

METHODOLOGY

Visualization methods for differential expression analysis

Lindsay Rutter^{1*}, Adrienne N Moran Lauter², Michelle A Graham³ and Dianne Cook⁴

*Correspondence:

lindsayannerutter@gmail.com

¹Bioinformatics and

Computational Biology Program,

Iowa State University, Ames, USA

Full list of author information is available at the end of the article

Abstract

Background: Despite the availability of many ready-made testing software, reliable detection of differentially expressed genes in RNA-seq data is not a trivial task. Even though the data collection is considered high-throughput, data analysis has intricacies that require careful human attention. Researchers should use modern data analysis techniques that incorporate visual feedback to verify the appropriateness of their models. While some RNA-seq packages provide static visualization tools, their capabilities should be expanded and their meaningfulness should be explicitly demonstrated to users.

Results: In this paper, we 1) introduce new interactive RNA-seq visualization tools, 2) compile a collection of examples that demonstrate to biologists *why* visualization should be an integral component of differential expression analysis. We use public RNA-seq datasets to show that our new visualization tools can detect normalization issues, differential expression designation problems, and common analysis errors. We also show that our new visualization tools can identify genes of interest in ways undetectable with models. Our R package “bigPint” includes the plotting tools introduced in this paper, many of which are unique additions to what is currently available. The “bigPint” website is located at <https://lindsayrutter.github.io/bigPint> and contains short vignette articles that introduce new users to our package, all written in reproducible code.

Conclusions: We emphasize that interactive graphics should be an indispensable component of modern RNA-seq analysis, which is currently not the case. This paper and its corresponding software aim to 1) persuade users to slightly modify their differential expression analyses by incorporating statistical graphics into their usual analysis pipelines, 2) persuade developers to create additional complex and interactive plotting methods for RNA-seq data, possibly using lessons learned from our open-source codes. We hope our work will serve a small part in upgrading the RNA-seq analysis world into one that more holistically extracts biological information using both models and visuals.

Keywords: Interactive; RNA-sequencing; Statistical graphics; Visualization

1 Background

RNA-sequencing (RNA-seq) uses next-generation sequencing (NGS) to estimate the quantity of RNA in biological samples at given timepoints. In recent years, decreasing cost and increasing throughput has rendered RNA-seq an attractive form of transcriptome profiling. Prior to RNA-seq, gene expression studies were performed

with microarray techniques, which required prior knowledge of reference sequences. RNA-seq does not have this limitation, and has enabled a new range of applications such as *de novo* transcriptome assembly [1] and detection of alternative splicing processes [2, 3]. Coupled with its high resolution and sensitivity, RNA-seq is revolutionizing our understanding of the intricacies of eukaryotic transcriptomes [4, 5].

One common format of RNA-seq data is a matrix containing mapped read counts for n rows of genes and p columns of samples. These mapped read counts provide gene expression level estimations across samples. Researchers often conduct RNA-seq studies to identify differentially expressed genes (DEGs) between treatment groups. In most popular RNA-seq analysis packages, this objective is approached with models, such as the negative binomial model [6, 7, 8, 9] and linear regression models [10].

Initially, it was widely claimed that RNA-seq produced unbiased data that did not require sophisticated normalization [4, 11, 12]. However, numerous studies have since revealed that RNA-seq data is replete with biases and that accurate detection of DEGs is not a negligible task. Problems that complicate RNA-seq data analysis include nucleotide and read-position biases [13], biases related to gene lengths and sequencing depths [14, 15], biases introduced during library preparation [16], and confounding combinations of technical and biological variability [17].

In light of these complications, researchers should analyze RNA-seq data like they would any other biased multivariate data. Solely applying models to such data is problematic because models hold assumptions that must be verified to ensure statistical soundness. Fortunately, data visualization enables researchers to see patterns and problems they may not otherwise detect with traditional modeling. As a result, the most effective approach to data analysis is to iterate between models and visuals, and enhance the appropriateness of applied models based on feedback from visuals [18]. With differential expression data, we primarily want to compare the

variability between replicates and between treatment groups. This is visually best achieved by drawing the mapped read count distributions across all genes and samples. To our knowledge, the few plotting tools offered in popular RNA-seq packages do *not* often allow users to effectively view their data in this manner.

In this paper, we strive to remedy this problem by highlighting the utility of new and effective differential expression plotting tools. We use real RNA-seq data to show that our tools can detect normalization problems, DEG designation problems, and common errors in the analysis pipeline. We also show that our tools can identify genes of interest that cannot otherwise be obtained by models. We emphasize that interactive graphics should be an indispensable component of modern RNA-seq analysis. Here, we do not propose that users drastically change their approach to differential expression analysis. Instead, we propose that users simply modify their approach to differential expression analysis by assessing the sensibility of their models with multivariate graphical tools, namely with parallel coordinate plots, scatterplot matrices, and litre plots.

Results

Parallel coordinate plots

Parallel coordinate plots are essential to inform the relationships between variables in multivariate data. A parallel coordinate plot draws each row (gene) as a line. For a given gene, two samples with similar read counts will have a flat connection and two samples with dissimilar read counts will have a sloped connection. The ideal dataset has more variability between treatments than between replicates. Researchers can quickly confirm this with a parallel coordinate plot: There should be flat connections between replicates but crossed connections between treatments.

There are several packages within the Bioconductor software [19] that provide graphics for RNA-seq data analysis [20]. Two of the most common graphic techniques are side-by-side boxplots and Multidimensional Scaling (MDS) plots

[21, 22, 9, 23]. Unfortunately, these plots can hide problems that still exist in the data even after normalization and that could be better detected with parallel coordinate plots.

Figure 1 exemplifies this problem for two *simulated* datasets, one displayed on the left half and the other displayed on the right half of the figure. Each dataset contains two treatment groups (A and B) with three replicates. The side-by-side boxplots (subplots A) both show fairly consistent medians across the six samples in the left and right dataset; the most prominent difference is the smaller interquartile ranges in the right dataset. The left MDS plot separates the treatment groups distinctively; the right MDS plot suggests a similar separation but in a much subtler manner (subplots B). In addition, the first replicate from treatment A appears as an outlier in the right MDS plot.

While the boxplots and MDS plots provide useful information, the parallel coordinate plots (subplots C) show an additional meaningful difference between the left and right datasets. The left dataset has consistent (level) lines between replicates and inconsistent (crossed) lines between treatment groups. This suggests that some of the genes (lines) have consistently low values for treatment group A and consistently high values for treatment group B, while some genes have the opposite phenomenon. As a result, the majority of the plotted genes may be DEG candidates. In contrast, the right dataset does not possess this ideal structure and suggests that the majority of its genes may not be DEG candidates. We could not see this important distinction as clearly using the side-by-side boxplots or the MDS plots because they only provide data summarization at the sample resolution, while the parallel coordinate plots show the sample connections for each gene in the data.

Please note that the example above was simulated for didactic purposes. We will now examine the application of parallel coordinate plots to real data from an RNA-seq study that compared soybean leaves after 120 minutes of iron-sufficient (group

P) and iron-deficient (group N) hydroponic treatments [24]. We filtered genes with low means and/or variance, performed a hierarchical clustering analysis with a cluster size of four, retained only significant genes, and visualized the results using parallel coordinate lines (Figure 2). For these visualizations, we standardized each gene to have a mean of zero and standard deviation of unity [25, 26]. Then, we performed hierarchical clustering on the standardized DEGs using Ward's linkage. This process can divide large DEG lists into smaller clusters of similar patterns, which allows us to more effectively detect the various types of patterns within large DEG lists. We note that the number and quality of clusters can vary depending on the data.

The majority of significant genes were in Clusters 1 and 2, which for the most part captured the expected patterns of differential expression (consistent replicates and inconsistent treatments) in reverse directions. Only 17 significant genes belonged to Cluster 4 and they mostly showed messy patterns with low signal to noise ratios. Interestingly, Cluster 3 had a fairly large number of significant genes ($n=861$). These genes mostly showed clean differential expression profiles similar to Cluster 2 (large values for group N and small values for group P), except for unexpectedly large values for the third replicate of group P. The reasons for a different response by these genes on this replicate is unclear, but warrants further study.

Scatterplot matrices

Overview of scatterplot matrices

A scatterplot matrix is another effective multivariate visualization tool that plots read count distributions across all genes and samples. Specifically, it represents each row (gene) as a point in each scatterplot. With this method, users can quickly discover unexpected patterns, recognize geometric shapes, and assess the structure and association between multiple variables in a manner that is different from most common practices.

114 Clean data would be expected to have larger variability between treatment groups
115 than between replicates. As Figure 3 shows, researchers can quickly confirm this
116 with a scatterplot matrix. Within each scatterplot, most genes should fall along the
117 $x=y$ line (in red) as we expect only a small proportion of them to show differential
118 expression between samples. However, a fraction of the genes should have lower
119 variability between replicates than between treatments, and so we should expect
120 the spread of the scatterplot points to fall more closely along the $x=y$ relationship
121 between replicates than between treatments. Indeed, in Figure 3, we created a scat-
122 terplot matrix for a public RNA-seq dataset that contains three replicates for two
123 developmental stages of soybean cotyledon (S1 and S2) [27]. We can immediately
124 verify that the nine scatterplots between treatment pairs (the bottom-left corner of
125 the matrix encased in the blue square) have more spread around the $x=y$ line than
126 the six scatterplots between replicate pairs.

