



OPEN

## Adapted tensor decomposition and PCA based unsupervised feature extraction select more biologically reasonable differentially expressed genes than conventional methods

Y-h. Taguchi<sup>1✉</sup> & Turki Turki<sup>2</sup>

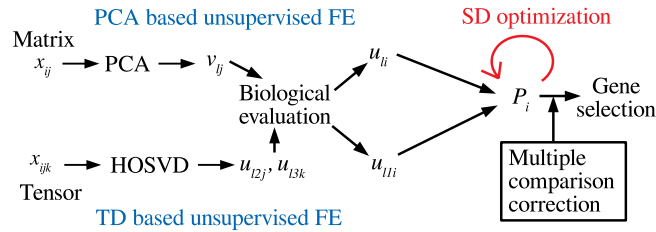
Tensor decomposition- and principal component analysis-based unsupervised feature extraction were proposed almost 5 and 10 years ago, respectively; although these methods have been successfully applied to a wide range of genome analyses, including drug repositioning, biomarker identification, and disease-causing genes' identification, some fundamental problems have been identified: the number of genes identified was too small to assume that there were no false negatives, and the histogram of  $P$  values derived was not fully coincident with the null hypothesis that principal component and singular value vectors follow the Gaussian distribution. Optimizing the standard deviation such that the histogram of  $P$  values is as much as possible coincident with the null hypothesis results in an increase in the number and biological reliability of the selected genes. Our contribution was that we improved these methods so as to be able to select biologically more reasonable differentially expressed genes than the state of art methods that must empirically assume negative binomial distributions and dispersion relation, which is required for the selecting more expressed genes than less expressed ones, which can be achieved by the proposed methods that do not have to assume these.

Identifying differentially expressed genes (DEGs) on the basis of comparative analyses<sup>1,2</sup> has always been difficult. This challenge is attributable to multiple reasons; however, the primary reason is being a *large  $p$  small  $n$*  problem. In a *large  $p$  small  $n$*  problem, it is difficult to select features based on statistical criteria because a small number of samples ( $= n$ ) have a tendency to lead to low significance; in reality, the obtained  $P$  values must be heavily corrected by considering a large number of features ( $= p$ ). This makes it difficult to find features with significance. To resolve this difficulty, many methods specific to gene expression analysis have been proposed. For example, significant analysis microarray (SAM)<sup>3</sup> adds a small amount of constancy to gene expression, thereby avoiding the misidentification of low expressed genes as DEGs. Limma<sup>4</sup> applied a Bayesian strategy to logarithmic gene expression. After high-throughput sequencing (HTS) became popular,  $P$  values are attributed to individual genes, assuming that gene expression follows a negative binomial (NB) distribution<sup>5,6</sup>, which is one of the simplest positively valued distributions with a tunable mean and variance. In addition to this, the so-called dispersion relation<sup>5,6</sup>,

$$\frac{\alpha(\mu)}{\mu^2} = \alpha_0 + \frac{\alpha_1}{\mu}, \quad (1)$$

has also been assumed, where  $\mu$  and  $\alpha$  are the mean and variance, respectively, and  $\alpha_0$  and  $\alpha_1$  are regression coefficients; to our knowledge, Eq. (1) is purely empirical and lacks rationalization. Despite these difficulties, many proposed state-of-art methods<sup>5-9</sup> have been widely employed and used in various studies.

<sup>1</sup>Department of Physics, Chuo University, 1-13-27 Kasuga, Bunkyo-ku, Tokyo 112-8551, Japan. <sup>2</sup>Department of Computer Science, King Abdulaziz University, Jeddah 21589, Saudi Arabia. ✉email: tag@granular.com



**Figure 1.** Schematic figure of TD- and PCA-based unsupervised FE with optimized SD.

Contrary to these empirical methods, we proposed tensor decomposition (TD)- and principal component analysis (PCA)-based unsupervised feature extraction (FE)<sup>10</sup> that only assumes that principal component (PC) and singular value vectors (SVVs) obey Gaussian distribution. Despite this simplicity, TD- and PCA-based unsupervised FE have been successfully applied to a wide range of genomic analyses. However, there have been two problems: 1. The histogram of the *P* values is not fully coincident with the null hypothesis that PC and SVV obey Gaussian distribution and 2. The number of genes selected is too small to have no false negatives. In this paper, we have shown that the optimization of standard deviation (SD) in Gaussian distribution can resolve these problems.

We tried optimizing SD for PCA-based unsupervised FE and applied this to two highly curated data sets—MAQC and SEQC. Then, we tested the optimization of SD for TD-based unsupervised FE and applied it to two more realistic problems: (1) drug repositioning for SARS-CoV-2 and (2) the analysis of gene expression of multiple organs treated with multiple drugs, to which TD-based unsupervised FE without SD optimization was already applied.

Our contributions are as follows. First, our methods allow more expressed genes to be more selected as DEGs without empirical dispersion relation, Eq. (1). Second, our methods can select significant DEGs without assuming not rationalized negative binomial distribution for individual gene expression. Third, our selected DEGs are much more biologically reasonable than those selected by other state of art methods.

## Results

**Outlines of TD and PCA based unsupervised FE.** In this section, we have briefly explained the algorithm of PCA- and TD-based unsupervised FE (Fig. 1) before explaining how we could improve them.

When a gene expression profile is formatted as a matrix,  $x_{ij} \in \mathbb{R}^{N \times M}$ , which represents the gene expression of the *i*th gene of the *j*th sample, we use PCA-based unsupervised FE. After standardizing  $x_{ij}$  as

$$\sum_i x_{ij} = 0 \tag{2}$$

$$\sum_i x_{ij}^2 = N, \tag{3}$$

a gram matrix  $\sum_j x_{ij}x_{i'j} \in \mathbb{R}^{N \times N}$  was diagonalized as

$$\sum_{i'} \left( \sum_j x_{ij}x_{i'j} \right) u_{\ell i'} = \lambda_\ell u_{\ell i} \tag{4}$$

where  $u_{\ell i} \in \mathbb{R}^{N \times N}$  is the  $\ell$ th PC score attributed to gene *i*. The  $\ell$ th PC loading attributed to the *j*th sample can be computed as

$$v_{\ell j} = \sum_i x_{ij}u_{\ell i} \in \mathbb{R}^{M \times M}. \tag{5}$$

After identifying  $v_{\ell j}$ , which is associated with a desired property, e.g., the district between control and treated samples, we attributed the *P* values to the gene *i* using the corresponding PC score,  $u_{\ell i}$ , as

$$P_i = P_{\chi^2} \left[ > \left( \frac{u_{\ell i}}{\sigma_\ell} \right)^2 \right] \tag{6}$$

assuming that  $u_{\ell i}$  obeys the Gaussian distribution, where  $P_{\chi^2}[> x]$  is cumulative  $\chi^2$  distribution when an argument larger than *x* and  $\sigma_\ell$  is the SD,

$$\sigma_\ell = \sqrt{\frac{1}{N} \sum_{i=1}^N (u_{\ell i} - \langle u_{\ell i} \rangle_i)^2} \tag{7}$$

$$\langle u_{\ell i} \rangle_i = \frac{1}{N} \sum_{i=1}^N u_{\ell i} \quad (8)$$

When we have gene expression that is formatted as a tensor,  $x_{ijk} \in \mathbb{R}^{N \times M \times K}$ , for the expression of the  $i$ th gene at  $j$ th sample with the  $k$ th condition, we used TD-based unsupervised FE. After standardizing  $x_{ijk}$  as

$$\sum_i x_{ijk} = 0 \quad (9)$$

$$\sum_i x_{ijk}^2 = N \quad (10)$$

Tucker decomposition of  $x_{ijk}$

$$x_{ijk} = \sum_{\ell_1=1}^N \sum_{\ell_2=1}^M \sum_{\ell_3=1}^K G(\ell_1 \ell_2 \ell_3) u_{\ell_1 i} u_{\ell_2 j} u_{\ell_3 k} \quad (11)$$

can be computed with a higher order singular value decomposition (HOSVD)<sup>10</sup>. After identifying which  $u_{\ell_2 j} \in \mathbb{R}^{M \times M}$  and  $u_{\ell_3 k} \in \mathbb{R}^{K \times K}$  are coincident with the target property, e.g., distinction between control and treated samples specifically under  $k$ th experimental condition, we try to find  $u_{\ell i} \in \mathbb{R}^{N \times N}$  associated with  $G(\ell_1 \ell_2 \ell_3) \in \mathbb{R}^{N \times M \times K}$  having the largest absolute value. Then, the  $P$  value is attributed to the  $i$ th gene as

$$P_i = P_{\chi^2} \left[ > \left( \frac{u_{\ell_1 i}}{\sigma_{\ell_1}} \right)^2 \right] \quad (12)$$

by also assuming that  $u_{\ell_1 i}$  obeys the Gaussian distribution and

$$\sigma_{\ell_1} = \sqrt{\frac{1}{N} \sum_{i=1}^N (u_{\ell_1 i} - \langle u_{\ell_1 i} \rangle_i)^2} \quad (13)$$

$$\langle u_{\ell_1 i} \rangle_i = \frac{1}{N} \sum_{i=1}^N u_{\ell_1 i} \quad (14)$$

For both PCA- and TD-based unsupervised FE,  $P_i$  is corrected with the Benjamini-Hochberg (BH) criterion<sup>10</sup>; further, the  $i$ th genes associated with adjusted  $P_i$  less than the threshold value, which is usually 0.01, are selected.

Although PCA- as well as TD-based unsupervised FE were successfully applied to a wide range of genomic analyses, there were two weak points:

- Too small a number of genes were selected to have no false negatives.
- The histogram of  $P_i$  did not fully obey the null assumption that  $u_{\ell i}$  and  $u_{\ell_1 i}$  obey the Gaussian distribution.

In this paper, by fixing these two problems, we have tried to establish a new method at least comparable to or even superior to state-of-art methods.

**Trials using highly curated data sets.** *Application to MAQC dataset.* Initially, to assess what the problem is, we compared the performance of PCA-based unsupervised FE with DESeq2, a state-of-art method, using the MAQC<sup>11</sup> data set, which has been carefully curated and frequently used for benchmark studies.

