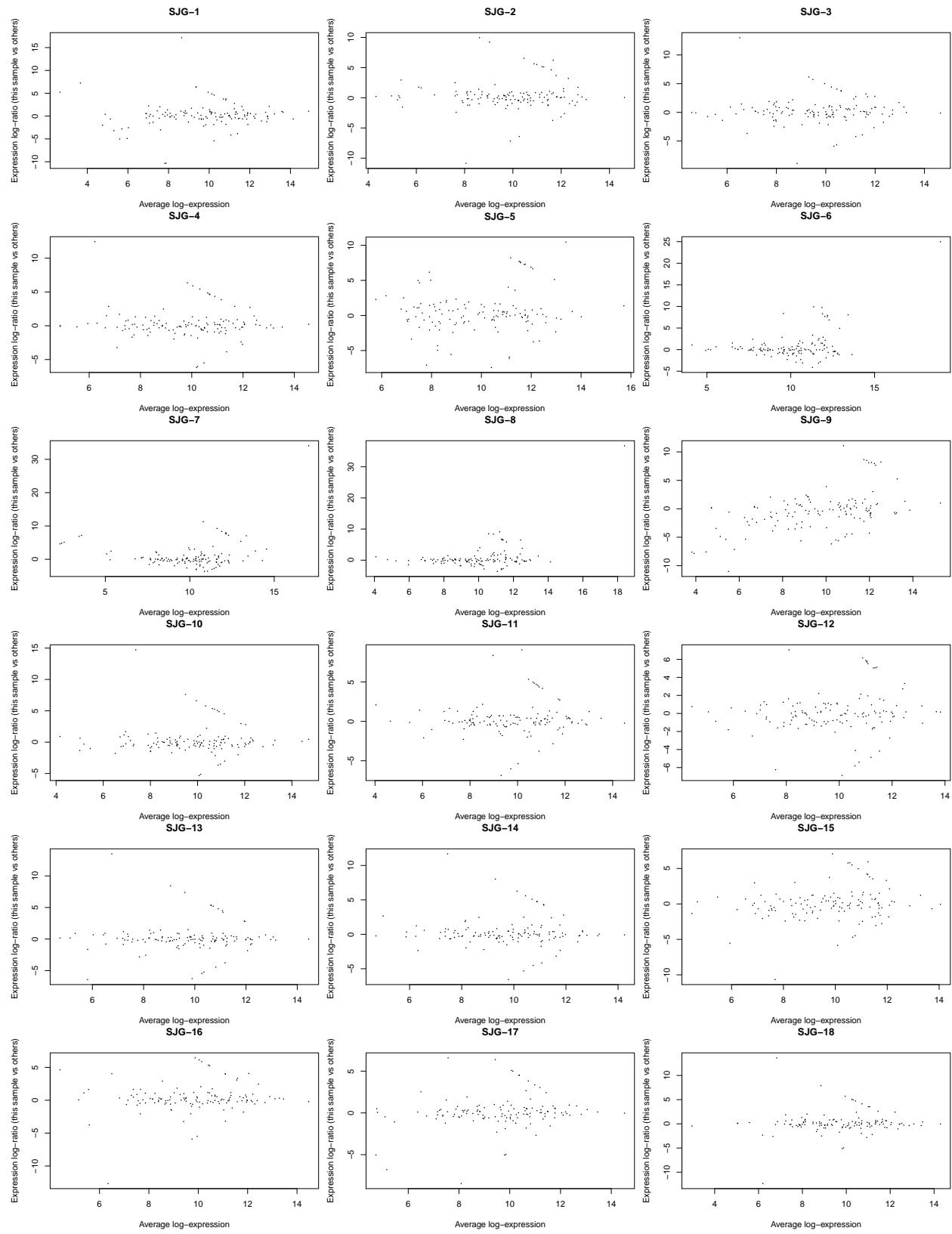


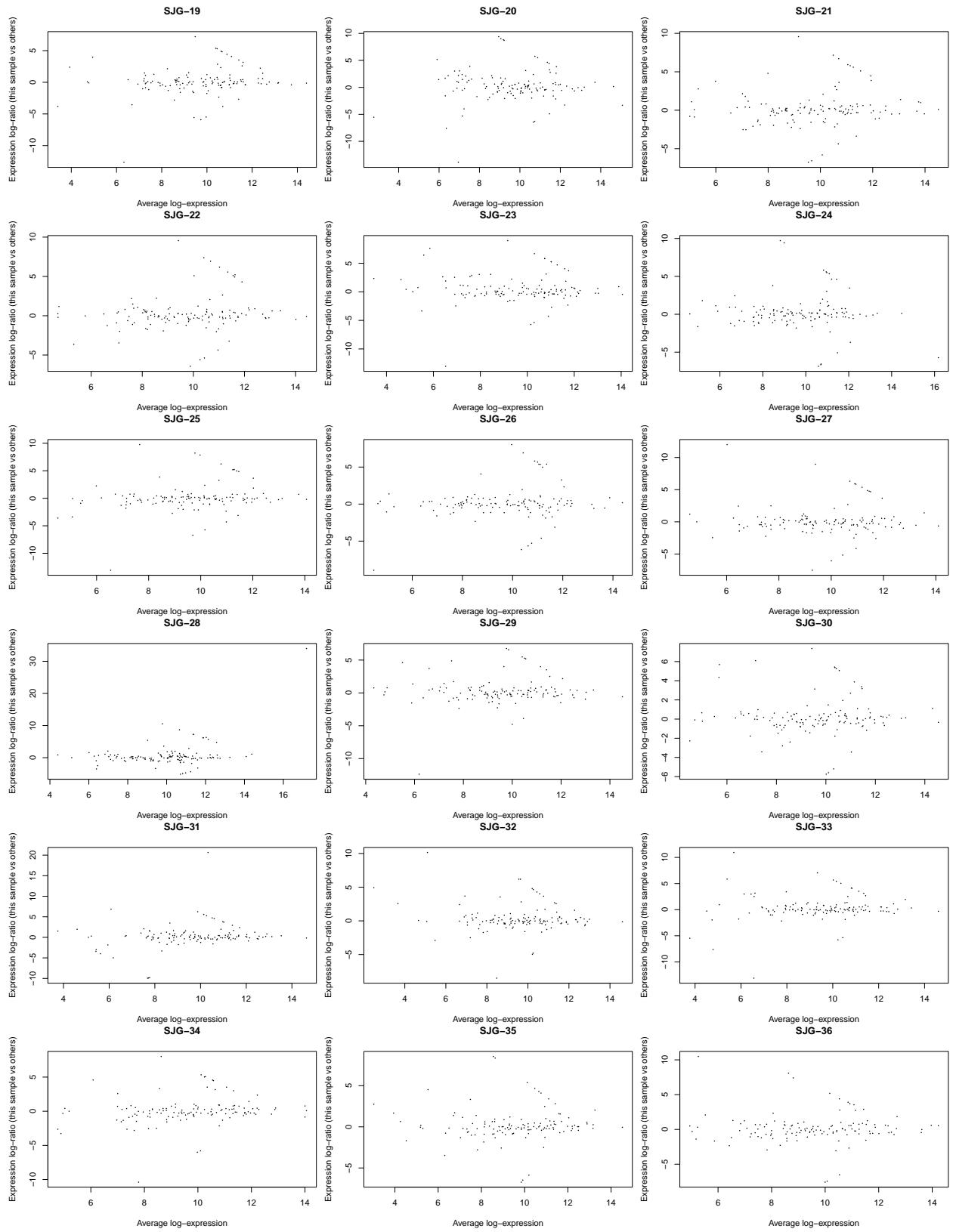


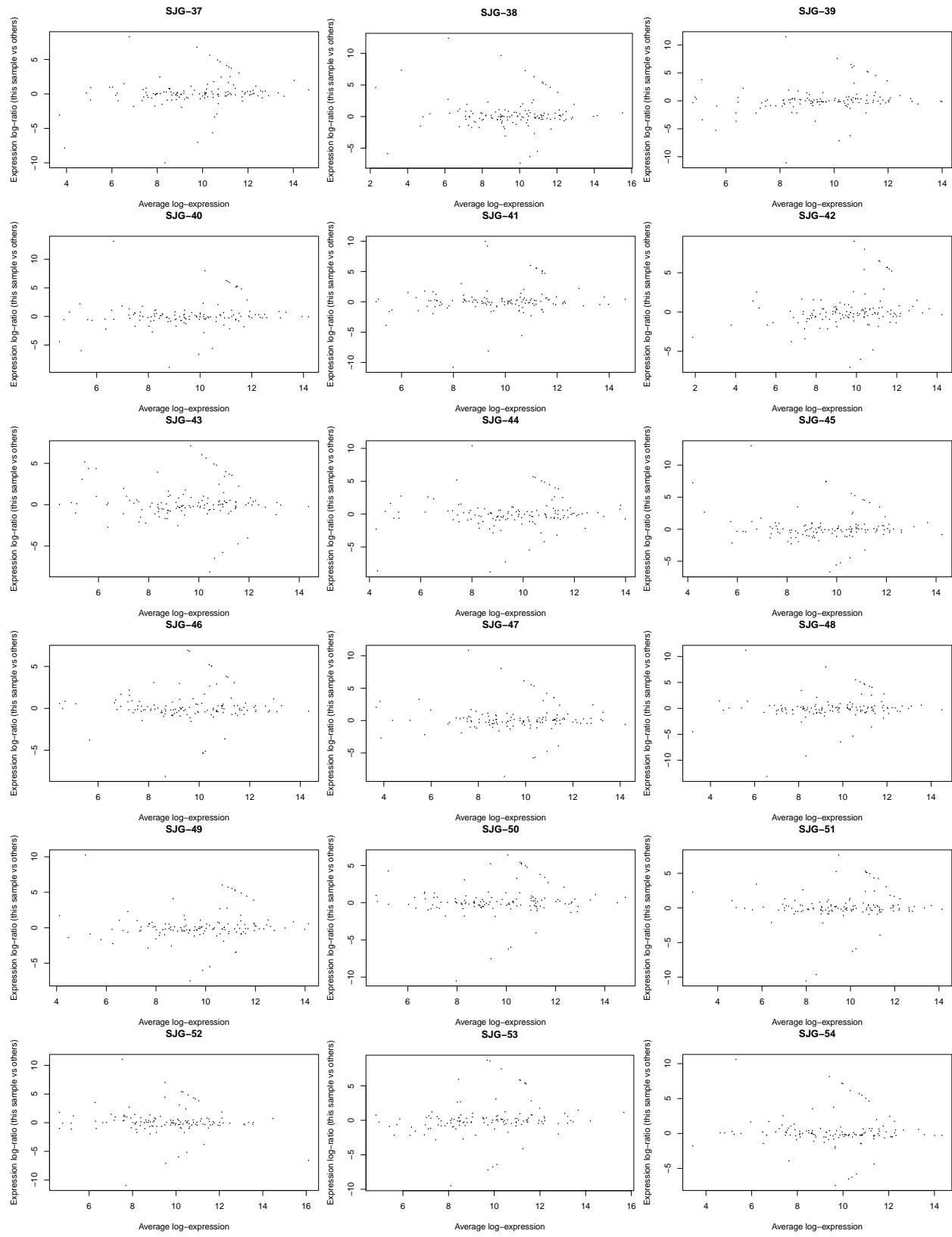


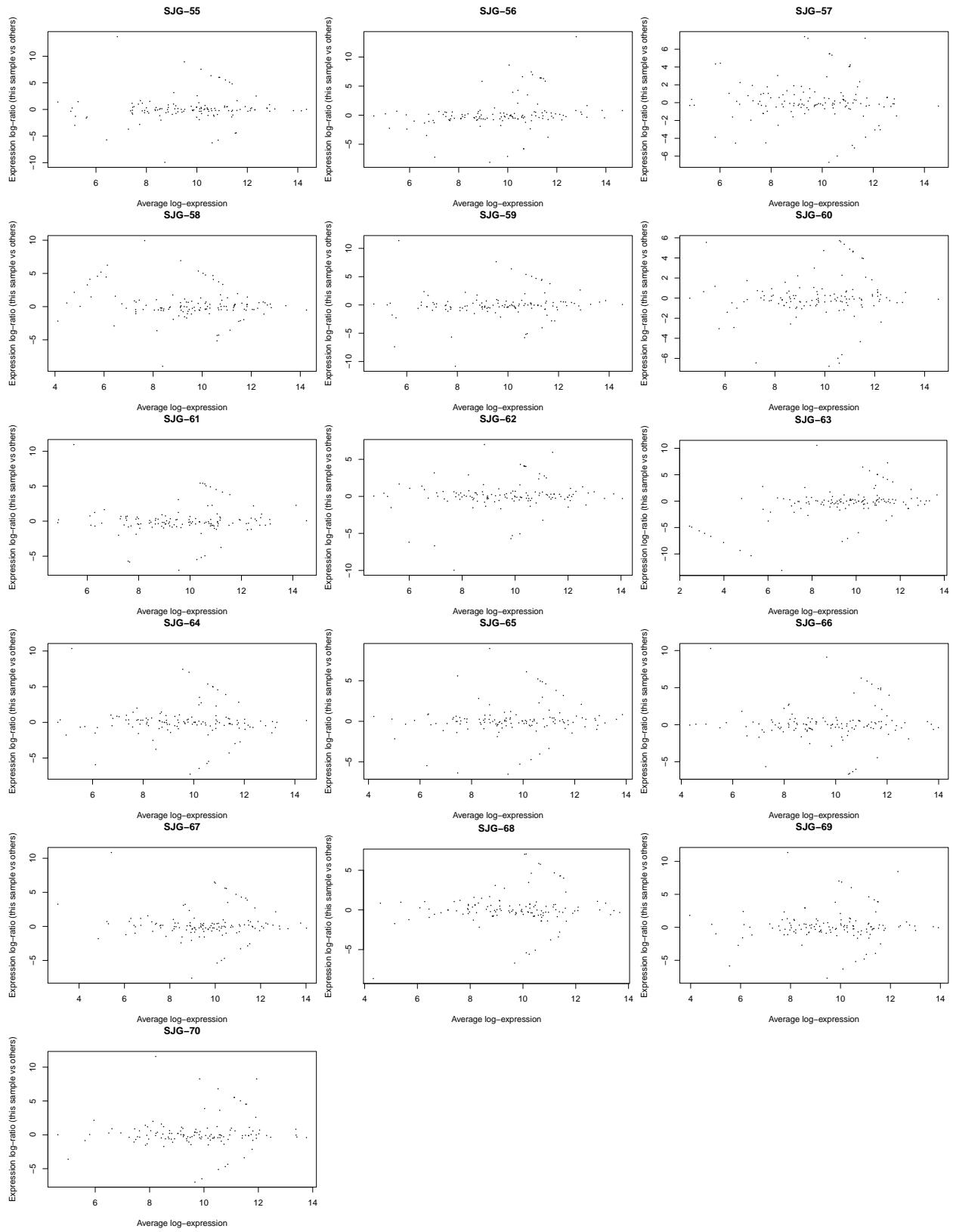


Global rank invariant



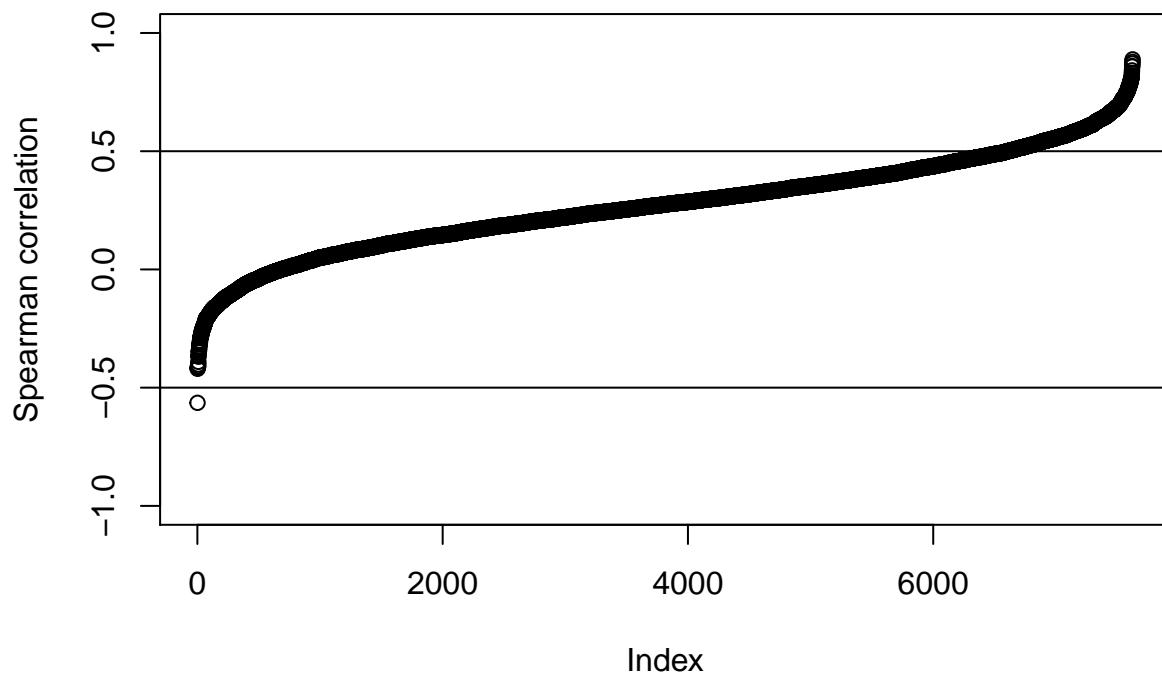




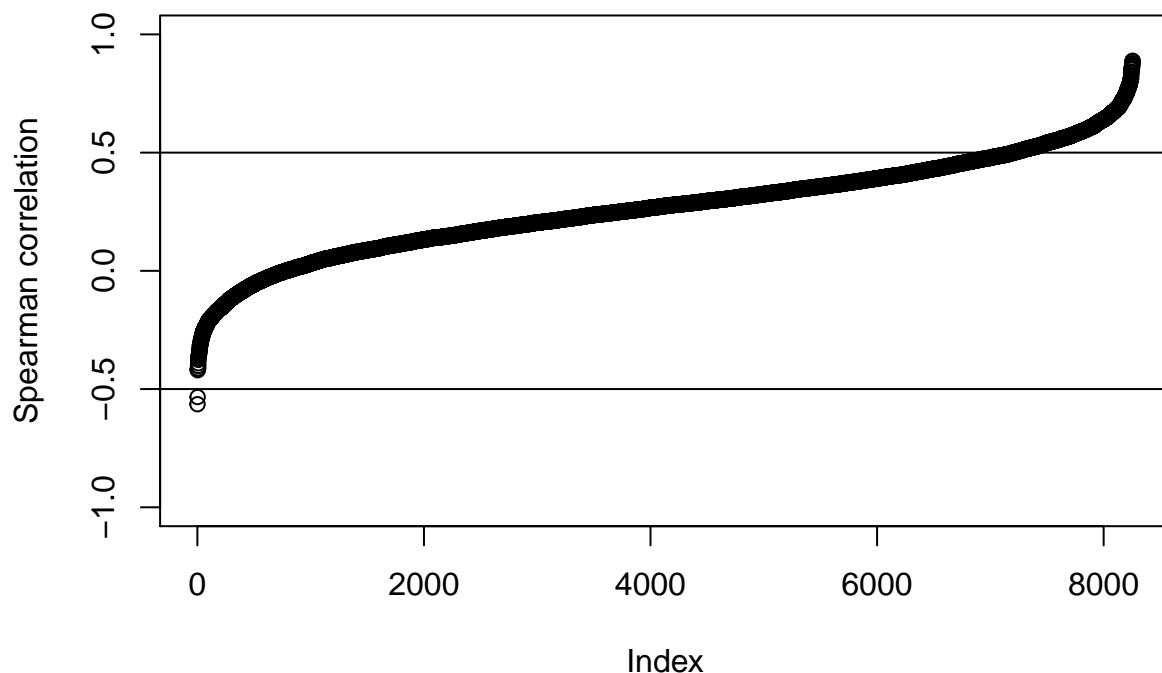


Pairwise correlation coefficient plots

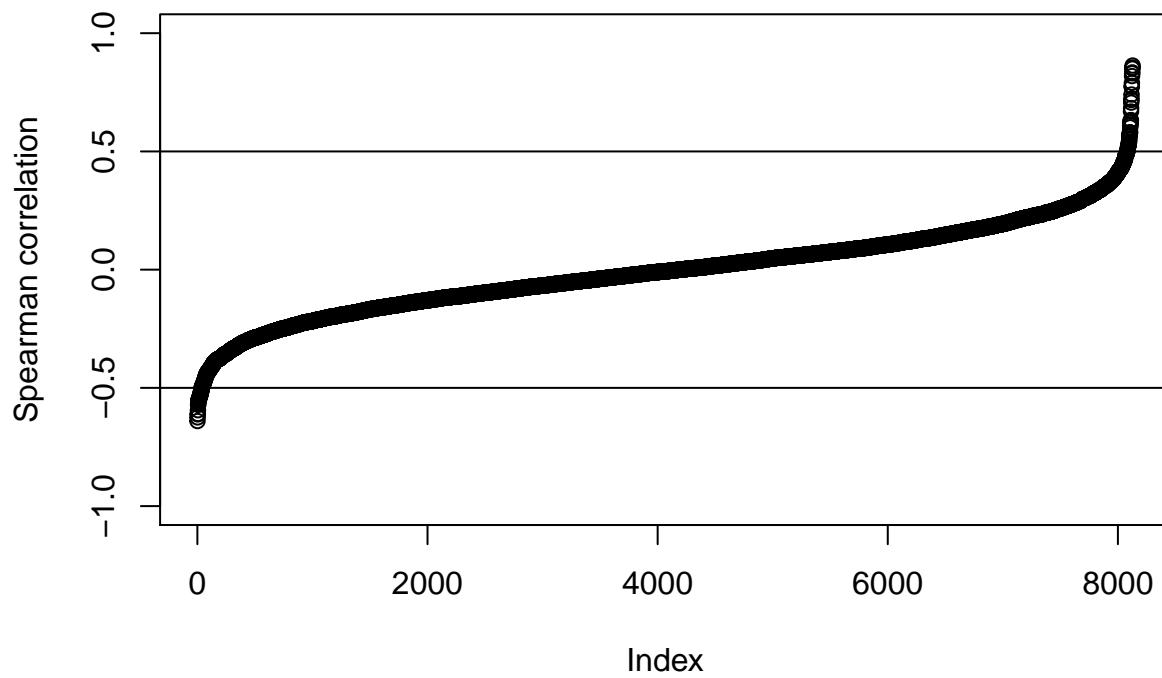
Raw data



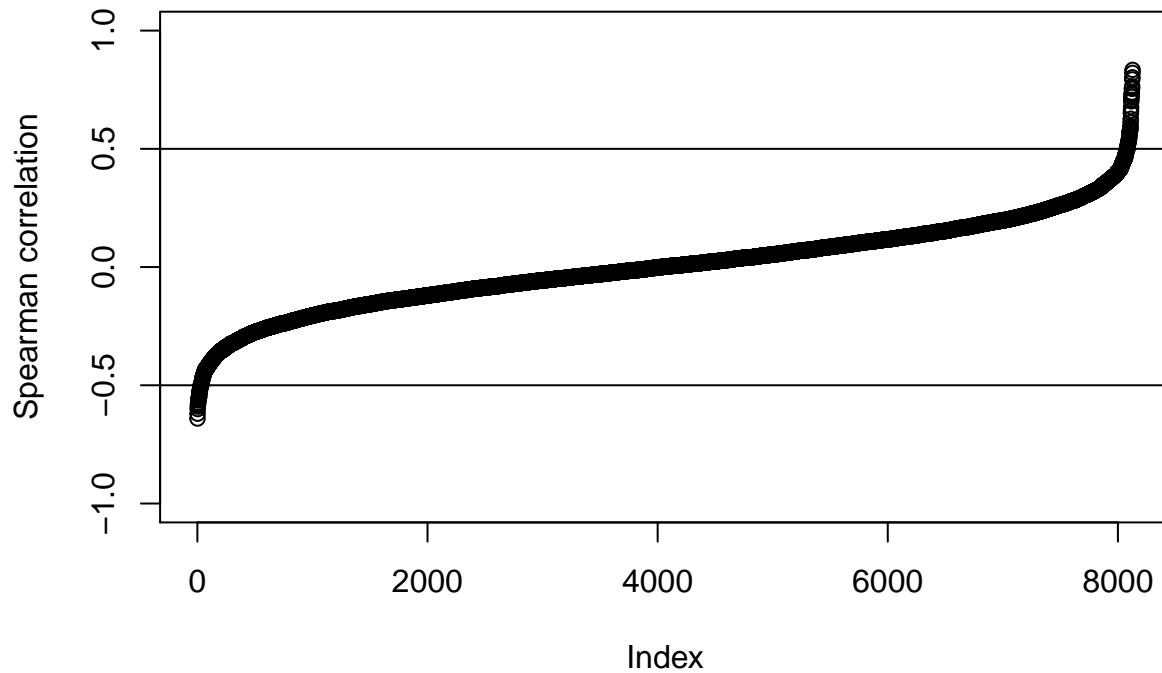
Median normalisation