127 After confirming this expected trend, users can use the scatterplot matrix to
128 focus on subsets of genes: Outlier genes that deviate from the $x=y$ line in replicate
129 scatterplots might be problematic, whereas outlier genes that deviate from the $x=y$
130 line in treatment scatterplots might be DEGs. In order to achieve this functionality,
131 the plots must be rendered interactive. This way, users can hover over and click on
132 gene subsets of interest and view their patterns from multiple perspectives while
133 also obtaining their identifiers.

134 Notice that each gene in our data is plotted once in each of the 15 scatterplots.
135 With 73,320 genes in our data, more than one million points must be plotted.
136 Rendering all points interactive would slow down the interactive capabilities of
137 the plot. To solve this, we can tailor the geometric object of the scatterplots to be
138 hexagon bins rather than points. This dramatically reduces the number of geometric
139 objects to be plotted, and increases the interactivity speed.

140 The interactive version of Figure 3 is available online [28]. Readers can read the
141 “About” Tab to fully understand how to use the application. Essentially, the user
142 can hover over a hexagon bin to see how many genes it contains. When the user
143 clicks on a hexagon bin, the names of the genes are listed and superimposed as
144 orange points across all scatterplots. The genes are also linked to a second plot that
145 superimposes them as parallel coordinate lines on a side-by-side boxplot of *all* gene
146 counts in the dataset. This interactivity and linking allows users to quickly examine
147 genes of interest from multiple viewpoints superimposed onto the summary of *all*
148 genes in the dataset.

149 Assessing normalization with scatterplot matrices

150 There is still substantial discussion about the normalization of RNA-seq data, and
151 the scatterplot matrix can be used to understand and assess various algorithms. To
152 exemplify this point, we will use a publicly-available RNA-seq dataset on *Saccha-*
153 *romyces cerevisiae* (yeast) grown in YP-Glucose (YPD) [22]. The data contained
154 four cultures from independent libraries that were sequenced using two library
155 preparation protocols and either one or two lanes in a total of three flow-cells.
156 This experimental design allowed researchers to examine various levels and com-
157 binations of technical effects (library preparation and protocol and flow cell) and
158 biological effects (culture).

159 The four cultures (Y1, Y2, Y4, and Y7) were treated as biological replicates for
160 which differential expression was not expected. Hence, the authors could estab-
161 lish a false positive rate in relation to the number of DEGs called between these
162 groups. They then demonstrated that within-lane regression alone was insufficient
163 in effectively removing biases. Instead, aggressive corrections for both within-lane
164 (GC-content and gene length) and between-lane (count distribution and sequencing
165 depth) biases were needed to effectively reduce the false-positive rate of DEG calls.

Figure 4A shows the scatterplot matrix of the read counts from the Y1 and Y4 treatments after within-lane normalization. As we stated earlier, we expect most genes to show similar expression between samples, except for the handful that are differentially expressed. However, it is immediately clear that the data still was not sufficiently normalized as the distribution of genes is not centered around the $x=y$ lines. In contrast, Figure 4B shows the scatterplot matrix of the read counts from the Y1 and Y4 treatments after *both* within-lane and between-lane normalization, as was recommended by the authors due to its reduced false-positive rate. Indeed, the scatterplot matrix now follows the expected structure with most genes falling along the $x=y$ line with thicker deviations from it between treatment groups than between replicate groups.

Additionally, we can also confirm from Figure 4B that the read counts fall closer to the $x=y$ line between the Y4 replicates (bottom-right scatterplot) than between the Y1 replicates (top-left scatterplot). This is expected because the Y1 replicates had additional technical variability as they used two different flow cells, whereas the Y4 replicates were from the same flow cell. As such, the scatterplot matrix can also be used to quickly inspect patterns of biological and technical variability in the dataset.

Checking for common errors with scatterplot matrices

Irreproducibility is prevalent in high-throughput biological studies. A study in *Nature Genetics* surveyed eighteen published microarray expression analyses and reported that only two were exactly reproducible [29]. The extent of the problem has spawned a field called “forensic bioinformatics” whereby researchers attempt to reverse-engineer reported results back into the raw datasets simply to derive the methodologies used in published studies [30].

Even though irreproducibility is merely cumbersome when it masks methods, it is unquestionably hazardous when it masks errors. With regards to personalized

193 medicine, for example, obscured errors may cause well-intentioned researchers to
194 present evidence for drugs that are ineffective or even harmful to patients [30].
195 Forensic bioinformaticians who have actively investigated common errors in high-
196 throughput biological studies have concluded that the largeness of the data itself
197 may hinder our ability to detect errors [30]. They also discovered that the most
198 common errors are simple errors, such as mixing up sample labels [30]. Collectively,
199 these findings suggest that simple errors can be difficult to detect using common
200 practices in high-throughput studies.

201 Fortunately, scatterplot matrices are a convenient tool to check for common er-
202 rors like sample mislabeling. Figure 5 shows the resulting scatterplot matrix after
203 we deliberately swapped the labels of the third replicate of the first treatment group
204 (S1.3) with the first replicate of the second treatment group (S2.1) in the previously-
205 mentioned cotyledon dataset [27]. We can immediately see that, as expected, there
206 are nine scatterplots with thicker distributions around the $x=y$ line and six scat-
207 terplots with thinner distributions around the $x=y$ line. However, we notice that a
208 subset of these thick and thin scatterplots appear outside of their expected locations
209 given the expected variability between treatments versus replicates. Rearranging the
210 columns of the two samples that appear suspicious in Figure 5 would indeed lead
211 us back to the clean-looking scatterplot matrix we saw in Figure 3. The scatterplot
212 matrix provides us convincing evidence of a mislabeling problem even down to the
213 gene level, which cannot be confirmed with such detail using traditional plots like
214 the boxplots and MDS plots before sample switching (left side of Figure 6) and after
215 sample switching (right side of Figure 6). While this method can inform suspicious
216 patterns in more detail than other means, users must still perform extra steps to
217 determine if these patterns more likely relate to mislabeling or some real biological
218 phenomenon. In the case of suspected mislabeling, the user would still need to sub-

stantiate this suspicion with decisive evidence and should only use the visualization as a guide.

Finding unexpected patterns in scatterplot matrices

Most popular RNA-seq plotting tools display summaries about the read counts, such as fold change summaries, principal component summaries, five number summaries, and dispersion summaries. In contrast to this trend, scatterplot matrices display the non-summarized read counts for all genes. This trait allows for geometric shapes and patterns relevant to the read count distribution to be readily visible in the scatterplot matrix.

An example of how geometric shapes in the scatterplot matrix can provide applicable information to researchers is shown in Figure 7, which uses the iron-metabolism soybean dataset [24]. After normalizing the data, we see the expected pattern of a scatterplot matrix, with more variation around the $x=y$ line between treatments than between replicates (Figure 7). However, one streak structure in the bottom right scatterplot stands out. A small subset of transcripts between replicates of the iron-sufficient group sharply deviates from the $x=y$ line. By interacting with the plot, we identified the five transcripts that deviated the most from the expected pattern, and searched for their putative functions. We discovered that these transcripts are reportedly involved in biotic and abiotic stress responses as well as the production of superoxides to combat microbial infections. It should be noted that these five transcripts did not reach significance unless the third replicate of the P group was removed. Therefore, these genes will still be reported as non-significant in this study.

Discussion with the authors of the study revealed that a lab biologist documented a clean data collection process. In the study, the authors determined the DEGs across three time points (30 minutes, 60 minutes, and 120 minutes) after exposure to the two iron condition levels. In order to reduce variability caused by plant

246 handling by different researchers, the same researcher collected the samples in suc-
247 cession. One major finding from their study was a vast change in gene expression
248 responses between these three time points (Figure 8). In light of these discoveries,
249 the authors tentatively suggest that the streak of genes shown in Figure 7 may be
250 due to the timing differences between replicate handling.

251 In any case, scientists cannot observe such interesting structures from any mod-
252 els. Hypothetically, these structures could lead to interesting post hoc analyses. For
253 instance, if a similar structure existed in data where the authors had noted an inad-
254 vertent experimental or biological discrepancy between those replicates, then a post
255 hoc hypothesis that these genes might respond to that discrepant condition could
256 be generated. We note this would only serve as a hypothesis generator; conventional
257 genetic studies and additional evidence would be needed to confirm any possible
258 role these genes have on this biological activity.

259 Assessing DEG calls in scatterplot matrices

260 The scatterplot matrix can also be used to quickly examine the DEGs returned
261 from a given model. Figure 9 shows the DEGs from the soybean cotyledon dataset
262 superimposed as orange points onto the scatterplot matrix [27]. We expect for DEGs
263 to fall along the $x=y$ line for scatterplots between replicates and deviate from the
264 $x=y$ line for scatterplots between treatment groups, as is confirmed in Figure 9.
265 As a side note, we could also link these DEGs as parallel coordinate lines on a
266 side-by-side boxplot to confirm the expected pattern of differential expression from
267 a second viewpoint. If we do not observe what should be expected of DEGs, then
268 the DEG calls from the model may need to be scrutinized further.