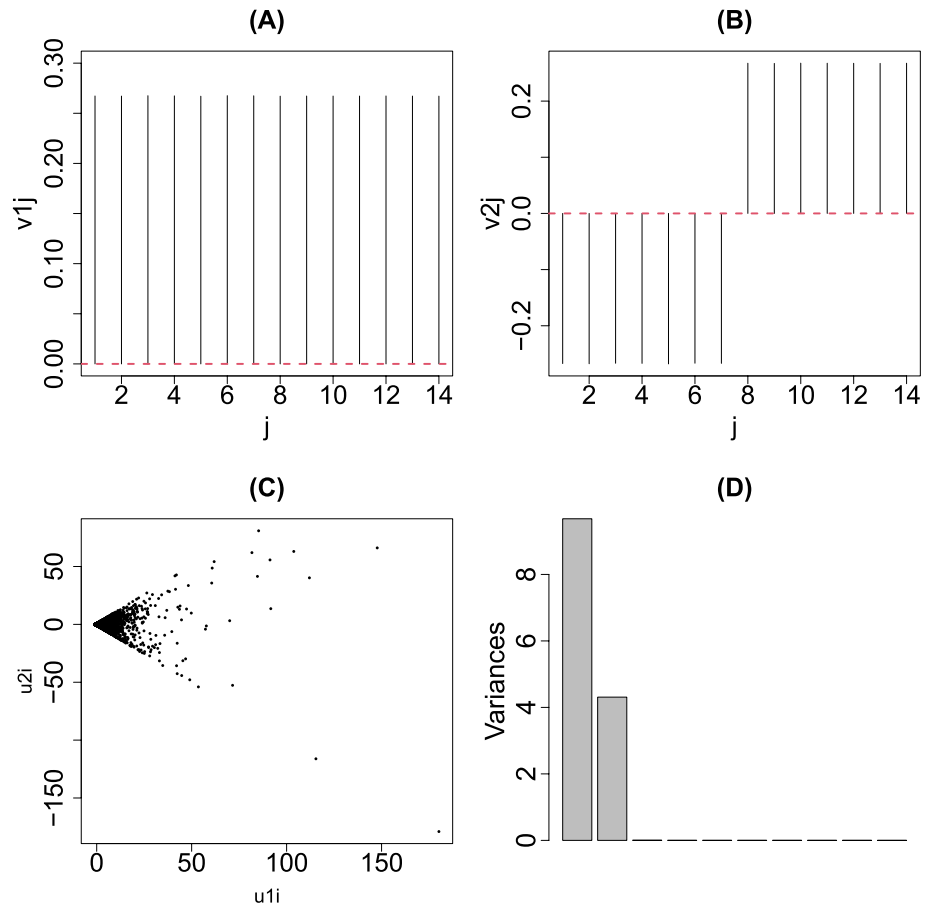
Figure 2C shows a scatter plot of genes using  $u_{1j}$  and  $u_{2j}$ . Figure 2A,B show the PC loading  $v_{1j}$  and  $v_{2j}$ ;  $v_{1j}$  represents the mean gene expression and  $v_{2j}$  represents the differential expression between universal human reference (UHR) and brain. Occasionally, this reminds us of the horizontal and vertical axes of an MAPlot; the horizontal axis of an MAPlot represents the mean expression of individual genes, typically the mean logarithmic expression,

$$\frac{1}{M} \sum_{j=1}^M \log_2 x_{ij} \quad (15)$$

whereas the vertical axis of an MAPlot represents the differential expression between the two classes, typically the mean logarithmic fold change (LFC),

$$\frac{1}{M_A} \sum_{j \in A} \log_2 x_{ij} - \frac{1}{M_B} \sum_{j \in B} \log_2 x_{ij} \quad (16)$$

where  $M_A$  and  $M_B (= M - M_A)$  are sample numbers within one of the two classes, A and B, respectively, and summations are taken within individual classes. As can be seen in Fig. 2D, which represents the contribution of PC loading,  $x_{ij}$  can be expressed almost fully in the 2-dimensional space spanned by the first two PCs. Thus, PCA can derive, in a fully unsupervised manner, something that qualitatively corresponds to an MAPlot (Fig. 8), which



**Figure 2.** PCA applied to MAQC data (A)  $v_{1j}$  (B)  $v_{2j}$  (C) Scatter plot of  $u_{1i}$  and  $u_{2i}$  (D) Contributions of individual PCs.

is usually drawn artificially. In spite of that, unfortunately, the genes selected by the adjusted  $P_i$  are too small to have no false negatives (Table 3) and an histogram of  $P_i$  is hardly regarded to obey the null hypothesis; the left panel of Fig. 3 shows the histogram of  $1 - P_i$ , where  $P_i$ s were computed from  $u_{2i}$  by Eq. (6) using  $\sigma_2$  defined as

$$\sigma_2 = \sqrt{\frac{1}{N} \sum_i (u_{2i} - \langle u_{2i} \rangle)^2} \quad (17)$$

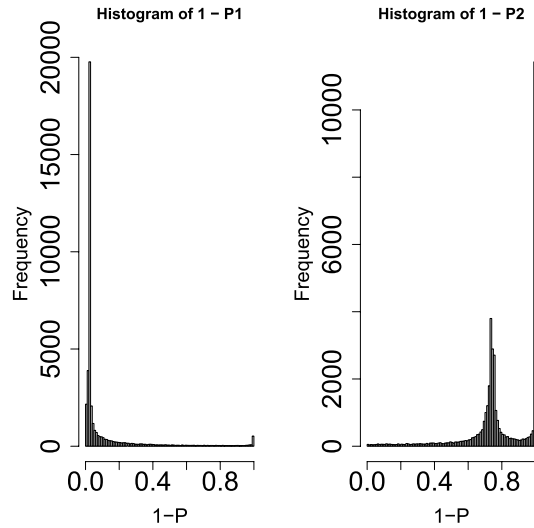
$$\langle u_{2i} \rangle = \frac{1}{N} \sum_i u_{2i}. \quad (18)$$

If  $1 - P_i$  is coincident with the null hypothesis; the histogram of  $1 - P_i < 1$  should have a flat distribution and that of  $1 - P_i \sim 1$  should have a sharp peak.

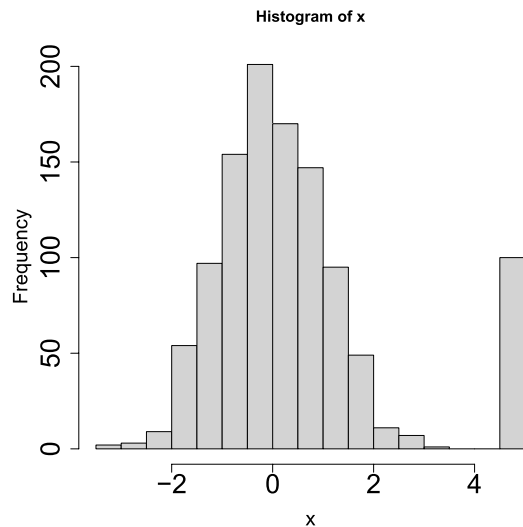
**Top ranked genes are coincident with DESeq2.** To understand the problem of  $P_i$ s computed by PCA-based unsupervised FE, we compared  $P_i$ s computed by PCA-based unsupervised FE with those computed by DESeq2, a state-of-art method. At first, AUC was computed to predict the top 1000 genes based on  $P_i$  derived with DESeq2 using  $P_i$ s computed by PCA-based unsupervised FE; the area under the curve (AUC) was 0.97. Next, in contrast, the AUC was computed to predict the top 1000 genes based on  $P_i$  derived with PCA-based unsupervised FE using  $P_i$ s computed using DESeq2; the AUC was 0.98. This indicated that the top-ranked genes were suitably shared between PCA-based unsupervised FE and DESeq2. Thus, the problem of PCA-based unsupervised FE is not the genes' ranking but the absolute value of  $P_i$ s.

**Optimization of SD.** Based on the observations at the end of the previous subsection, we arrived at optimizing  $\sigma_\ell$  such that  $u_{\ell i}$  and  $u_{\ell+1 i}$  obeyed the Gaussian distribution. Generally, optimizing SD to be fitted to the null hypothesis is not easy. For example, Mudge et al<sup>12</sup> had to assume the equivalence between Type I and II errors, which we cannot assume because of an imbalance of numbers between DEGs and the other genes; typically,





**Figure 3.** Histogram of  $1 - P_i$  of the MAQC data set with PCA-based unsupervised FE Left:  $P_i$ s by Eq. (6) using SD  $\sigma_2$  directly computed from  $u_{2i}$ , right: using SD optimized to obey the Gaussian distribution as much as possible.



**Figure 4.** A histogram of Gaussian distribution with outliers.

DEGs are expected to be minorities. Next, we decided to employ an alternative and more empirical approach. To visualize the idea, we have shown some illustrative examples.

Figure 4 shows a histogram of the variable  $x_i$  derived from the Gaussian distribution and outliers. If we attribute the  $P$  values to the  $i$ th variable with  $x_i$

$$P_i = P_{\chi^2} \left[ > \left( \frac{x_i}{\sigma} \right)^2 \right] \tag{19}$$

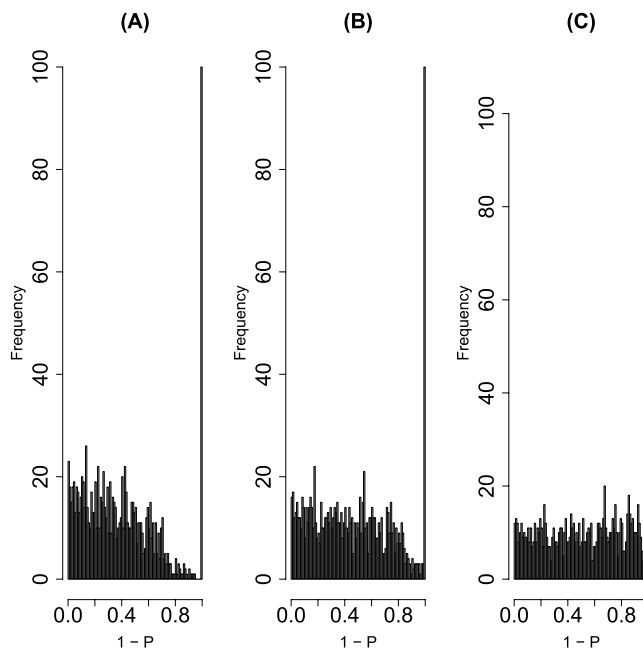
using the SD,  $\sigma$ , directly computed by all points

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \langle x_i \rangle)^2} \tag{20}$$

$$\langle x_i \rangle = \frac{1}{N} \sum_{i=1}^N x_i \tag{21}$$

|           | True                            | Not outliers | Outliers |
|-----------|---------------------------------|--------------|----------|
| Predicted | Adjusted $P$ values $> 0.01$    | 1000         | 100      |
|           | Adjusted $P$ values $\leq 0.01$ | 0            | 0        |

**Table 1.** Confusion matrix of the Gaussian distribution with outliers and prediction for  $x_i$ , the histogram for which is given in Fig. 4.