Loading control normalisation



Global rank invariant



Normalised data

Average replicates then normalise

Normalise data using Global rank invariant and remove the 5 replicates with much poorer performance than their corresponding pair. Write out the data into a csv file. The steps are as follows:

1. Remove the 5 replicates that performed poorly. There were 5 samples where 1 replicate performed much worse than the other replicate.
2. Take the mean of the RFI values from the two replicates, for each sample. For samples where 1 replicate was removed, the RFI from the remaining replicate will be used.
3. Use Global rank invariant to normalise the data.
4. Write to a csv file.

Normalise both replicates

Normalise both replicates. We first look at how many NA's each AB has:

Antibody.Name	countNAs
Beta Catenin_P Ser33, Ser37, Thr41	70
Jak2_P Tyr1007, Tyr1008	70
MEK1/2	70
FLT3_P Tyr591	65
HER2/ErbB2_P Tyr1221, Tyr1222	65
JAK1_P Tyr1022, Tyr1023	65
Lyn_P Tyr507	65
N-Cadherin (D4R1H) XP	65
NDRG1	65
S6 Ribosomal Protein	65
Smad2_P Ser465, Ser467; Smad3_P Ser423, Ser425	65
Zap-70	65
Bad_P Ser112	28
Caspase 7_Cleaved Asp198	28
PKC pan_P bII Ser660	28
Caspase 3_Cleaved Asp175	23
Cyclin E1	15
BRCA1	14
Caveolin-1	14
cdc2_P Tyr15	14
cdc25c_P Ser216	14

Looking at the number of NA values, I have set the threshold at 30, removing all ABs with more than 30 NA values. The steps are:

1. Remove the 5 replicates that performed poorly. There were 5 samples where 1 replicate performed much worse than the other replicate.
2. Remove ABs with more than 30 missing values.
3. Use Global rank invariant to normalise the data.
4. Write to a tsv file.

References

- Gandolfo, Luke C., and Terence P. Speed. 2018. “RLE Plots: Visualizing Unwanted Variation in High Dimensional Data.” *PLOS ONE* 13 (2): e0191629. doi:10.1371/journal.pone.0191629.
- Liu, Wenbin, Zhenlin Ju, Yiling Lu, Gordon B. Mills, and Rehan Akbani. 2014. “A Comprehensive Comparison of Normalization Methods for Loading Control and Variance Stabilization of Reverse-Phase Protein Array Data , A Comprehensive Comparison of Normalization Methods for Loading Control and Variance Stabilization of Reverse-Phase Protein Array Data.” *Cancer Informatics* 13 (January): CIN.S13329. doi:10.4137/CIN.S13329.
- Pelz, Carl R., Molly Kulesz-Martin, Grover Bagby, and Rosalie C. Sears. 2008. “Global Rank-Invariant Set Normalization (GRSN) to Reduce Systematic Distortions in Microarray Data.” *BMC Bioinformatics* 9 (1): 520. doi:10.1186/1471-2105-9-520.
- Wachter, Astrid, Stephan Bernhardt, Tim Beissbarth, and Ulrike Korf. 2015. “Analysis of Reverse Phase Protein Array Data: From Experimental Design Towards Targeted Biomarker Discovery.” *Microarrays* 4 (4): 520–39. doi:10.3390/microarrays4040520.