269 Litre plots

270 We demonstrated how to view DEGs onto the Cartesian coordinates of the scat-
271 terplot matrix in Figure 9. Unfortunately, this figure becomes limited when we
272 investigate treatment groups that contain a large number of replicates because we

273 then have too many small scatterplots for it to remain an effective visualization
274 tool. Moreover, researchers could benefit from additional plotting tools that allow
275 them to quickly verify individual DEGs returned from a model. As a result, we de-
276 veloped a plot that allows users to visualize *one* DEG of interest onto the Cartesian
277 coordinates of *one* scatterplot matrix.

278 The “replicate line plot” was developed by researchers who demonstrated it could
279 detect model scaling problems in microarray data [31]. Unfortunately, this plot is
280 only applicable on datasets where treatment groups contain exactly two replicates.
281 The plot we now introduce is an extension of the “replicate line plot” that can be
282 applied to datasets with two or more replicates. We call this new plot a replIcate
283 TREatment (“litre”) plot.

284 In the litre plot, each gene is plotted once for each possible combination of repli-
285 cates between treatment groups. For example, there are nine ways to pair a replicate
286 from one treatment group with a replicate from the other treatment group in the
287 soybean iron-metabolism dataset (N.1 and P.1, N.1 and P.2, N.1 and P.3, N.2 and
288 P.1, N.2. and P.2, N.2 and P.3, N.3 and P.1, N.3 and P.2, and N.3 and P.3) [24].
289 Hence, each gene in this dataset is plotted as nine points in the litre plot. With
290 56,044 genes in this data, we would need to plot 504,396 points. This would reduce
291 the speed of interactive functionality as well as cause overplotting problems. As a
292 result, we again use hexagon bins to summarize this massive information (Figure 10
293 shows eight example litre plots).

294 Once the background of hexagons has been drawn to give us a sense of the distri-
295 bution of all between-treatment sample pair combinations for all genes, the user can
296 superimpose the nine points of one gene of interest. We can examine and compare
297 litre plots using the clusters we created in Figure 2. Subplots A and B of Figure 10
298 each show a significant gene from Cluster 1 plotted as nine green points, subplots
299 C and D each show a significant gene from Cluster 2 plotted as nine dark yellow

points, subplots E and F each show a significant gene from Cluster 3 plotted as nine pink points, and subplots G and H each show a significant genes from Cluster 4 plotted as nine orange points.

For the case of Figures 10 A and B, the nine overlaid points are superimposed in a manner we would expect from a DEG: They are located far from the $x=y$ line (difference between treatments) and are close to each other (similarity between replicates). Figures 10 C and D also show expected patterns for DEGs, although the genes are now overexpressed in the other treatment (group N). The replicates in subplot D are so precise that the overlaid points almost entirely overlap each other. In contrast, Figures 10 E and F do not seem to show as much replicate consistency. Now, there seems to be a pattern in which one replicate from the P group is larger than (and visually distanced from) the other two replicates. In other words, litre plots are able to capture the pattern differences in the significant genes from Cluster 2 and 3 that we saw back in Figure 2.

Moreover, in the case of Figures 10 G and H, the nine overlaid points are not clearly superimposed in the distinct pattern we expect of significant genes. While subplot G shows a gene that has consistent replications, the difference between the treatment groups is so small that the overlaid points cluster around the $x=y$ line. Additionally, the gene displayed in subplot H shows inconsistent replications and consistent treatment groups, as the spread-out overlaid points center on the $x=y$ line. Despite these genes being deemed significant by the model, the litre plots call into question whether the genes from this cluster show an expected profile of differential expression. This is similar to the messy-looking parallel coordinate plots we saw from these genes in Cluster 4 back in Figure 2. As a result, litre plots can detect odd and questionable patterns in individual “significant genes” that cannot be detected numerically through models. If this happens, the user may wish to further investigate these DEG calls.

Interactive litre plots are available online for the Cluster 1 significant genes (Figure 10 A and B) [32], Cluster 2 significant genes (Figure 10 C and D) [33], Cluster 3 significant genes (Figure 10 E and F) [34], and Cluster 4 significant genes (Figure 10 G and H) [35]. As can be verified in the interactive versions of the litre plot, users are provided several input fields that tailor the plot functionality. For instance, the user can easily select which treatment pair to explore (for data that contains more than two treatment groups) and can quickly scroll through significant genes one by one in order of increasing FDR values. Please read the “About” tab in the interactive links for more information.

Corresponding scatterplot matrices with the DEGs from these four clusters overlaid can be viewed in Figures 11- 14. Readers can verify that the parallel coordinate plots, litre plots, and scatterplot matrices tell a similar story about the DEG patterns in these four clusters.

Closing case study

We briefly discuss an additional example that merges many of the topics addressed in this paper. The publicly available data for this example contain technical replicates of liver and kidney RNA samples from one human male [12]. We first calculate DEG calls for this data using the normalization method of library size scaling, where the number of total reads in each sample are normalized to a common value across all samples. This process leads to 9,018 DEGs, with most of them ($\sim 78\%$) showing higher expression in the kidney group.

Although we could finish our analysis at this point and draw conclusions based on this list of DEGs that came from the model, it would be wise to also visualize this dataset. Viewing this data as a scatterplot matrix confirms the expected pattern with treatment scatterplots showing larger variation than technical replicate scatterplots (Figure 15). However, it also uncovers a hidden pattern in the treatment plots: There is a pronounced streak of genes with higher expressions in the

354 liver group (highlighted with a blue oval in Figure 15). We should also view the
355 DEGs from the model using parallel coordinate plots: Upon doing so, we notice
356 that while the 1,968 liver-specific DEGs follow the expected pattern of significant
357 calls, a substantial fraction of the 7,050 kidney-specific DEGs appear comparatively
358 noisy (Figure 16A).

359 Taking both of these observations into account, we may need to reconsider our
360 normalization technique. Some authors have argued that library size scaling method
361 is not adequate in all cases, especially when the underlying distribution of reads be-
362 tween samples is inconsistent. In the current data, the observed streak of outlier
363 genes that are highly expressed in the liver samples (Figure 15) reduces the sequenc-
364 ing quota available to the remaning genes in these samples, which could create an
365 articial inflation of the kidney-specific DEG calls. These authors have recommended
366 TMM normalization for such cases (including for this particular dataset) as this
367 technique generates sample scaling factors that consider sample distributions [15].

368 In light of all this, we re-start the analysis and now apply TMM normalization
369 to this data. This process leads to 7,520 DEGs that have a more level distribution
370 between the kidney ($\sim 53\%$) and liver ($\sim 47\%$) groups. The scatterplot matrix did
371 not appear differently from what we saw in Figure 15 as both of these normalization
372 methods are scaling procedures. However, we should visualize the new DEG calls.
373 Plotting these DEGs as parallel coordinate lines paints a much cleaner picture from
374 what we saw earlier, with most genes following the expected pattern of significance
375 (Figure 16B). Of the 7,050 kidney-specific DEGs we saw previously with library
376 size scaling normalization, only a much cleaner-looking subset ($n=3,974$) of them
377 remained as such using TMM normalization. TMM normalization kept the original
378 1,968 liver-specific DEGs from library size scaling but added 1,578 more for a total
379 of 3,546 liver-specific DEGs. As such, it appears that the liver-specific DEGs may be
380 slightly less clean-looking with TMM normalization. We emphasize that the 3,974

381 kidney-specific DEGs from TMM normalization are a proper subset of the 7,050
382 kidney-specific DEGs from library scale normalization, and the 1,968 liver-specific
383 DEGs from library scale normalization are a proper subset of the 3,546 liver-specific
384 DEGs from TMM normalization.

385 We therefore perform a deeper investigation of the effects of normalization on this
386 data. To do this, we thoroughly explore four subsets of genes from this case study in
387 the form of parallel coordinate plots, scatterplot matrices, and litre plots. We also
388 demonstrate the use of data standardization for scatterplot matrices and litre plots
389 as a means to magnify certain informative patterns. In this thorough examination,
390 we will use consistent color-coding when plotting example genes from each of the
391 four gene subsets. The four gene subsets and their color-codes are as follows:

392 1 The 3,974 kidney-specific DEGs from library size scale normalization that
393 remained as DEGs even after TMM normalization. These DEGs will be plot-
394 ted in purple. As these genes were declared significant with both library size
395 scale normalization and TMM normalization, we expect them to follow the
396 expected patterns of DEGs.

397 2 The 1,968 liver-specific DEGs from library size scale normalization that re-
398 mained as DEGs even after TMM normalization. These DEGs will be plotted
399 in orange. As these genes were declared significant with both library size scale
400 normalization and TMM normalization, we expect them to follow the expected
401 patterns of DEGs.