**Figure 5.** Histograms of  $1 - P_i$ ,  $h(1 - P_i)$ , for (A)  $P_i$  computed by Eq. (6) with  $\sigma$  defined in Eq. (20), (B) that with optimized SD, (C) that with true SD,  $\sigma = 1$ .

and select outliers associated with adjusted  $P$  values  $< 0.01$ , we cannot select any of the outliers (Table 1); this is because the SD computed,  $\sigma = \frac{1000 \times 1 + 100 \times 5^2}{1000 + 100} = 1.75$ , is larger than that of the Gaussian distribution,  $\sigma = 1$ , because of outliers. Because  $P_i$ s computed with  $\sigma = 1.75$  is larger than that with  $\sigma = 1$ , it fails to recognize outliers correctly.

We computed the histogram of  $1 - P_i$ , Fig. 5A, which is far being idealized, Fig. 5C, that should have a constant histogram  $h(1 - P_i)$  up to  $1 - P_i$  very close to 1 and has one with a narrow peak near  $1 - P_i \sim 1$ . To optimize the SD, we tried to find an optimal SD such that the histogram for those not recognized as outliers was as flat as possible, i.e., obeying the null hypothesis of the Gaussian distribution; we decided to find the optimal SD that results in the most flat  $h(1 - P_i)$  for  $1 - P_i$  adjusted  $P_i$  less than threshold value  $1 - P_0$  (adjusted  $P_0$  should be small enough). To minimize the SD of binned  $h_i = h(1 - P_i)$ ,  $\sigma_h$ ,

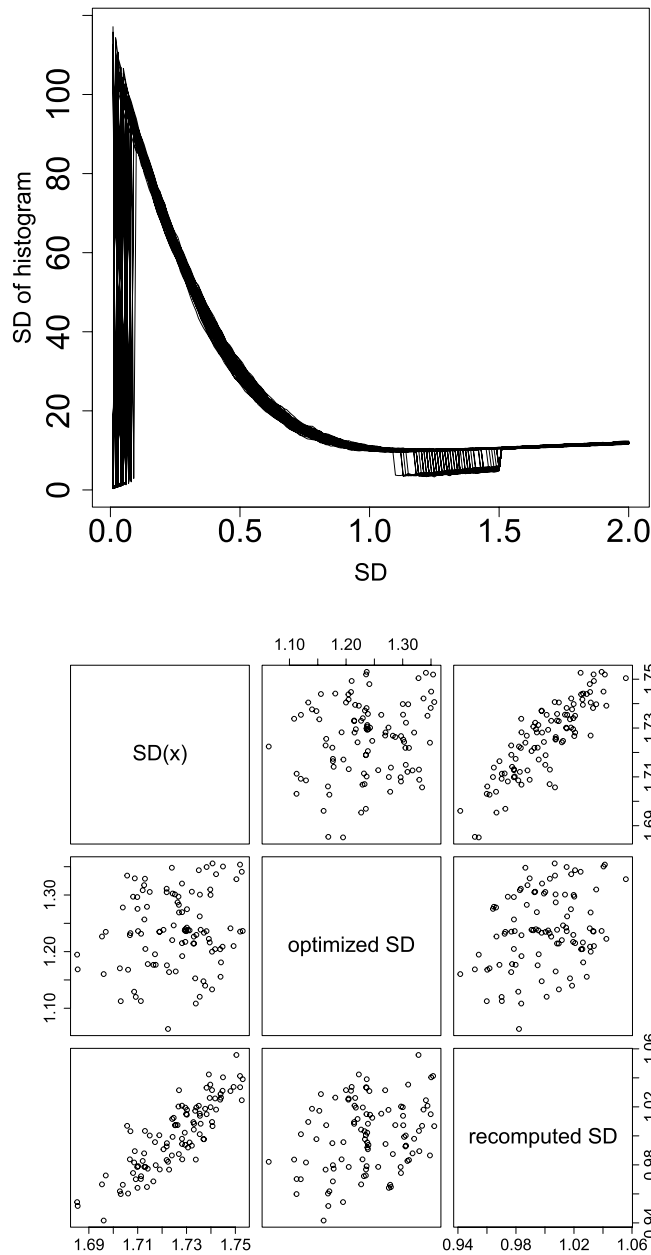
$$\sigma_h = \sqrt{\frac{\sum_{\text{adjusted } P_i < \text{adjusted } P_0} (h_i - \langle h_i \rangle)^2}{N(\text{adjusted } P_0)}} \quad (22)$$

$$\langle h_i \rangle = \frac{\sum_{\text{adjusted } P_i < \text{adjusted } P_0} h_i}{N(\text{adjusted } P_0)} \quad (23)$$

with respect to  $\sigma$ , where  $N(\text{adjusted } P_0)$  is the number of  $h_i$ s associated with adjusted  $P_i > \text{adjusted } P_0$ , i.e., not recognized as outliers and recognized as a part of the Gaussian distribution. After optimizing  $\sigma_\ell$ , we recomputed  $P_i$ . Figure 5A,B show the histogram of  $1 - P_i$  using  $\sigma = 1.75$  and optimized SD, respectively; the latter is closer to an idealized histogram of  $P_i$ , Fig. 5C, than the former.

To validate the effectiveness of the optimization of SD, we repeated this procedure 100 times.

Figure 6 shows the dependence of  $\sigma_h$  on SD (upper panel) and the comparison between SD in Eq. (20), optimized SD, and SD computed using  $is$  for adjusted  $P_i < \text{adjusted } P_0$  (lower panel). In the lower panel, the optimized SD was approximately 1.2, which is much closer to 1 than 1.75, computed by Eq. (20). In addition, the fact that SD computed using  $is$  for adjusted  $P_i < \text{adjusted } P_0$ , which is expected to correspond to the Gaussian distribution part in Fig. 4, is almost 1 helps justify our optimization procedure (Fig. 6, lower panel). The reason



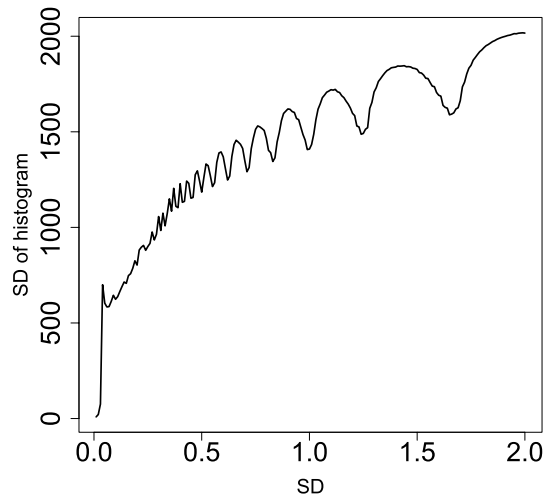
**Figure 6.** Scatter plot of SDs. Upper:  $\sigma_h$ , defined in Eq. (22) as a function of SD used for computing  $P_i$  in Eq. (19). Lower: Scatter plot  $\sigma$  of Eq. (20), optimized SD, and SD computed using  $is$  with adjusted  $P_i < \text{adjusted } P_0$  (recomputed SD).

|         | True                            | Not outliers | Outliers |
|---------|---------------------------------|--------------|----------|
| Predict | Adjusted $P$ values $> 0.01$    | 1000         | 0        |
|         | Adjusted $P$ values $\leq 0.01$ | 0            | 100      |

**Table 2.** Averaged confusion matrix of Gaussian distribution with outliers and prediction using optimized SD.

why  $SD = 0$  with  $\sigma_h = 0$  in the upper panel of Fig. 6 was not selected as optimal (as having the smallest  $\sigma_h$ ) is because  $\sigma = 0$  corresponds to nothing selected and is thus meaningless. Using  $P_i$  computed by optimized SD, we can discriminate the outliers almost perfectly (Table 2).

Next, we applied this strategy to the MAQC data set. Figure 7 shows  $\sigma_h$ , defined in Eq. (22), as a function of SD to compute  $P_i$  in Eq. (19) using the MAQC data set; the optimal SD was 0.05557979. It is close to the SD



**Figure 7.**  $\sigma_h$ , defined in Eq. (22) as a function of SD used for computing  $P_i$  in Eq. (19) using MAQC data.

|                               | Adjusted $P_i$ |             |
|-------------------------------|----------------|-------------|
|                               | > 0.01         | $\leq$ 0.01 |
| PCA based unsupervised FE     | ---            | ---         |
| Original (without optimal SD) | 40589          | 344         |
| With optimal SD               | 28681          | 12252       |
| DESeq2                        | 8789           | 20546       |

**Table 3.** The number of genes selected with original PCA-based unsupervised FE, that with optimal SD, and DESeq2.

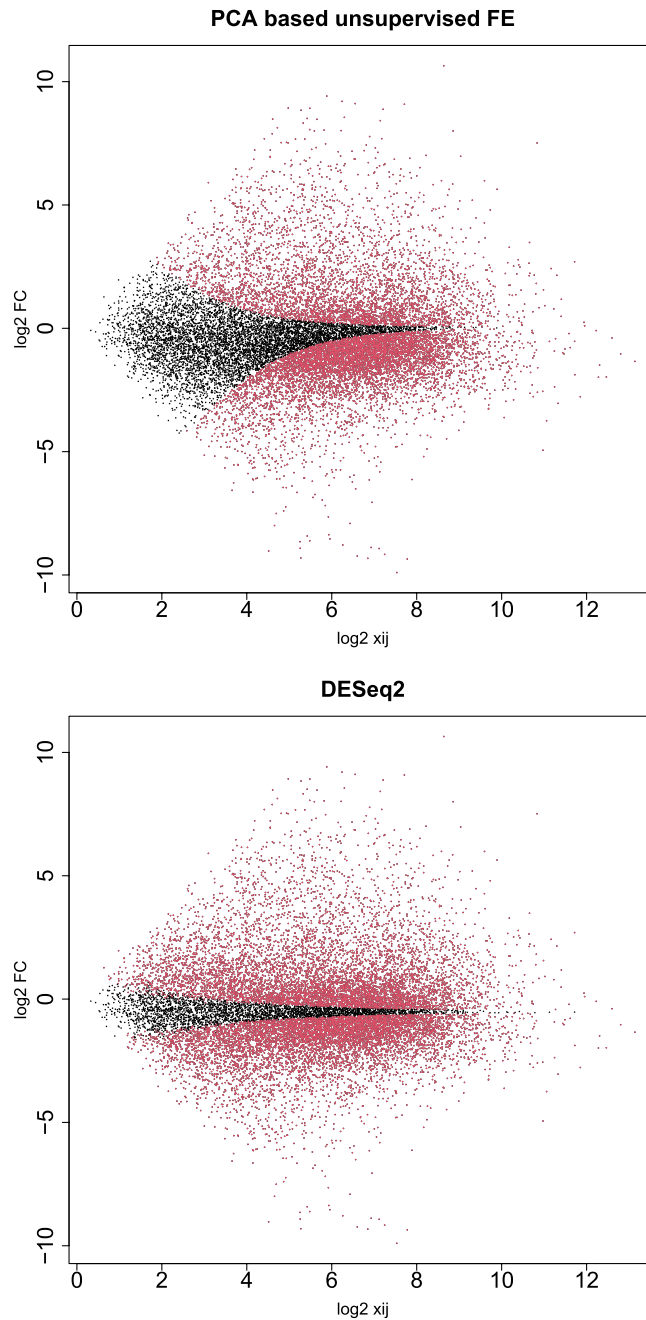
recomputed using  $is$  with adjusted  $P_i < \text{adjusted } P_0, 0.03871846$ ; moreover,  $h(1 - P_i)$  derived from optimal SD looks more idealized (the right panel of Fig. 3). Thus, the optimal SD improved PCA-based unsupervised FE.