402 3 The 3,076 kidney-specific DEGs from library size scale normalization that were
403 *removed* as DEGs using TMM normalization. These DEGs will be plotted
404 in red. As these genes were removed from DEG designation with the more-
405 appropriate TMM normalization, we expect them to *not* convincingly follow
406 the expected patterns of DEGs.

407 4 The 1,578 liver-specific genes that were not detected as DEGs with library
408 size scale normalization but were then *added* as such using TMM normaliza-
409 tion. These DEGs will be plotted in pink. As these genes were not declared
410 significant with library size scale normalization but were then declared as sig-
411 nificant using the more-appropriate TMM normalization, we expect them to
412 *somewhat* convincingly follow the expected patterns of DEGs.

413 We begin by plotting the four gene subsets in the form of parallel coordinate
414 plots after application of hierarchical clustering analysis (Figures 17 through 20).
415 Each subset is grouped into eight clusters, not only to separate the genes into any
416 subtle pattern differences, but also to reduce any overplotting that would occur
417 should they all be viewed together as one large cluster. Figures 17 and 18 show
418 that the genes designated as DEGs in both forms of normalization (purple and
419 orange) have clean-looking patterns (especially in their largest cluster), Figure 19
420 shows that the genes removed with TMM normalization (red) have messy-looking
421 parallel coordinate plots, and Figure 20 shows that the genes added with TMM
422 normalization (pink) have parallel coordinate plots that are less clean than those
423 in Figures 17 and 18 but more clean than those in Figure 19.

424 We continue our visualization study by overlaying genes from the largest cluster
425 of the four gene subsets in the form of *standardized* scatterplot matrices (Figures 21
426 through 24). Notice that standardization causes the whole dataset to appear as
427 oval-shapes that are almost identical across all scatterplots. In other words, when
428 we standardize our scatterplot matrices, we lose geometric structures that can elicit
429 meaningful information about the dataset as a whole like we saw in Figures 3, 4,
430 5, 7, 9, and 15. However, in compensation for losing useful information about the
431 whole dataset, standardization often amplifies meaningful patterns in the overlaid
432 DEGs. Should the reader be interested, Additional files 1-4 show the same scat-
433 terplot matrices as the current case study below (Figures 21 through 24) only not

434 standardized. The reader can verify that the overlaid DEG patterns are more spread
435 out in the standardized version, allowing for better interpretation.

436 In general, we see that the genes that were called DEGs in both forms of nor-
437 malization (purple and orange) have the expected differential expression profiles in
438 the standardized scatterplot matrices, deviating from the $x=y$ line in the treatment
439 scatterplots in the anticipated direction (Figures 21 and 22). The standardized red
440 gene profiles show widely dispersed genes that sometimes deviate from the $x=y$ line
441 in the replicate scatterplots and cross both sides of and sometimes stick to the $x=y$
442 line in the treatment scatterplots (Figure 23). In other words, the red gene profiles
443 often show patterns not akin to differential expression, which we would expect from
444 genes that were *removed* as DEGs with TMM normalization. In contrast, the stan-
445 dardized pink gene profiles show less-widely dispersed genes that deviate less from
446 the $x=y$ line in the replicate scatterplots and deviate more from the $x=y$ line in
447 the treatment scatterplots (Figure 24). In other words, the pink gene profiles show
448 patterns more akin to differential expression than the red genes, which we would
449 expect from genes that were *added* as DEGs with TMM normalization. At the same
450 time, the pink gene profiles are not as clean-looking as the purple and orange genes
451 that were designated as DEGs in both forms of normalization. Overall, in these
452 standardized scatterplot matrices, the pink genes appear as an intermediate be-
453 tween the clean-looking purple and orange genes and the messy-looking red genes,
454 which we might expect.

455 We end our investigation by overlaying example genes from the largest cluster of
456 the four gene subsets in the form of *standardized* litre plots (Figures 25 through 28).
457 Similar to what we saw earlier, standardization causes the dataset to appear as an
458 oval-shape and removes the original geometric structure in the hexagonal binning.
459 Should the reader be interested, Additional files 5-8 show the same litre plots as the
460 current case study below (Figures 25 through 28) only not standardized. The reader

461 can verify that the overlaid DEG patterns are more spread out in the standardized
462 version in the current case study below, allowing for better interpretation.

463 Overall, we see that the example genes that were called DEGs in both forms
464 of normalization (purple and orange) have the expected profiles in the litre plots,
465 deviating as concentrated bundles away from the $x=y$ line (Figures 25 and 26). The
466 standardized litre plots for the nine genes with the lowest FDR values for both
467 the red (Figure 27) and pink (Figure 28) groups allow us to quickly determine
468 that the pink profiles show patterns more akin to differential expression than the
469 red groups. Namely, the overlaid pink points deviate more from the $x=y$ line in
470 a tight cluster than the overlaid red points. At the same time, the overlaid pink
471 points show patterns less akin to differential expression than the purple and orange
472 points. All together, the pink gene profiles again appear as intermediates between
473 the clean-looking purple and orange genes and the messy-looking red genes in the
474 standardized litre plots, which is to be expected if TMM normalization is the more
475 appropriate technique.

476 Our in-depth analyses in this case study collectively suggest that this dataset
477 indeed requires more than just library size scaling for reliable analysis. This case
478 study was meant to underscore the overarching theme of this paper that iteration
479 between models and visualizations is crucial to achieve the most convincing results
480 and conclusions in RNA-seq studies.

481 **Plot scalability**

482 All visualization plots discussed in this paper have limitations based on the number
483 of samples in the data. Plots that appear messy, regardless of sample numbers,
484 indicate the presence of data quality problems. In general, MDS plots, boxplots,
485 and parallel coordinate plots can remain effective with fairly large sample numbers.
486 We note that parallel coordinate plots should be sorted with logic, especially when
487 scaling to larger data sets.

Scatterplot matrices usually lose their efficiency at smaller sample numbers due to restricted space: n^2 scatterplots must be drawn for n -dimensional data. One remedy is for users to subset their data and plot several smaller scatterplot matrices. We used this technique in our recent honey bee RNA-seq paper where we investigated 2 groups of 12 replicates (24 samples total) [36]. Plotting all samples onto one scatterplot matrix would have required a prohibitive $24 \times 24 = 576$ scatterplots. Instead, we divided the data into four subsets, each with 2 groups of 3 replicates (6 samples total) so that each scatterplot matrix only required $6 \times 6 = 36$ scatterplots. See additional files 11-14 of [36].

Litre plots are another remedy for large data sets and can often accommodate more samples than scatterplot matrices. Indeed, in our honey bee RNA-seq paper, we successfully applied litre plots to our full data that contained 2 groups of 12 replicates (24 samples total). In cases where there are two treatment groups with an equal number of replicates, the litre plot draws n^2 number of points, where n is the number of replicates in each group. Hence, we were able to draw $12 \times 12 = 144$ dots on the litre plot successfully in our previous paper. See additional file 4 of [36]. We note that the litre plot is more suitable for large datasets than the replicate line plot, which is ideal for 2 groups of 2 replicates (4 samples total).

Discussion

In this paper, we strived to convince readers that effective visualization should be a crucial part of two-group differential expression analysis. We used real data to demonstrate that scatterplot matrices, parallel coordinate plots, and litre plots help users check for normalization problems, catch common errors in analysis pipelines, and confirm that the variation between replicates and treatments is as expected. We also showed that these graphical tools allow researchers to quickly explore DEG lists that come out of models and ensure which ones make sense from an additional and arguably more intuitive vantage point. Moreover, we demonstrated that our

515 plotting tools allow researchers to discover genes of interest through visual geometric
516 patterns that would otherwise remain undiscovered with models.

517 In general, scientists might uncover surprising patterns lurking in their data with
518 plots in ways that cannot be achieved with any formulas or models. Researchers
519 from all statistical backgrounds can use graphical tools to better understand (if not
520 demystify) how the application of various normalization techniques and/or models
521 affect their results. All in all, scientists can gain more confidence in the data analysis
522 pipelines they choose and in the results they draw at the mere cost of briefly creating
523 and exploring graphical outputs during their analyses.

524 Conclusions

525 Modern data analysis is most reliable when models and visuals are used congru-
526 ently. Unfortunately, there is a propensity for researchers to overtrust model results
527 without confirming them with graphics. This, as we have shown, calls into ques-
528 tion the soundness of results derived from differential expression studies. Solving
529 this problem is straightforward and does not require scientists to drastically change
530 their differential expression analyses. Instead, scientists simply need to incorporate
531 effective plotting tools during their usual analysis pipelines. The main motivation
532 of this paper was to provide a collection of examples that show *why* visualization
533 tools should be an integral part of two-group differential expression analysis. We
534 hope our work may motivate researchers to take into account plotting tools that are
535 conveniently and freely available for differential expression analysis. We also hope
536 our work may influence developers to create additional RNA-seq plotting tools that
537 can be applied outside of the case of two-group differential expression. This could
538 include plotting tools for cases that contain many groups, cases of single-cell anal-
539 ysis, and cases where researchers are looking for specificity rather than differential
540 expression.