Table 3 shows the number of genes selected using DESeq2 (list of genes available as Data S1), the original PCA-based unsupervised FE, than by using optimal SD (list of genes available as Data S2). Although the number of genes selected by original PCA-based unsupervised FE, 344, is too small to regard no false negatives, that of genes selected by PCA-based unsupervised FE with optimal SD, 12252, is large enough to regard no false negatives. Furthermore, that of DESeq2, 20546, seems to be too large to have no false positives, because it is unlikely true that more than half the genes (40933) are distinctly expressed between the brain and controls.

*Less expressed genes are less likely to be DEGs.* Figure 8 shows the selected genes in MAPlot. Although we assumed neither NB distribution nor dispersion relation, Eq. (1), the distribution of selected genes in the MAPlot is reasonable; genes with the same LFC (vertical axis) are less likely selected when associated with smaller mean expression (horizontal axis). Although this property is explicitly assumed in DESeq2 with dispersion relation, Eq. (1), PCA-based unsupervised FE seems to possess the property without assuming dispersion relation explicitly (see the “Discussion” section). On the other hand, DESeq2 selects too many genes and is less likely reasonable. This suggests that PCA-based unsupervised FE with optimized  $\sigma_\ell$  is a promising method.

*Confirmation using the SEQC dataset.* To see if it occurs only occasionally, we repeated all computations on as many as 13 data sets in SEQC<sup>13</sup>, which is yet another curated data set. Coincidence between DESeq2 and PCA-based unsupervised FE (Fig. 9), a reasonable number of selected genes ( $\sim 10^3$ , Fig. 10), and a lower opportunity of less expressed genes to be DEGs (Fig. 11) are also observed, as in the case of MAQC. In addition to this, although the number of genes selected by DESeq2 are too large ( $\sim 10^4$ ) and heavily dependent upon sample numbers ( $\sim 10^3$  for the smallest sample number  $\sim 10^0$ ), that by PCA-based unsupervised FE is not and is always  $\sim 10^3$ , regardless of sample numbers. Thus, PCA-based unsupervised FE is seemingly superior to DESeq2.

*Biological validation.* Based on the above results, PCA-based unsupervised FE is seemingly better than DESeq2. Nonetheless, PCA-based unsupervised FE can select a reasonable number of genes regardless of sample numbers (Fig. 10), and less expressed genes are unlikely to be DEGs when genes are selected by PCA-based unsupervised FE with optimized SD (Figs. 8, 11), even without assuming NB distribution and dispersion relations, Eq. (1), which DESeq2 requires, if the selected genes are not biological, it is meaningless. To evaluate the selected genes biologically, we uploaded the genes selected using MAQC to Enrichr. As can be seen in Fig. 12, the genes selected

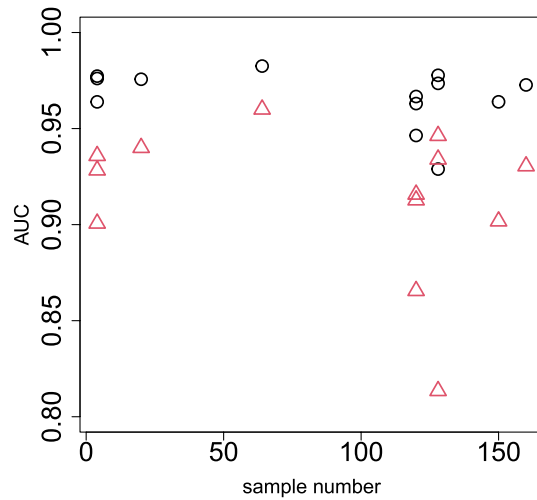


**Figure 8.** MAPlot with selected genes colored in red Upper: PCA-based unsupervised FE with optimized SD, lower: DESeq2.

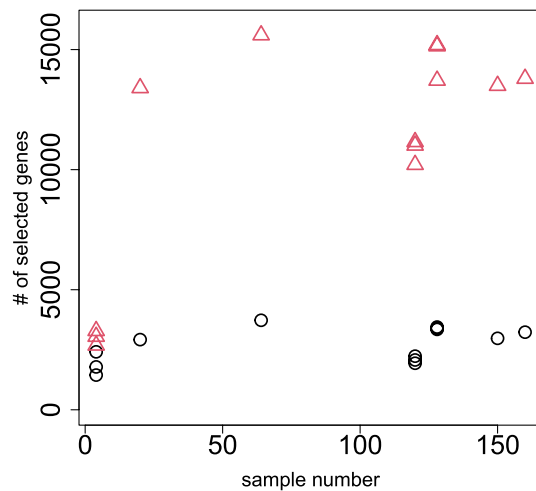
by PCA-based unsupervised FE were better than those selected by DESeq2 (Full list of enrichment analysis is available in Data S1 and S2).

One may still wonder the other state-of-art methods might be better than PCA-based unsupervised FE. To deny this possibility, we biologically evaluated the genes selected for MAQC using edgeR<sup>6</sup> (full list of enrichment analysis available in Data S3), voom<sup>8</sup> (full list of enrichment analysis available in Data S4), and NOISeq<sup>9</sup> (full list of enrichment analysis available in Data S5); it is obvious that these three methods are even inferior to DESeq2 biologically (Fig. 13).

**Drug discovery for SARS-CoV-2.** Although we have demonstrated that PCA-based unsupervised FE with optimized SD can outperform other state-of-art methods in highly curated data, one might wonder that it is not the case for a realistic and more noisy case. To check if PCA-based unsupervised FE with optimized SD can



**Figure 9.** Coincidence of top-ranked genes between DESeq2 and PCA-based unsupervised FE using the SEQC data set. Open circles: AUC when  $P$  values computed by PCA-based unsupervised FE with optimized SD discriminates top 1000 genes ranked by  $P$  values computed by DESeq2. Open red triangles: AUC when  $P$  values computed by DESeq2 discriminating top 1000 genes ranked by  $P$  values computed by PCA-based unsupervised FE with optimized SD.



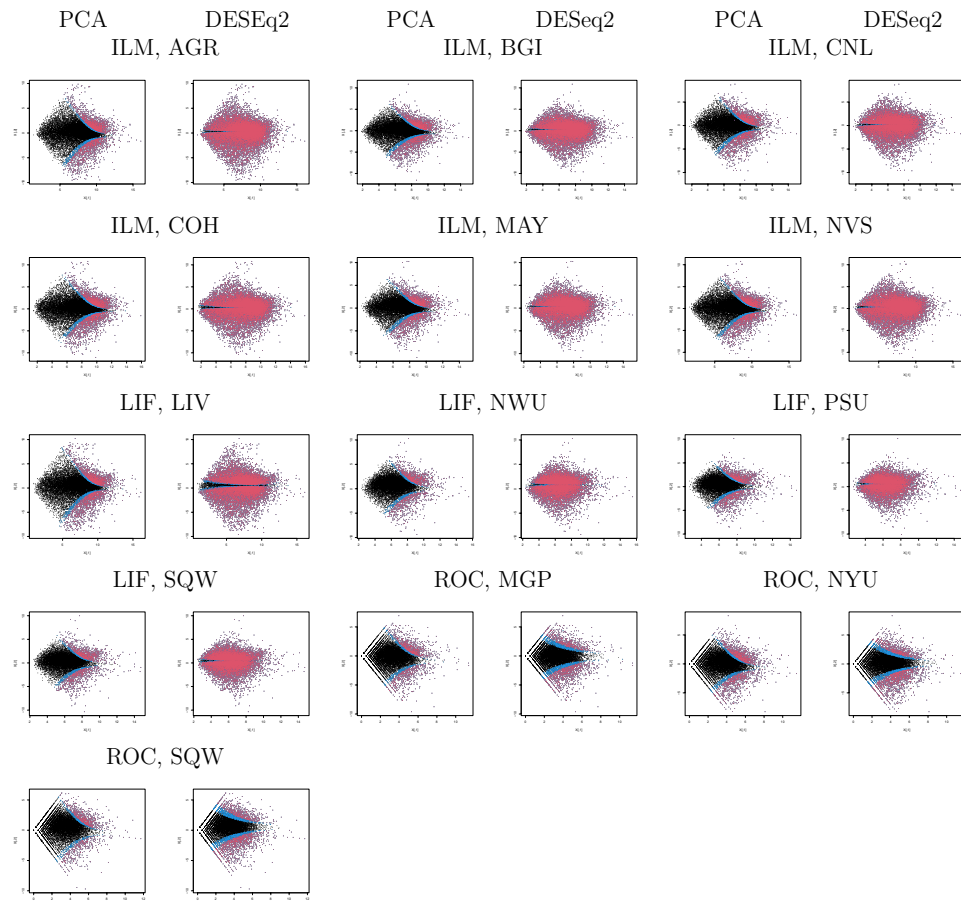
**Figure 10.** Dependence of the number of DEGs on sample numbers using the SEQC data set. Open circles: the number of genes selected by PCA-based unsupervised FE with optimized SD. Open red triangles: the number of genes selected by DESeq2.

outperform DESeq2 in more realistic data sets, we considered the drug repositioning of SARS-CoV-2, to which we applied TD-based unsupervised FE<sup>14</sup> and its kernelized version<sup>15</sup>.

In our implementation, we employed HOSVD to obtain the tensor decomposition, Eq. (11); because HOSVD is equivalent to SVD applied to a matrix obtained by unfolding a tensor, we can obtain the identical  $u_{\ell j}$  independent of which of PCA or HOSVD is used; SD used in Eq. (12) can be optimized too. Next, we applied the optimization of SD and could select 3627 genes associated with adjusted  $P$  values of less than 0.1 (list of genes available as Data S6), which is a much higher number of genes than 163 genes than that selected in previous studies<sup>14,15</sup>.

**Overlap with human genes known to interact with SARS-CoV-2 protein.** We evaluated the selected 3627 genes based on the overlap with the human genes known to interact with SARS-CoV-2, as has been done in previous studies<sup>14,15</sup> (Fig. 14). It is obvious that TD-based unsupervised FE with an optimized SD can outperform kernel TD-based unsupervised FE, original (without optimized SD) TD-based unsupervised FE as well as DESeq2 (list of overlap available in Data S7). Thus, it is indeed an outstanding method.