We strive to serve another small role in this solution with our R software package called “bigPint” that includes the plotting techniques introduced in this paper, many of which are unique additions to the array of plotting tools currently available in differential expression analysis packages. The “bigPint” website is available online [37]. To encourage scientists to use our resource, we include short vignette articles on our website that introduce users to our package. One article provides a recommended pipeline for iterating between models and visualizations when performing differential expression analysis [38]. There is a need to make it easier for researchers to use models and visuals in a complimentary fashion when analyzing RNA-seq data. Our software incorporates data structures that allow users to transition smoothly between our plots and popular models from packages like edgeR [9], DESeq2 [21], and limma [23]. We demonstrate the ease of transition between models and visualizations in the articles of our website. All articles are written using reproducible code that new users can follow. It is our hope that such work will serve a small part in upgrading the RNA-seq analysis world into one that more holistically extracts meaningful biological information using both models and visuals.

Methods

Four public RNA-sequencing datasets were studied in this paper. R [39] was used to conduct analyses. Packages htmlwidgets [40], ggplot2 [41], shiny [42], and plotly [43] were used to build the graphics. The pkgdown [44] package was used to construct the “bigPint” software webpage. Our interactive applications were deployed on shinyapps.io [45].

Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Availability of data and materials

The bigPint methods used in this paper are available on the software website [37]. The bigPint package is available on Bioconductor at <https://bioconductor.org/packages/devel/bioc/html/bigPint.html>. The four public datasets discussed in this publication are available online: Three are deposited on the NCBI Sequence Read Archive with

572 accession numbers SRA000299 [12], PRJNA318409 [24], and SRA048710 [22]. One is deposited on the NCBI Gene
 573 Expression Omnibus with accession number GSE61857 [27]. Reproducible scripts for all figures in this manuscript
 574 are available online at <https://github.com/lindsayrutter/VisualizationMethods>.

575 Competing interests

576 The authors declare that they have no competing interests.

577 Funding

578 Graham and Moran Lauter were financed by the United States Department of Agriculture, Agricultural Research
 579 Service (USDA-ARS) CRIS Project 5030-21220-005-00D and the Iowa Soybean Association.

580 Author's contributions

581 LR produced the visualization tools and wrote the manuscript. LR and DC analyzed the four case study datasets.
 582 AML and MAG planned, designed, and conducted the experiment to identify iron-stress responsive genes in soybean
 583 and contributed to its analysis. All authors revised and approved the final manuscript.

584 Acknowledgements

585 We would like to thank Dr. Amy L. Toth for her helpful discussions and thorough feedback for this project.

586 Author details

587 ¹Bioinformatics and Computational Biology Program, Iowa State University, Ames, USA. ²USDA-ARS, Corn
 588 Insects and Crop Genetics Research Unit, Ames, USA. ³Department of Agronomy, Iowa State University, Ames,
 589 USA. ⁴Econometrics and Business Statistics, Monash University, Clayton VIC, Australia.

590 References

- 591 1. Robertson, G., Schein, J., Chiu, R., Corbett, R., Field, M., Jackman, S.D., Mungall, K., Lee, S., Okada, H.M.,
 592 Qian, J.Q., *et al.*: De novo assembly and analysis of rna-seq data. *Nature methods* **7**(11), 909 (2010)
- 593 2. Anders, S., Reyes, A., Huber, W.: Detecting differential usage of exons from rna-seq data. *Genome research*,
 594 133744 (2012)
- 595 3. Pan, Q., Shai, O., Lee, L.J., Frey, B.J., Blencowe, B.J.: Deep surveying of alternative splicing complexity in the
 596 human transcriptome by high-throughput sequencing. *Nature genetics* **40**(12), 1413 (2008)
- 597 4. Wang, Z., Gerstein, M., Snyder, M.: Rna-seq: a revolutionary tool for transcriptomics. *Nature reviews genetics*
 598 **10**(1), 57 (2009)
- 599 5. Zhao, S., Fung-Leung, W.-P., Bittner, A., Ngo, K., Liu, X.: Comparison of rna-seq and microarray in
 600 transcriptome profiling of activated t cells. *PLoS one* **9**(1), 78644 (2014)
- 601 6. Anders, S., Huber, W.: Differential expression analysis for sequence count data. *Genome biology* **11**(10), 106
 602 (2010)
- 603 7. Trapnell, C., Roberts, A., Goff, L., Pertea, G., Kim, D., Kelley, D.R., Pimentel, H., Salzberg, S.L., Rinn, J.L.,
 604 Pachter, L.: Differential gene and transcript expression analysis of rna-seq experiments with tophat and
 605 cufflinks. *Nature protocols* **7**(3), 562 (2012)
- 606 8. Trapnell, C., Hendrickson, D.G., Sauvageau, M., Goff, L., Rinn, J.L., Pachter, L.: Differential analysis of gene
 607 regulation at transcript resolution with rna-seq. *Nature biotechnology* **31**(1), 46 (2013)
- 608 9. Robinson, M.D., McCarthy, D.J., Smyth, G.K.: edgeR: a bioconductor package for differential expression
 609 analysis of digital gene expression data. *Bioinformatics* **26**(1), 139–140 (2010)
- 610 10. Law, C.W., Chen, Y., Shi, W., Smyth, G.K.: voom: Precision weights unlock linear model analysis tools for
 611 rna-seq read counts. *Genome biology* **15**(2), 29 (2014)
- 612 11. Morin, R.D., Bainbridge, M., Fejes, A., Hirst, M., Krzywinski, M., Pugh, T.J., McDonald, H., Varhol, R.,
 613 Jones, S.J., Marra, M.A.: Profiling the hela s3 transcriptome using randomly primed cdna and massively parallel
 614 short-read sequencing. *Biotechniques* **45**(1), 81–94 (2008)
- 615 12. Marioni, J.C., Mason, C.E., Mane, S.M., Stephens, M., Gilad, Y.: Rna-seq: an assessment of technical
 616 reproducibility and comparison with gene expression arrays. *Genome research* (2008)
- 617 13. Hansen, K.D., Brenner, S.E., Dudoit, S.: Biases in illumina transcriptome sequencing caused by random
 618 hexamer priming. *Nucleic acids research* **38**(12), 131–131 (2010)
- 619 14. Oshlack, A., Robinson, M.D., Young, M.D.: From rna-seq reads to differential expression results. *Genome*
 620 *biology* **11**(12), 220 (2010)
- 621 15. Robinson, M.D., Oshlack, A.: A scaling normalization method for differential expression analysis of rna-seq
 622 data. *Genome biology* **11**(3), 25 (2010)
- 623 16. McIntyre, L.M., Lopiano, K.K., Morse, A.M., Amin, V., Oberg, A.L., Young, L.J., Nuzhdin, S.V.: Rna-seq:
 624 technical variability and sampling. *BMC genomics* **12**(1), 293 (2011)
- 625 17. Bullard, J.H., Purdom, E., Hansen, K.D., Dudoit, S.: Evaluation of statistical methods for normalization and
 626 differential expression in mrna-seq experiments. *BMC bioinformatics* **11**(1), 94 (2010)
- 627 18. Shneiderman, B.: Inventing discovery tools: combining information visualization with data mining. *Information*
 628 *visualization* **1**(1), 5–12 (2002)
- 629 19. Gentleman, R.C., Carey, V.J., Bates, D.M., Bolstad, B., Dettling, M., Dudoit, S., Ellis, B., Gautier, L., Ge, Y.,
 630 Gentry, J., *et al.*: Bioconductor: open software development for computational biology and bioinformatics.
 631 *Genome biology* **5**(10), 80 (2004)
- 632 20. Huber, W., Carey, V.J., Gentleman, R., Anders, S., Carlson, M., Carvalho, B.S., Bravo, H.C., Davis, S., Gatto,
 633 L., Girke, T., *et al.*: Orchestrating high-throughput genomic analysis with bioconductor. *Nature methods* **12**(2),
 634 115 (2015)
- 635 21. Love, M.I., Huber, W., Anders, S.: Moderated estimation of fold change and dispersion for rna-seq data with
 636 deseq2. *Genome biology* **15**(12), 550 (2014)
- 637 22. Risso, D., Schwartz, K., Sherlock, G., Dudoit, S.: Gc-content normalization for rna-seq data. *BMC*
 638 *bioinformatics* **12**(1), 480 (2011)