**Drug repositioning.** We also tried drug discovery using the genes selected by TD-based unsupervised FE with optimized SD. See Table 4 (Full list of drug repositioning available as Data S6). The first one, imatinib, was once



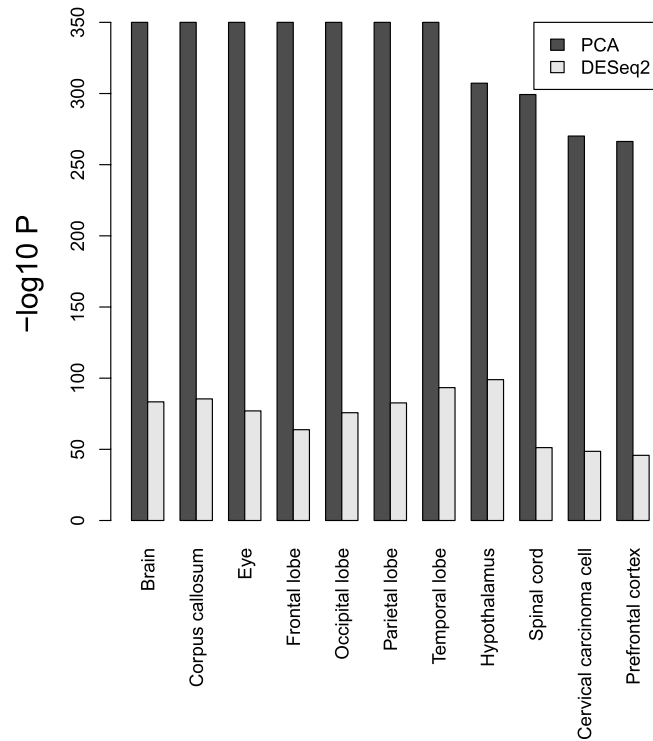
**Figure 11.** MAPlot for SEQC PCA-based unsupervised FE with optimized SD: the first, third, and fifth columns, DESeq2: the second, fourth, and sixth columns. Three character IDs represent platform and sites. Blue: genes associated with adjusted  $P$  values less than 0.1 but greater than 0.01. Red: genes associated with adjusted  $P$  values less than 0.01.

identified as a promising drug toward COVID-19, although it was rejected later<sup>16</sup>. The second one, apratoxin A, was reported to be a promising compound based on its protein binding affinity<sup>17</sup>. The third and fourth one, doxycycline, was supposed to be a promising drug toward COVID-19<sup>18</sup>. The seventh one, trovafloxacin, was reported to be a promising compound based on its protein binding affinity<sup>19</sup>. The eighth one, doxorubicin, was also reported to be a promising compound based on its protein binding affinity<sup>20</sup>. The ninth one, cisplatin, and the tenth one, carboplatin, were proposed as a result of drug repositioning<sup>21</sup>. Seven of the nine compounds identified as the top 10 compounds have been previously reported as drugs toward SARS-CoV-2.

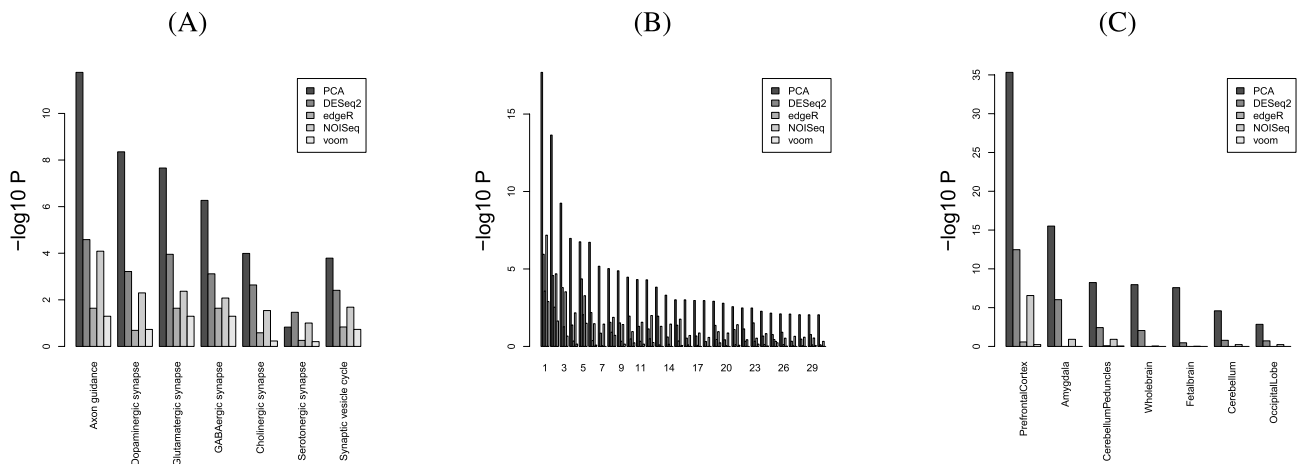
See Table 5. The first, fourth, and tenth one, estradiol, was reported as a promising compound<sup>22</sup>. The second one, tamoxifen, was reported to inhibit SARS-CoV-2 infection by suppressing viral entry<sup>23</sup>. The third one, apratoxin A, has been listed in Table 4, too. The fifth one, MK-886, was reported to be an inhibitor of 3CL protease<sup>24</sup>, although its efficiency was limited to 40%. The sixth one, IFN-alphacon1, was reported to be an inhibitor of SARS-CoV<sup>25</sup> but not for SARS-CoV-2. The seventh one, arachidonic acid, was generally expected to inhibit SARS-CoV-2 infection<sup>26</sup>. The eighth one, arsenic, was also generally expected to act against the RdRp of coronavirus<sup>27</sup>. The ninth one, metoprolol, was reported to be a promising drug toward COVID-19<sup>28</sup>. Thus, all the top 10 compounds were reported to be promising.

On the other hand, for DESeq2, see Table 6 (full list of drug repositioning is available in Data S8). The use of the second and third one, dexamethasone, resulted in lower 28-day mortality among those who received either invasive mechanical ventilation or oxygen alone at randomization but not among those receiving no respiratory support.<sup>29</sup> The seventh one, metformin, suppressed SARS-CoV-2 in cell culture<sup>30</sup>. The eighth one, etanercept, significantly decreased the risk of developing COVID-19 in patients with rheumatoid arthritis or spondyloarthropathies<sup>31</sup>. The tenth one, lipopolysaccharide, is not a compound but a bacterial protein reported to bind to the SARS-CoV-2 spike protein<sup>32</sup>.

See Table 7. The first and fourth one, resveratrol, inhibits HCoV-229E and SARS-CoV-2 coronavirus replication in vitro<sup>33</sup>. The second, third, and fifth one, carboplatin, was proposed as a result of drug repositioning<sup>21</sup>. The seventh one, lipopolysaccharide, is listed in Table 6, too.

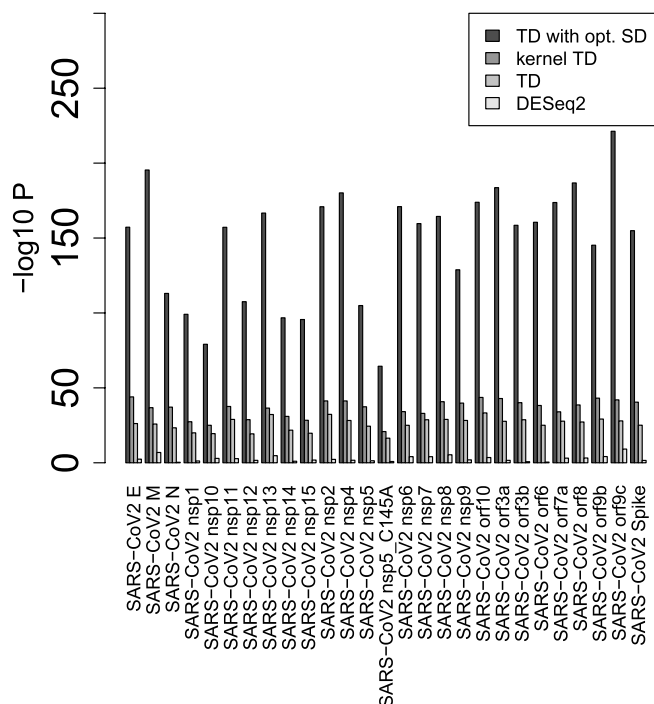


**Figure 12.** Enrichment analysis of the selected genes, whose numbers in Table 3 *P* values are adjusted *P* values (based upon “Jensen Tissues” category in Enrichr). Seven terms associated with  $-\log_{10} P = 350$  are linked with  $\infty$ , since  $P = 0$ .



**Figure 13.** Enrichment analysis for MAQC with other methods in Enrichr (A) KEGG (B) GO BP (C) Human gene atlas. Numbers in (B) correspond to 1. “axonogenesis,” 2. “axon guidance,” 3. “axon development,” 4. “regulation of axonogenesis,” 5. “synapse organization,” 6. “modulation of chemical synaptic transmission,” 7. “positive regulation of axonogenesis,” 8. “modulation of excitatory postsynaptic potential,” 9. “regulation of axon extension,” 10. “positive regulation of synaptic transmission,” 11. “axon extension,” 12. “negative regulation of axonogenesis,” 13. “chemical synaptic transmission,” 14. “signal release from synapse,” 15. “synapse assembly,” 16. “regulation of neuronal synaptic plasticity,” 17. “positive regulation of axon extension,” 18. “regulation of trans-synaptic signaling,” 19. “positive regulation of excitatory postsynaptic potential,” 20. “negative regulation of axon extension,” 21. “regulation of synapse assembly,” 22. “retrograde axonal transport,” 23. “synaptic vesicle endocytosis,” 24. “synaptic transmission, GABAergic,” 25. “synaptic transmission, glutamatergic,” 26. “regulation of long-term synaptic potentiation,” 27. “regulation of axon extension involved in axon guidance,” 28. “synaptic membrane adhesion,” 29. “regulation of synaptic transmission, glutamatergic,” 30. “regulation of postsynaptic neurotransmitter receptor activity.” *P* values are adjusted *P* values.