- 639 23. Ritchie, M.E., Phipson, B., Wu, D., Hu, Y., Law, C.W., Shi, W., Smyth, G.K.: limma powers differential
640 expression analyses for rna-sequencing and microarray studies. *Nucleic acids research* **43**(7), 47–47 (2015)
- 641 24. Moran Lauter, A.N., Graham, M.A.: NCBI SRA Bioproject Accession: PRJNA318409.
642 <https://www.ncbi.nlm.nih.gov/bioproject/PRJNA318409> Accessed 2018-12-20
- 643 25. Chandrasekhar, T., Thangavel, K., Elayaraja, E.: Effective clustering algorithms for gene expression data. *arXiv*
644 preprint arXiv:1201.4914 (2012)
- 645 26. de Souto, M.C., de Araujo, D.S., Costa, I.G., Soares, R.G., Luderemir, T.B., Schliep, A.: Comparative study on
646 normalization procedures for cluster analysis of gene expression datasets. In: *Neural Networks, 2008. IJCNN*
647 *2008. (IEEE World Congress on Computational Intelligence). IEEE International Joint Conference On*, pp.
648 2792–2798 (2008). IEEE
- 649 27. Brown, A.V., Hudson, K.A.: Developmental profiling of gene expression in soybean trifoliate leaves and
650 cotyledons. *BMC plant biology* **15**(1), 169 (2015)
- 651 28. Rutter, L., Cook, D. <https://rnaseqvisualization.shinyapps.io/scatmat> Accessed 2018-12-20
- 652 29. Ioannidis, J.P., Allison, D.B., Ball, C.A., Coulibaly, I., Cui, X., Culhane, A.C., Falchi, M., Furlanello, C., Game,
653 L., Jurman, G., *et al.*: Repeatability of published microarray gene expression analyses. *Nature genetics* **41**(2),
654 149 (2009)
- 655 30. Baggerly, K.A., Coombes, K.R.: Deriving chemosensitivity from cell lines: Forensic bioinformatics and
656 reproducible research in high-throughput biology. *The Annals of Applied Statistics*, 1309–1334 (2009)
- 657 31. Cook, D., Hofmann, H., Lee, E.-K., Yang, H., Nikolau, B., Wurtele, E.: Exploring gene expression data, using
658 plots. *Journal of Data Science* **5**(2), 151 (2007)
- 659 32. Rutter, L., Cook, D. <https://rnaseqvisualization.shinyapps.io/litrecluster1> Accessed 2018-12-20
- 660 33. Rutter, L., Cook, D. <https://rnaseqvisualization.shinyapps.io/litrecluster2> Accessed 2018-12-20
- 661 34. Rutter, L., Cook, D. <https://rnaseqvisualization.shinyapps.io/litrecluster3> Accessed 2018-12-20
- 662 35. Rutter, L., Cook, D. <https://rnaseqvisualization.shinyapps.io/litrecluster4> Accessed 2018-12-20
- 663 36. Rutter, L., Carrillo-Tripp, J., Bonning, B.C., Cook, D., Toth, A.L., Dolezal, A.G.: Transcriptomic responses to
664 diet quality and viral infection in *apis mellifera*. *BMC genomics* **20**(1), 412 (2019)
- 665 37. Rutter, L., Cook, D.: bigPint: Make Big Data Pint-sized. <https://lindsayrutter.github.io/bigPint>
666 Accessed 2018-12-20
- 667 38. Rutter, L., Cook, D.: Recommended RNA-seq Pipeline.
668 <https://lindsayrutter.github.io/bigPint/articles/pipeline.html> Accessed 2018-12-20
- 669 39. Team, R.C.: A Language and Environment for Statistical Computing. <https://www.R-project.org> Accessed
670 2018-12-20
- 671 40. Vaidyanathan, R., Xie, Y., Allaire, J., Cheng, J., Russell, K.: *Htmlwidgets: HTML Widgets for R*, 2016.
672 <https://cran.r-project.org/package=htmlwidgets> Accessed 2018-12-20
- 673 41. Wickham, H.: *Ggplot2: Elegant Graphics for Data Analysis*. Springer, ??? (2016)
- 674 42. Chang, W., Cheng, J., Allaire, J., Xie, Y., McPherson, J.: Shiny: Web Application Framework for R [Computer
675 Software]. <https://cran.r-project.org/package=shiny> Accessed 2018-12-20
- 676 43. Sievert, C., Parmer, C., Hocking, T., Chamberlain, S., Ram, K., Corvellec, M., Despouy, P.: Plotly: Create
677 Interactive Web Graphics Via 'plotly.js'. <https://cran.r-project.org/package=plotly> Accessed 2018-12-20
- 678 44. Wickham, H., Hesselberth, J.: Pkgdown: Make Static HTML Documentation for a Package.
679 <https://cran.r-project.org/package=pkgdown> Accessed 2018-12-20
- 680 45. RStudio: Integrated Development for R. <http://www.rstudio.com> Accessed 2018-12-20

681 **Figures**

Figure 1 Comparison of plotting methods using simulated data. One simulated dataset is shown on the left half and another simulated dataset is shown on the right half of the figure. The parallel coordinate plots (subplots C) show a critical difference at the gene-level between the datasets. Namely, the left dataset is composed of genes with small replicate variation and large treatment group variation (suggesting DEGs), while the right dataset is composed of genes with similar variation between replicates and treatment groups (not suggesting DEGs). We cannot see this gene-level difference with the boxplots and MDS plots.

Figure 2 Parallel coordinate plots of clustered significant genes in the soybean iron metabolism data. Parallel coordinate plots of significant genes after hierarchical clustering of the soybean iron metabolism data [24]. We can quickly confirm that Clusters 1 and 2 show the typical pattern for significant genes. Cluster 4 does not distinctively show the usual profile for significant genes. Cluster 3 looks similar to Cluster 2, except for unexpectedly large P.3 values.

Figure 3 Expected structure of RNA-seq data plotted as a scatterplot matrix. Example of the expected structure of an RNA-seq dataset, using soybean cotyledon data from [27]. Within a given scatterplot, most genes (points) should fall along the $x=y$ line. We should see genes deviate more strongly from the $x=y$ line in treatment scatterplots (the nine scatterplots enclosed in the blue square) than in replicate scatterplots (the remaining six scatterplots).

Figure 4 Assessing normalization of RNA-seq data using scatterplot matrices. Illustrating normalization checks with data from [22]. The collective deviation of genes from the $x=y$ line instantly reveals that the RNA-seq dataset was not thoroughly normalized using within-lane normalization (subplot A). However, within-lane normalization followed by between-lane normalization sufficiently normalized the data (subplot B). The authors who developed these normalization methods showed that the later approach generated a lower false-positive DEG call rate in this dataset.

Figure 5 Checking common errors of RNA-seq data analysis using scatterplot matrices. As expected, the scatterplot matrix of the coytedon data [27] contains nine scatterplots with thicker distributions (should be treatment pairs) and six scatterplots with thinner distributions (should be replicate pairs). However, a subset of scatterplots unexpectedly show thicker distributions between replicate pairs and thinner distributions between treatment pairs. If we switch the labels of two suspicious samples (S1.3 and S2.1), the scatterplot matrix displays the anticipated structure we saw in Figure 3. At this point, we have evidence that these two samples may have been mislabeled, and we may wish to confirm this suspicion and correct it before continuing with the analysis.

Figure 6 Checking common errors of RNA-seq data analysis using side-by-side boxplots and MDS plots. Side-by-side boxplots and MDS plots are popular plotting tools for RNA-seq analysis. This figure shows these traditional visualization methods applied to the soybean cotyledon data before sample switching (left half) and after sample switching (right half) [27]. We cannot suspect from the right boxplot that samples S1.3 and S2.1 have been swapped (subplots A). This is because all six samples have similar five number summaries. For the MDS plots, we do see a cleaner separation of the two treatment groups across the first dimension in the left plot than in the right plot (subplots B). However, taking into account the second dimension, both MDS plots contain three clusters, with sample S1.1 appearing in its own cluster. Without seeing one distinct cluster for each of the two treatment groups, it is difficult to suspect that samples S1.3 and S2.1 have been swapped in the right MDS plot (subplots B). We can only derive clear suspicion that the samples may have been switched by using less-popular plots that provide gene-level resolution like with the scatterplot matrix from Figure 5.

Figure 7 Finding unexpected patterns in RNA-seq data using scatterplot matrices. Scatterplot matrix of RNA-seq read counts from soybean leaves after exposure to iron-sufficient (treatment group P) and iron-deficient (treatment group N) hydroponic conditions [24]. We observe the expected structure of treatment pairs showing larger variability around the $x=y$ line than replicate pairs. However, we notice a pronounced streak structure in the bottom-right scatterplot (green arrow) that compares two replicate samples from the iron-sufficient group. The genes in the streak structure have large read counts that deviate in a parallel fashion from the $x=y$ line.

Figure 8 Gene expression responses across time points. The authors of the soybean iron metabolism study [24] determined the DEGs across three time points (30 minutes, 60 minutes, and 120 minutes) in the leaves after onset of iron sufficient and deficient hydroponic conditions. They used the same researcher to collect the samples in succession. One major finding from their study was a vast change in gene expression responses between these three time points. As a result, the streak observed in the scatterplot matrix containing the subset of data at the 120 minute time point (Figure 7) may be due to the timing differences between replicate handling.

Figure 9 Assessing differential expression in RNA-seq data using scatterplot matrices. Example of the expected structure of DEG calls (in orange) from the soybean cotyledon dataset [27]. In the scatterplot matrix (subplot A), DEGs should fall along the $x=y$ line for replicates and deviate from it for treatments. In the parallel coordinate plot (subplot B), DEGs should show levelness between replicates and crosses between treatments. These two plotting types can be linked to quickly provide users multiple perspectives of their DEG calls.