**Figure 14.** P values computed by Fishers’ exact test to evaluate the overlap between human genes known to interact with SARS-Cov-2 proteins and genes selected by various methods. DESeq2 is only for A549 cell lines.

| Rank | Term   | Overlap | P value                 | Adjusted P value        | Odds Ratio |
|------|--|---------|-------------------------|-------------------------|------------|
| 1    | imatinib (glivec) 123596 human GSE12211 sample 2518      | 316/442 | $7.81 \times 10^{-137}$ | $7.06 \times 10^{-134}$ | 12.3       |
| 2    | apratoxin A 6326668 human GSE2742 sample 3071            | 279/389 | $3.77 \times 10^{-121}$ | $1.57 \times 10^{-118}$ | 12.3       |
| 3    | doxycycline DB00254 human GSE2624 sample 3074            | 294/425 | $5.22 \times 10^{-121}$ | $1.57 \times 10^{-118}$ | 10.9       |
| 4    | doxycycline DB00254 human GSE2624 sample 3077            | 278/391 | $3.83 \times 10^{-119}$ | $8.64 \times 10^{-117}$ | 11.9       |
| 5    | grepafloxacin 72474 human GSE9166 sample 2627            | 320/495 | $5.62 \times 10^{-119}$ | $1.02 \times 10^{-116}$ | 8.96       |
| 6    | clinafloxacin 60063 human GSE9166 sample 2625            | 309/470 | $8.04 \times 10^{-118}$ | $1.21 \times 10^{-115}$ | 9.38       |
| 7    | trovafloxacin 62959 human GSE9166 sample 2629            | 302/459 | $3.05 \times 10^{-115}$ | $3.94 \times 10^{-113}$ | 9.38       |
| 8    | doxorubicin, 2xEC50, 5 d 31703 human GSE6930 sample 3265 | 314/493 | $4.76 \times 10^{-114}$ | $5.37 \times 10^{-112}$ | 8.57       |
| 9    | cisplatin DB00515 human GSE6410 sample 2532              | 239/315 | $1.06 \times 10^{-112}$ | $1.07 \times 10^{-110}$ | 15.1       |
| 10   | carboplatin DB00958 human GSE7035 sample 3060            | 284/422 | $4.57 \times 10^{-112}$ | $4.13 \times 10^{-110}$ | 9.99       |

**Table 4.** Drug perturbations from GEO down.

| Rank | Term   | Overlap | P value                 | Adjusted P value        | Odds Ratio |
|------|--|---------|-------------------------|-------------------------|------------|
| 1    | estradiol 5757 human GSE4668 sample 3063           | 276/367 | $1.26 \times 10^{-128}$ | $1.14 \times 10^{-125}$ | 14.74      |
| 2    | tamoxifen DB00675 human GSE4025 sample 2820        | 271/361 | $6.30 \times 10^{-126}$ | $2.85 \times 10^{-123}$ | 14.61      |
| 3    | apratoxin A 6326668 human GSE2742 sample 3068      | 278/389 | $4.61 \times 10^{-120}$ | $1.12 \times 10^{-117}$ | 12.16      |
| 4    | estradiol DB00783 human GSE4668 sample 2727        | 261/350 | $4.96 \times 10^{-120}$ | $1.12 \times 10^{-117}$ | 14.19      |
| 5    | MK-886 CID 3651377 human GSE3202 sample 3193       | 268/368 | $5.29 \times 10^{-119}$ | $9.59 \times 10^{-117}$ | 12.98      |
| 6    | IFN-alphacon1 DB05258 human GSE5542 sample 2474    | 242/313 | $2.21 \times 10^{-117}$ | $3.34 \times 10^{-115}$ | 16.41      |
| 7    | Arachidonic acid DB04557 human GSE3737 sample 3171 | 277/395 | $2.80 \times 10^{-116}$ | $3.63 \times 10^{-114}$ | 11.39      |
| 8    | ARSENIC 5359596 human GSE6907 sample 3529          | 276/394 | $1.15 \times 10^{-115}$ | $1.30 \times 10^{-113}$ | 11.35      |
| 9    | metoprolol DB00264 human GSE3356 sample 2786       | 306/469 | $2.67 \times 10^{-115}$ | $2.68 \times 10^{-113}$ | 9.16       |
| 10   | estradiol 5757 human GSE4668 sample 3062           | 245/325 | $1.92 \times 10^{-114}$ | $1.74 \times 10^{-112}$ | 14.75      |

**Table 5.** Drug perturbations from GEO up.

| Rank | Term  | Overlap | P value                | Adjusted P value       | Odds Ratio |
|------|---|---------|------------------------|------------------------|------------|
| 1    | PLX4032 DB05238 human GSE24862 sample 2568            | 65/318  | $1.59 \times 10^{-29}$ | $1.42 \times 10^{-26}$ | 7.06       |
| 2    | dexamethasone DB01234 human GSE34313 sample 2714      | 51/297  | $7.68 \times 10^{-20}$ | $3.44 \times 10^{-17}$ | 5.59       |
| 3    | dexamethasone DB01234 human GSE54608 sample 3093      | 52/322  | $5.45 \times 10^{-19}$ | $1.63 \times 10^{-16}$ | 5.19       |
| 4    | VX 39793 human GSE33606 sample 3376                   | 54/367  | $8.17 \times 10^{-18}$ | $1.58 \times 10^{-15}$ | 4.65       |
| 5    | PLX4032 DB05238 human GSE24862 sample 2570            | 56/393  | $8.78 \times 10^{-18}$ | $1.58 \times 10^{-15}$ | 4.49       |
| 6    | formoterol DB00983 human GSE30242 sample 2631         | 49/315  | $2.83 \times 10^{-17}$ | $4.23 \times 10^{-15}$ | 4.94       |
| 7    | metformin DB00331 human GSE33612 sample 2483          | 50/343  | $2.07 \times 10^{-16}$ | $2.65 \times 10^{-14}$ | 4.58       |
| 8    | etanercept DB00005 human GSE41663 sample 2605         | 45/322  | $3.29 \times 10^{-14}$ | $3.69 \times 10^{-12}$ | 4.33       |
| 9    | cisplatin DB00515 human GSE47856 sample 3145          | 40/267  | $8.93 \times 10^{-14}$ | $8.91 \times 10^{-12}$ | 4.68       |
| 10   | Lipopolysaccharide 11970143 human GSE5504 sample 3486 | 35/224  | $9.25 \times 10^{-13}$ | $8.30 \times 10^{-11}$ | 4.89       |

**Table 6.** Drug perturbations from GEO down for A549 by DESeq2.

| Rank | Term   | Overlap | P value                | Adjusted P value       | Odds Ratio |
|------|--|---------|------------------------|------------------------|------------|
| 1    | resveratrol DB02709 human GSE25412 sample 3500         | 70/250  | $2.90 \times 10^{-41}$ | $2.63 \times 10^{-38}$ | 10.81      |
| 2    | carboplatin (30 h) 10339178 human GSE13525 sample 3031 | 85/423  | $7.47 \times 10^{-38}$ | $3.38 \times 10^{-35}$ | 7.09       |
| 3    | carboplatin (36 h) 10339178 human GSE13525 sample 3032 | 74/392  | $3.93 \times 10^{-31}$ | $1.19 \times 10^{-28}$ | 6.46       |
| 4    | resveratrol DB02709 human GSE25412 sample 3501         | 51/194  | $7.59 \times 10^{-29}$ | $1.72 \times 10^{-26}$ | 9.66       |
| 5    | Carboplatin DB00958 human GSE13525 sample 3089         | 65/357  | $1.69 \times 10^{-26}$ | $3.07 \times 10^{-24}$ | 6.11       |
| 6    | NSC319726 5351307 human GSE35972 sample 2479           | 59/309  | $2.99 \times 10^{-25}$ | $4.52 \times 10^{-23}$ | 6.43       |
| 7    | Lipopolysaccharide 11970143 human GSE5504 sample 3483  | 72/468  | $1.29 \times 10^{-24}$ | $1.67 \times 10^{-22}$ | 5.01       |
| 8    | dasatinib DB01254 human GSE59357 sample 3306           | 57/298  | $1.81 \times 10^{-24}$ | $1.98 \times 10^{-22}$ | 6.43       |
| 9    | thapsigargin 446378 human GSE19519 sample 3236         | 66/399  | $1.97 \times 10^{-24}$ | $1.98 \times 10^{-22}$ | 5.43       |
| 10   | Y15 23627197 human GSE43452 sample 2554                | 64/390  | $1.59 \times 10^{-23}$ | $1.44 \times 10^{-21}$ | 5.37       |

**Table 7.** Drug perturbations from GEO up for A549 by DESeq2.

The proposed method can predict effective drugs for COVID-19 based on gene expression analysis, at least, comparatively to DESeq2. Nevertheless, DESeq2 has less significance and has a tendency to list the same compounds multiple times. The proposed method can identify more convincing and diverse candidate compounds than DESeq2.

Based on the overlap between human genes known to interact with SARS-CoV-2 proteins and selected genes (Fig. 14) and from the point of drug repositioning, TD-based unsupervised FE with optimized SD is, at least, competitive with DESeq2.

**Comparison of methods using multi-organ measurements with multiple drug treatments.** One might wonder if the proposed methods, TD- and PCA-based unsupervised FE with optimized SD, are applicable to a more complicated set-up. To investigate this point, we checked the case where multiple drugs are applied to mice whose gene expression of multiple tissues are measured, to which we applied TD-based unsupervised FE<sup>34</sup>.

*Enrichment of tissue-specific genes.* In the previous study<sup>34</sup>, although we applied TD-based unsupervised FE to gene expression profiles, there existed some problems. First of all, the number of genes selected was too small to have no false negatives. Using the optimized SD, the number of selected genes increased (Table 8; for more details, e.g., the definition of the four gene sets, neurons and testis, muscle, gastrointestinal 1 and 2, see the previous study<sup>34</sup>. This topic has not been discussed herein as it is not directly related to the comparison of the performance between the original TD-based unsupervised FE and that with the optimised SD. The full list of the selected genes is available in Data S9). Although an increased number of genes is meaningless if the biological reliability is less, the biological reliability of selected genes is also improved (lower panel of Fig. 15, which corresponds to a present study and is associated with a greater number of cell lines and tissue specificity than that in the upper panel of Fig. 15, which corresponds to a previous study). Thus, the employment of optimized SD is also effective to a more complicated data set than simple pairwise comparisons between the treated and control samples investigated in the previous sections.

*Coincidence with drug treatment.* We have also performed additional validation of the genes selected by TD-based unsupervised FE with optimized SD associated with adjusted *P* values less than 0.1 (Table 8, full list is available in Data S10–S13). We have uploaded selected genes to Enrichr<sup>36</sup> and evaluated the overlaps between

| Adjusted <i>P</i> values | TD-based unsupervised FE <sup>34</sup> |        | TD-based unsupervised FE with optimized SD |
|--------------------------|--|--------|--|
|                          | ≤ 0.01                                 | ≤ 0.01 | ≤ 0.1                                      |
| Neuron                   | 18                                     | 356    | 472  |
| Muscle                   | 51                                     | 547    | 663  |
| Gastrointestine 1        | 97                                     | 1026   | 1322                                       |
| Gastrointestine 2        | 128                                    | 574    | 722  |

**Table 8.** Comparison of selected genes between TD-based unsupervised FE<sup>34</sup> and optimal SD with multi-organ data sets.

the genes selected and those whose expression was altered with the treatment of the 15 drugs used in this study. Then, we found that all four gene sets in Table 8 had a significant overlap with the genes whose expression was altered with the treatment of 5 of the drugs (acetaminophen, cisplatin, clozapine, doxycycline, and olanzapine) in DrugMatrix, which does not include other drug treatments (Supplementary material). This suggests that TD-based unsupervised FE with optimal SD can correctly recognize drug treatments based on gene expression; this was impossible in the previous study<sup>34</sup> because of the very small number of genes selected (Table 8). Thus, considering the optimization of SD enables TD-based unsupervised FE to recognize a greater number of biologically reliable genes than the original TD-based unsupervised FE, which did not include the optimization of SD.

## Discussion

In this study, we have introduced the optimization of SD to TD- and PCA-based unsupervised FE and have improved their performance by increasing the identified DEGs associated with greater biological reliability. One of the striking features is that DEGs with lesser gene expression are less likely recognized even with the same LFC, if the genes are selected by TD- and PCA-based unsupervised FE with optimized SD. In DESeq2, the tendency that less expressed genes are hardly recognized as DEGs is artificially introduced by assuming dispersion relation, Eq. (1). Nevertheless, in PCA- and TD-based unsupervised FE, it is automatically introduced. Generally, there exists a relationship between difference,  $\Delta$  of two variables,  $x$  and  $y$ , and LFC as

$$\Delta \equiv x - y \quad (24)$$

$$\text{LFC} \equiv \log_2 \frac{x}{y} = \log_2 \left( 1 + \frac{\Delta}{y} \right) \quad (25)$$

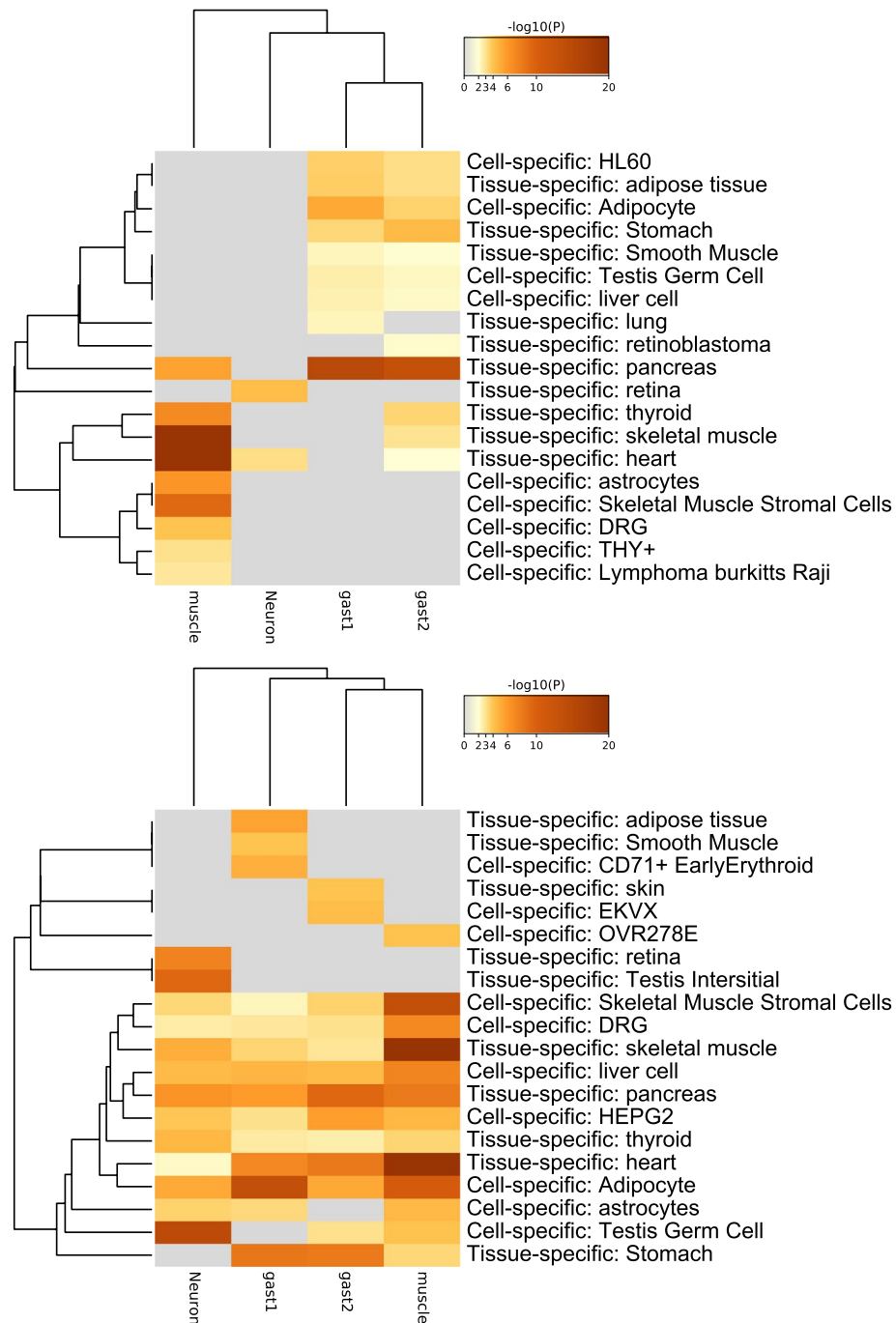
Then

$$\Delta = y(2^{\text{LFC}} - 1) \quad (26)$$

Because  $v_{2j}$  (Fig. 2B) corresponds to  $\Delta$ , if DEGs are identified using  $u_{2i}$  that corresponds to  $v_{2j}$  as in TD- and PCA-based unsupervised FE (see Eqs. (6) and (12)), DEGs associated with the same LFC are less likely selected for the smaller  $y$  that corresponds to  $\mu$ . This results in the distribution of DEGs in MAPlot (Fig. 8), where genes with the same LFC (vertical axis) are less likely identified as DEGs with smaller gene expression (horizontal axis). Figure 16 shows the MAPlot drawn using two independent random variables obeying the same positive uniform distribution; the red colored region associated with  $|\Delta|$  larger than some threshold values qualitatively represents the tendency that indicates that a smaller  $x + y$  is less likely selected even with the same LFC,  $\log_2 \frac{x}{y}$ . Thus, TD- and PCA-based unsupervised FE can introduce the tendency that genes with less expression are less likely to be DEGs, even with the same amount of LFC more naturally than DESeq2, which has to manually introduce a dispersion relation, Eq. (1).

In addition to this, although DESeq2 assumes NB distribution that does not have any rationalization other than that it takes only positive values and has a tunable mean as well as variance simultaneously, TD- and PCA-based unsupervised FE assume only that  $u_{\ell i}$  obeys the Gaussian distribution (Eqs. (6) and (12)), which is more reasonable because Gaussian distributions can generally appear when independent random variables are summed up. Actually, NOISeq does not assume NB distribution as well but achieves comparative performance with DESeq2 (Fig. 13). In this sense, TD- and PCA-based unsupervised FE can realize DEG distribution in an MAPlot more naturally than DESeq2.

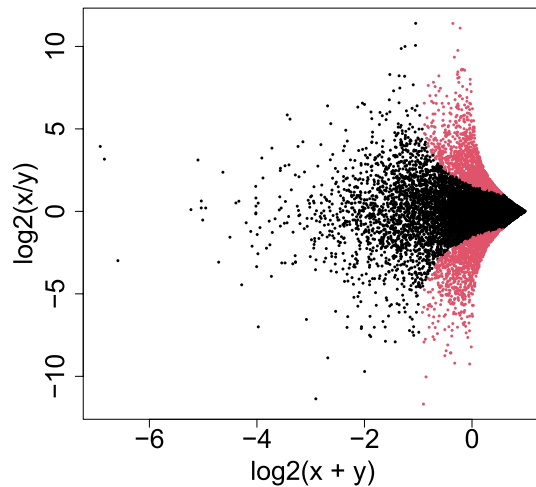
Another remarkable point of TD- and PCA-based unsupervised FE with optimized SD is that it does not have to screen for selected genes by LFC after the genes are selected using  $P$  values. As can be seen in Fig. 10, state-of-art methods, including DESeq2, often identify too many DEGs. In these circumstances, LFC is often used to reduce the number of DEGs. Nevertheless, Stupnikov et al<sup>37</sup> found that the coincidence of the selected genes among the various state-of-art methods drastically decreases if the genes selected based on  $P$  values are further screened with LFC. In this sense, TD- and PCA-based unsupervised FE with optimized SD are more promising methods than state-of-art methods that need screening by LFC to yield a reasonable number of DEGs.



**Figure 15.** Enrichment analysis of cell and tissue specificity with Metascape<sup>35</sup>. Upper: original TD-based unsupervised FE (using genes with adjusted  $P \leq 0.01$  in Table 8), lower; the present study with optimized SD (using genes with adjusted  $P \leq 0.1$  in Table 8). Metascape 3.5, <https://metascape.org/>.

Yet another advantage is that TD- and PCA-based unsupervised FE have already been applied to a wide range of problems. Not only can optimized SD improve the performance of PCA- and TD-based unsupervised FE, as can be seen in Figs. 14 and 15, but also the alteration is limited to the last stage, i.e.,  $P$  value computation, Eqs. (6) and (12). Thus, the optimized SD is expected to improve the performance in a wide range of problems, to which TD- and PCA-based unsupervised FE have been applied.

One might wonder if the validation should be based upon ground truth. Nevertheless, we do not think that there are ground truth for DEGs; DEGs are depend upon the definition of DEGs since the amount of differential expression is not discrete variable but continuous one. We need to decide threshold values for DEGs which affects which genes are DEGs. In contrast, biological significance is more trustable. In addition to this, the purpose of



**Figure 16.** “MAPlot” using two independent variables,  $x$  and  $y$ , drawn from uniform distribution  $\in [0, 1]$ . Red dots are associated with  $|x - y| > 0.5$ .

identification of DEGs is to further make use of them as biological studies. Thus, we believe that the proposed methods that can select biologically more reasonable genes than state-of-art methods is worthwhile publishing.

## Conclusions

In this study, we optimized SD to improve TD- and PCA-based unsupervised FE. As a result, not only the obtained DEGs increased and became reasonable in number but also the histogram of  $1-P$  became more reliable, i.e., more coincident with the null hypothesis that SVV and PC obey Gaussian distribution. In addition to this, TD- and PCA-based unsupervised FE provide reliable distribution of DEGs in MAPlot, i.e., less expressed genes are less likely selected as DEGs even if they are associated with the same LFC; this property was implemented manually by assuming dispersion relation, Eq.(1), in DESeq2. The biological reliability of the selected genes is also much better by this method than by other state-of-art methods. These points suggest that TD- and PCA-based unsupervised FE are superior than state-of-art methods in terms of achieving better performance with less assumption.