Figure 10 Example litre plots for clustered significant genes in the soybean iron metabolism data. Litre plots for representative genes from clusters created in Figure 2 [24]. Subplots A and B each show a gene from Cluster 1 overlaid as green points. Subplots C and D each show a gene from Cluster 2 overlaid as dark yellow points. Subplots E and F each show a gene from Cluster 3 overlaid as pink points. Subplots G and H each show a gene from Cluster 4 overlaid as orange points.

Figure 11 Cluster 1 significant genes from the soybean iron metabolism data overlaid on a scatterplot matrix. Example of using a scatterplot matrix to assess DEG calls from a model in the iron-metabolism soybean dataset. There were 2,751 significant genes in Cluster 1 after performing a hierarchical clustering analysis with a cluster size of four (Figure 2). These significant genes are overlaid in green on the scatterplot matrix. They follow the expected patterns of differential expression with most green points falling along the $x=y$ line in the scatterplots between replicates, but deviating from the $x=y$ line in the scatterplots between treatments. The deviation consistently demonstrates higher expression in the P group than in the N group. Hence, these green points seem to represent DEGs that were significantly overexpressed in the P group, which draws the same conclusion with what we derived using the parallel coordinate plots in Figure 2. One difficulty with plotting such a large number of DEGs onto the scatterplot matrix is that overplotting can obscure our inability to determine how many DEGs are in a given location. This is why we should also view these genes individually in litre plots (Figure 10 A and B).

Figure 12 Cluster 2 significant genes from the soybean iron metabolism data overlaid on a scatterplot matrix. Example of using a scatterplot matrix to assess DEG calls from a model in the iron-metabolism soybean dataset. There were 2,009 significant genes in Cluster 2 after performing a hierarchical clustering analysis with a cluster size of four (Figure 2). These significant genes are overlaid in dark yellow on the scatterplot matrix. They follow the expected patterns of differential expression with most dark yellow points falling along the $x=y$ line in the scatterplots between replicates, but deviating from the $x=y$ line in the scatterplots between treatments. The deviation consistently demonstrates higher expression in the N group than in the P group. Hence, these dark yellow points seem to represent genes that were significantly overexpressed in the N group, which draws the same conclusion with what we derived using the parallel coordinate plots in Figure 2. One difficulty with plotting such a large number of DEGs onto the scatterplot matrix is that overplotting can obscure our inability to determine how many DEGs are in a given location. This is why we might also view these genes individually in litre plots (Figure 10 C and D).

Figure 13 Cluster 3 significant genes from the soybean iron metabolism data overlaid on a scatterplot matrix. Example of using a scatterplot matrix to assess DEG calls from a model in the iron-metabolism soybean dataset. There were 861 significant genes in Cluster 3 after performing a hierarchical clustering analysis with a cluster size of four (Figure 2). These significant genes are overlaid in pink on the scatterplot matrix. For the most part, they follow the expected patterns of differential expression with pink points falling along the $x=y$ line in the scatterplots between replicates, but deviating from the $x=y$ line in the scatterplots between treatments. The deviation consistently demonstrates higher expression in the N group than in the P group. The scatterplot between replicates P.1 and P.3 shows slightly higher expression in P.3, and the scatterplot between replicates P.2 and P.3 also shows slightly higher expression in P.3. Hence, these pink points seem to represent genes that were significantly overexpressed in the N group, but with slight inconsistencies in the replicates in the P group, which matches what we saw in the parallel coordinate plots in Figure 2. One difficulty with plotting such a large number of DEGs onto the scatterplot matrix is that overplotting can obscure our inability to determine how many DEGs are in a given location. This is why we might also view these genes individually in litre plots (Figure 10 E and F).

Figure 14 Cluster 4 significant genes from the soybean iron metabolism data overlaid on a scatterplot matrix. Example of using a scatterplot matrix to assess DEG calls from a model in the iron-metabolism soybean dataset. There were 17 significant genes in Cluster 4 after performing a hierarchical clustering analysis with a cluster size of four (Figure 2). These significant genes are overlaid in orange on the scatterplot matrix. For the most part, they do not seem to follow the expected patterns of differential expression: In many of the scatterplots between treatments, the orange points do not seem to deviate much from the $x=y$ line. Moreover, in the scatterplots between P.1 and P.2 as well as P.1 and P.3, the orange points seem to indicate an underexpression of the P.1 replicate. We similarly saw somewhat messy looking DEG calls in Cluster 4 in the form of parallel coordinate plots (Figure 2) and litre plots (Figure 10 G and H).

Figure 15 Scatterplot matrix detects unexpected structure in liver and kidney technical replicate RNA-sequencing dataset. Scatterplot matrix of liver and kidney technical replicates [12]. The technical replicate scatterplots look precise as is expected, with little variability around the $x=y$ line. The treatment group scatterplots have much more variability around the $x=y$ line, as we would expect. However, each treatment group scatterplot contains a pronounced streak of highly-expressed liver-specific genes, which deviates from the expected distribution (shown in blue oval in one example scatterplot). Some researchers have suggested that differences in the distribution of reads between groups may require particularly stringent normalization.

Figure 16 Comparing normalization method effect on significance designation using parallel coordinate plots. Subplot A shows parallel coordinate plots of the DEGs from liver and kidney technical replicates [12] after library size scale normalization. The division of DEGs between the two groups was rather disparate, with $\sim 78\%$ of the DEGs being kidney-specific and only $\sim 22\%$ of the DEGs being liver-specific. Also of note, while the parallel coordinate patterns of the liver-specific DEGs appear as expected, the patterns of the kidney-specific DEGs seem to show comparatively larger variability between the replicates. Subplot B shows parallel coordinate plots of the DEGs from liver and kidney technical replicates after TMM normalization. The division of DEGs between the two groups is more balanced than in Subplot A, with $\sim 53\%$ of the DEGs being kidney-specific and $\sim 47\%$ of the DEGs being liver-specific. Additionally, the parallel coordinate patterns of the kidney DEGs is vastly improved. However, the parallel coordinate patterns of the liver DEGs is slightly more messy looking. As a result, we investigate the effects of normalization on this data more thoroughly.

Figure 17 Parallel coordinate plots for gene clusters that remained as kidney-specific DEGs after TMM normalization. Parallel coordinate plots showing eight hierarchical clusters from the 3,974 genes that remained in the kidney-specific DEGs after TMM normalization. We see that, for the most part, the parallel coordinate patterns follow the expected patterns across the clusters. The ideal pattern of DEGs is especially captured in the first cluster (the largest one with 1,136 genes). We applied ombre coloring across the clusters in order of cluster size. We used hierarchical clustering to tease apart subtle pattern differences and to mitigate additional overplotting that would occur if we were to plot all genes onto only one parallel coordinate plot. The side-by-side boxplots represent *all* gene counts in the dataset.

Figure 18 Parallel coordinate plots for gene clusters that remained as liver-specific DEGs after TMM normalization. Parallel coordinate plots showing eight hierarchical clusters from the 1,968 genes that remained in the liver-specific DEGs after TMM normalization. We see that, for the most part, the parallel coordinate patterns follow the expected patterns across the clusters. The ideal pattern of DEGs is especially captured in the first cluster (the largest one with 933 genes). We applied ombre coloring across the clusters in order of cluster size. We used hierarchical clustering to tease apart subtle pattern differences and to mitigate additional overplotting that would occur if we were to plot all genes onto only one parallel coordinate plot. The side-by-side boxplots represent *all* gene counts in the dataset.

Figure 19 Parallel coordinate plots for gene clusters that were removed from kidney-specific DEGs after TMM normalization. Parallel coordinate plots showing eight hierarchical clusters from the 3,076 genes that were removed from the kidney-specific DEGs after TMM normalization. The patterns in almost all clusters do not resemble the expected DEG format; instead, they show large variability between replicates and small variability between treatments. This plot provides additional statistical evidence that the application of TMM normalization successfully removed genes that were previously mislabeled as kidney-specific DEGs with library size scaling normalization. We used hierarchical clustering to tease apart subtle pattern differences and to mitigate overplotting. The side-by-side boxplots represent *all* gene counts in the dataset.

Figure 20 Parallel coordinate plots for gene clusters that were added as liver-specific DEGs after TMM normalization. Parallel coordinate plots showing eight hierarchical clusters from the 1,578 genes that were added as liver-specific DEGs after TMM normalization. We see that the parallel coordinate lines *somewhat* follow the expected patterns across the clusters, better than what we saw in the red (Figure 19) gene subsets, but not as precisely as we saw with the purple (Figure 17) and orange (Figure 18) gene subsets. We applied ombre coloring across the clusters in order of cluster size. We used hierarchical clustering to tease apart subtle pattern differences and to mitigate additional overplotting that would occur if we were to plot all genes onto only one parallel coordinate plot. The side-by-side boxplots represent *all* gene counts in the dataset.

Figure 21 Standardized scatterplot matrix for gene cluster that remained as kidney-specific DEGs after TMM normalization. Scatterplot matrix of the *standardized* 1,136 genes that were in the first cluster (Figure 17) from genes that remained as kidney-specific DEGs even after TMM normalization. Even though the standardization process removes the interesting geometrical features we would otherwise see, it amplifies DEG patterns more clearly. Here, the highlighted genes appear more clustered and separated from the $x=y$ line in the treatment scatterplots, and more clustered and connected to the $x=y$ line in the replicate scatterplots. We can also now see more clearly in the replicate scatterplots that the kidney expression is higher than the liver expression.

Figure 22 Standardized scatterplot matrix for gene cluster that remained as liver-specific DEGs after TMM normalization. Scatterplot matrix of the *standardized* 933 genes that were in the first cluster (Figure 18) from genes that remained as liver-specific DEGs even after TMM normalization. Even though the standardization process removes the interesting geometrical features we would otherwise see, it amplifies DEG patterns in meaningful ways. Here, the highlighted genes appear more clustered and separated from the $x=y$ line in the treatment scatterplots, and more clustered and connected to the $x=y$ line in the replicate scatterplots. We can also now see more clearly in the replicate scatterplots that the liver expression is higher than the kidney expression.

Figure 23 Standardized scatterplot matrix for gene cluster that were removed from kidney-specific DEGs after TMM normalization. Scatterplot matrix of the *standardized* 529 genes that were in the first cluster (Figure 19) from genes that no longer remained as kidney-specific DEGs after TMM normalization. Even though the standardization process removes the interesting geometrical features we would otherwise see, it amplifies DEG patterns in meaningful ways. Namely, the genes of interest are now spread out more, and the replicate and treatment scatterplots are almost indistinguishable from each other, with both of them showing genes of interest crossing both sides of the $x=y$ line. In other words, standardization of the data provides clear visualization evidence that TMM normalization was justified in removing these genes from DEG designation.

Figure 24 Standardized scatterplot matrix for gene cluster that were added as liver-specific DEGs after TMM normalization. Scatterplot matrix of the *standardized* 317 genes that were in the first cluster (Figure 20) from genes that were *added* as liver-specific DEGs after TMM normalization. Even though the standardization process removes the interesting geometrical features we would otherwise see, it amplifies DEG patterns in meaningful ways. Namely, the genes of interest are now spread out more, and we can now distinguish the replicate and treatment scatterplots more clearly. For the most part, the genes of interest deviate from the $x=y$ line in the treatment scatterplots more so than in the replicate scatterplots, and hence display somewhat of the pattern of differential expression. In fact, the pink genes again appear as an intermediate between the purple and orange genes that clearly display differential expression (Figures 21 and 22) and the red genes that clearly do *not* display differential expression (Figure 23). In other words, standardized scatterplot matrices provide additional visualization evidence that TMM normalization was justified in removing the red genes from and adding the pink genes to DEG designation.

Figure 25 Example standardized litre plots for genes that remained as kidney-specific DEGs after TMM normalization. Example *standardized* litre plots from the 1,136 genes that were in the first cluster (Figure 17) of genes that remained as kidney-specific DEGs even after TMM normalization. With standardization, we immediately note that meaningful information about the dataset as a whole (variation between treatments and replicates, normalization, sample mislabeling, and unexpected patterns like streaks) is now gone. In any case, we confirm that these standardized litre plots corroborate that these purple genes demonstrate the expected patterns of DEGs.

Figure 26 Example standardized litre plots for genes that remained as liver-specific DEGs after TMM normalization. Example litre plots from the 933 genes that were in the first cluster (Figure 18) from genes that remained as liver-specific DEGs even after TMM normalization. With standardization, we immediately note that meaningful information about the dataset as a whole (variation between treatments and replicates, normalization, sample mislabeling, and unexpected patterns like streaks) is now gone. In any case, we confirm that these standardized litre plots corroborate that these orange genes demonstrate the expected patterns of DEGs.

Figure 27 Example standardized litre plots for genes that were removed from kidney-specific DEGs after TMM normalization. *Standardized* litre plots for the nine genes with the lowest FDR values out of the 529 genes that were in the first cluster (Figure 19) of genes that no longer remained as kidney-specific DEGs after TMM normalization. We verify from an additional perspective that the red genes do not demonstrate the expected patterns of DEGs. The example red genes here are show much larger inconsistencies between replicates than what we saw with the purple (Figure 25) and orange (Figure 26) genes. This provides additional evidence that TMM normalization removing these genes from DEG status may be valid.

Figure 28 Example standardized litre plots for genes that were added as liver-specific DEGs after TMM normalization. *Standardized* litre plots for the nine genes with the lowest FDR values out of the 317 genes that were in the first cluster (Figure 20) from genes that were *added* as liver-specific DEGs after TMM normalization. We can quickly determine that the pink profiles in this figure show patterns more akin to differential expression than the red profiles in Figure 27. That is, the overlaid pink points deviate more from the $x=y$ line in a tight cluster than the overlaid red points. At the same time, the overlaid pink points here show patterns less akin to differential expression than the purple (Figure 25) and orange (Figure 26) points. In sum, the standardized litre plots again place the pink gene profiles as an intermediate.

Additional Files

Additional file 1 — Scatterplot matrix for gene cluster that remained as kidney-specific DEGs after TMM normalization.

Scatterplot matrix of the 1,136 genes that were in the first cluster (of Figure 17) from genes that remained as kidney-specific DEGs even after TMM normalization. With this scatterplot matrix, we verify from an additional perspective that these genes demonstrate the expected patterns of DEGs. (JPG).

Additional file 2 — Scatterplot matrix for gene cluster that remained as liver-specific DEGs after TMM normalization.

Scatterplot matrix of the 933 genes that were in the first cluster (of Figure 18) from genes that remained as liver-specific DEGs even after TMM normalization. With this scatterplot matrix, we verify from an additional perspective that these genes demonstrate the expected patterns of DEGs. (JPG).

Additional file 3 — Scatterplot matrix for gene cluster that were removed from kidney-specific DEGs after TMM normalization.

Scatterplot matrix of the 529 genes that were in the first cluster (of Figure 19) from genes that no longer remained as kidney-specific DEGs after TMM normalization. With this scatterplot matrix, we verify from an additional perspective that these genes do not demonstrate the expected patterns of DEGs too strongly (they do not deviate much from the $x=y$ line in the treatment scatterplots). This provides additional evidence that TMM normalization removing these genes from DEG status may be valid. (JPG).

Additional file 4 — Scatterplot matrix for gene cluster that were added as liver-specific DEGs after TMM normalization.

Scatterplot matrix of the 317 genes that were in the first cluster (of Figure 20) from genes that were added as liver-specific DEGs after TMM normalization. With this scatterplot matrix, we see that the genes do not demonstrate the expected patterns of DEGs too strongly (they do not deviate much from the $x=y$ line in the treatment scatterplots). In fact, these pink genes appear similarly to what we saw from the scatterplot matrix of the red genes (Additional file 3). This is somewhat of a surprise, given that the pink genes were *added* by TMM normalization, while the red genes were *removed* by TMM normalization. Stated differently, we would expect the pink genes to appear more like differentially expressed genes if TMM normalization is appropriate, but we could not confirm this expectation. We solved this problem using standardization techniques (Figures 23 and 24). (JPG).

Additional file 5 — Example litre plots for genes that remained as kidney-specific DEGs after TMM normalization.

Example litre plots from the 1,136 genes that were in the first cluster (Figure 17) of genes that remained as kidney-specific DEGs even after TMM normalization. With these litre plots, we verify from an additional perspective that these genes demonstrate the expected patterns of DEGs. (JPG).

Additional file 6 — Example litre plots for genes that remained as liver-specific DEGs after TMM normalization.

Example litre plots from the 933 genes that were in the first cluster (Figure 18) from genes that remained as liver-specific DEGs even after TMM normalization. With these litre plots, we verify from an additional perspective that these genes demonstrate the expected patterns of DEGs. (JPG).

Additional file 7 — Example litre plots for genes that were removed from kidney-specific DEGs after TMM normalization.

Example litre plots from the 529 genes that were in the first cluster (Figure 19) of genes that no longer remained as kidney-specific DEGs after TMM normalization. With these litre plots, we verify from an additional perspective that these genes do not demonstrate the expected patterns of DEGs. This provides additional evidence that TMM normalization removing these genes from DEG status may be valid. (JPG).

Additional file 8 — Example litre plots for genes that were added as liver-specific DEGs after TMM normalization.

Example litre plots from the 317 genes that were in the first cluster (Figure 20) from genes that were added as liver-specific DEGs after TMM normalization. With these litre plots, we see that the genes do not demonstrate the expected patterns of DEGs in a trustworthy manner. In fact, these pink genes appear similarly to what we saw from the example litre plots of the red genes (Additional file 7). This is somewhat of a surprise, given that the pink genes were *added* by TMM normalization, while the red genes were *removed* by TMM normalization. Stated differently, we would expect the pink genes to appear more like differentially expressed genes if TMM normalization is appropriate, but we could not confirm this expectation. We solved this problem using standardization techniques (Figures 27 and 28). (JPG).