## Methods

Sample R code to perform analyses in this study is available as Data S14.

**Gene expression profiles.** *MAQC.* Seven human brain expression profiles were downloaded from SRA<sup>38</sup> (ID SRX016359), and seven UHR expression profiles were downloaded from SRA (ID SRX016367). Fourteen FASTQ files were mapped to the hg38 human genome using rapmap<sup>39</sup>. htseq-count<sup>40</sup> was used to convert the obtained bam files to count data files using the gtf file taken from [ftp://ftp.ensembl.org/pub/release-105/gtf/homo\\_sapiens/Homo\\_sapiens.GRCh38.105.gtf.gz](ftp://ftp.ensembl.org/pub/release-105/gtf/homo_sapiens/Homo_sapiens.GRCh38.105.gtf.gz).

*SEQC.* SEQC<sup>13</sup> were obtained from bioconductor<sup>41</sup> as an experimental package, seqc. It includes thirteen profiles shown in Fig. 11. For more details, see Vignettes in the seqc experimental package.

**The histogram composed of Gaussian distribution and outliers in Fig. 4.** The Gaussian part is one thousand values drawn from Gaussian distribution with zero mean and an SD of one. Outliers are 100 values, which are equal to 5.

**PCA-based unsupervised FE applied.** *MAQC.* Genes not expressed in any of the 14 samples have been excluded. Four rows having annotations “\_\_no\_feature”, “\_\_ambiguous”, “\_\_not\_aligned”, and “\_\_alignment\_not\_unique” have also been excluded. As a result, we got  $x_{ij} \in \mathbb{R}^{40933 \times 14}$ . The  $x_{ij}$  was processed as described in the main text.

*SEQC.* Regardless of which of the 13 data sets was considered, only those genes expressed in all samples were considered. An individual data set has a distinct number of rows (genes) and columns (samples). The  $x_{ij}$  obtained from an individual data set was processed as described in the main text.

*SARS-CoV-2.* All processes used were exactly the same as those described in the previous study<sup>14</sup>. After obtaining  $u_{5i}$ , the SD was optimized as described in the main text.

*Multi-organ.* All processes used were exactly the same as those described in the previous study<sup>34</sup>. After getting  $u_{li}$ , the SD was optimized as described in the main text.

| Cell lines     | Adjusted $P$ values $\leq 0.01$ | Alternative conditions                     | The number of DEG2 |
|----------------|---------------------------------|--|--------------------|
| Calu3          | 16432                           | Adjusted $P$ value $\leq 0.05$ , LFC > 2.0 | 340                |
| NHBE           | 327                             | Adjusted $P$ value $\leq 0.05$ , LFC > 0.5 | 171                |
| A549           |                                 |  |                    |
| MOI 0.2        | 15852                           | Adjusted $P$ value $\leq 0.05$ , LFC > 2.0 | 176                |
| MOI 2.0        | 7431                            | Adjusted $P$ value $\leq 0.05$ , LFC > 2.0 | 547                |
| ACE2 expressed | 7509                            | Adjusted $P$ value $\leq 0.05$ , LFC > 1.0 | 756                |

**Table 9.** The number of DEGs in SARS-CoV-2 study by DESeq2 (based on author-provided supplementary material).

**Optimization of SD.** At first, a histogram of  $1 - P_i$  was computed using `hclust` function in R with the “break=100” option. Then, an SD of the binned histogram, `hc$count` associated with `hc$breaks` less than 1- $P$  whose adjusted  $P$  value was less than threshold value  $P_0$ , was minimized using `optim` function in R. The R code has been provided in Data S14 to show how to optimize SD in an individual data set.

**Coincidence between PCA-based unsupervised FE and DESeq2.** The coincidence between PCA-based unsupervised FE and DESeq2 was evaluated by AUC (Fig. 9) as follows. At first, the top 1000 genes based on  $P$  values computed by DESeq2 were regarded positive and the remaining genes were regarded negative. Then,  $P$  values computed by PCA-based unsupervised FE were used to predict positive genes. Using this result, AUC was computed. Next, on the contrary, the top 1000 genes based on  $P$  values computed by PCA-based unsupervised FE were regarded positive and the remaining genes were negatives. Then,  $P$  values computed by DESeq2 were used to predict positive genes. Using this result, AUC was computed.

**Enrichment analyses.** Enrichment analyses were performed using either Metascape<sup>35</sup> or Enrichr<sup>36</sup> by uploading gene symbols. If the gene ID was not a gene symbol in individual data sets, the gene ID conversion tool in Database for Annotation, Visualization, and Integrated Discovery (DAVID)<sup>42,43</sup> was used for conversion.

**DEG identification of SARS-CoV-2 data by DESeq2.** We used author-provided adjusted  $P$  values and LFC (in supplementary data in their paper) to identify DEGs. If we considered only adjusted  $P$  values to identify DEGs, DESeq2 would identify too many genes (Table 9). Thus, we had to consider LFC as well. Table 9 shows the number of DEGs used in this study.

The evaluation of the overlap with human genes known to interact with SARS-CoV-2 proteins is available in Supplementary materials. The best one, that for the ACE2-expressed A549 cell line, is also included in the main text as Fig. 14.

## Data availability

The sequencing datasets are available via the NIH/NCBI Sequence Read Archive (SRA) repository using accession number SRX016359 and SRX016367, via bioconductor with the package of `seqc` [https://doi.org/doi:10.18129/B9.bioc.seqc, accessed 10th July 2022], via the NIH/NCBI Gene Expression Omnibus (GEO) repository using accession number GSE147507 and GSE142068.

Received: 7 March 2022; Accepted: 27 September 2022

Published online: 19 October 2022

## References

- Taguchi, Y-h. Comparative transcriptomics analysis. In *Encyclopedia of Bioinformatics and Computational Biology* (eds Ranganathan, S. *et al.*) 814–818 (Academic Press, 2019). <https://doi.org/10.1016/B978-0-12-809633-8.20163-5>.
- Rapaport, F. *et al.* Comprehensive evaluation of differential gene expression analysis methods for RNA-seq data. *Genome Biol.* **14**, 3158. <https://doi.org/10.1186/gb-2013-14-9-r95> (2013).
- Tusher, V. G., Tibshirani, R. & Chu, G. Significance analysis of microarrays applied to the ionizing radiation response. *Proc. Natl. Acad. Sci.* **98**, 5116–5121. <https://doi.org/10.1073/pnas.091062498> (2001).
- Ritchie, M. E. *et al.* Limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res.* **43**, e47. <https://doi.org/10.1093/nar/gkv007> (2015).
- Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550. <https://doi.org/10.1186/s13059-014-0550-8> (2014).
- Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139–140. <https://doi.org/10.1093/bioinformatics/btp616> (2009).
- McCarthy, D. J., Chen, Y. & Smyth, G. K. Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Res.* **40**, 4288–4297. <https://doi.org/10.1093/nar/gks042> (2012).
- Law, C. W., Chen, Y., Shi, W. & Smyth, G. K. voom: Precision weights unlock linear model analysis tools for RNA-seq read counts. *Genome Biol.* **15**, R29. <https://doi.org/10.1186/gb-2014-15-2-r29> (2014).
- Tarazona, S., Garcia, F., Ferrer, A., Dopazo, J. & Conesa, A. NOIseq: a RNA-seq differential expression method robust for sequencing depth biases. *EMBnet.journal* **17**, 18–19. <https://doi.org/10.14806/ej.17.B.265>
- Taguchi, Y-h. *Unsupervised Feature Extraction Applied to Bioinformatics* (Springer International Publishing, 2020).
- Shi, L. *et al.* The MicroArray quality control (MAQC) project shows inter- and intraplatform reproducibility of gene expression measurements. *Nat. Biotechnol.* **24**, 1151–1161. <https://doi.org/10.1038/nbt1239> (2006).



12. Mudge, J. F., Baker, L. F., Edge, C. B. & Houlahan, J. E. Setting an optimal  $\alpha$  that minimizes errors in null hypothesis significance tests. *PLoS ONE* **7**, 1–7. <https://doi.org/10.1371/journal.pone.0032734> (2012).
13. SEQC/MAQC-III Consortium, A comprehensive assessment of RNA-seq accuracy, reproducibility and information content by the sequencing quality control consortium. *Nature Biotechnology* **32**, 903–914. <https://doi.org/10.1038/nbt.2957> (2014).
14. Taguchi, Y.-H. & Turki, T. A new advanced in silico drug discovery method for novel coronavirus (SARS-CoV-2) with tensor decomposition-based unsupervised feature extraction. *PLoS ONE* **15**, 1–16. <https://doi.org/10.1371/journal.pone.0238907> (2020).
15. Taguchi, Y.-H. & Turki, T. Application of tensor decomposition to gene expression of infection of mouse hepatitis virus can identify critical human genes and effective drugs for SARS-CoV-2 infection. *IEEE J. Sel. Top. Signal Process.* **15**, 746–758. <https://doi.org/10.1109/JSTSP.2021.3061251> (2021).
16. Zhao, H., Mendenhall, M. & Deininger, M. W. Imatinib is not a potent anti-SARS-CoV-2 drug. *Leukemia* **34**, 3085–3087. <https://doi.org/10.1038/s41375-020-01045-9> (2020).
17. Naidoo, D., Roy, A., Kar, P., Mutanda, T. & Anandraj, A. Cyanobacterial metabolites as promising drug leads against the mpro and ppro of SARS-CoV-2: An in silico analysis. *J. Biomol. Struct. Dyn.* **39**, 6218–6230. <https://doi.org/10.1080/07391102.2020.1794972> (2021).
18. Dorobisz, K., Dorobisz, T., Janczak, D. & Zatoński, T. Doxycycline in the coronavirus disease 2019 therapy. *Ther. Clin. Risk Manag.* **17**, 1023–1026. <https://doi.org/10.2147/tcrm.s314923> (2021).
19. Gimeno, A. *et al.* Prediction of novel inhibitors of the main protease (M-pro) of SARS-CoV-2 through consensus docking and drug reposition. *Int. J. Mol. Sci.* **21**, 3793. <https://doi.org/10.3390/ijms21113793> (2020).
20. Jamal, Q. M. S., Alharbi, A. H. & Ahmad, V. Identification of doxorubicin as a potential therapeutic against SARS-CoV-2 (COVID-19) protease: a molecular docking and dynamics simulation studies. *J. Biomol. Struct. Dyn.* **40**, 7960–7974. <https://doi.org/10.1080/07391102.2021.1905551> (2021).
21. MotieGhader, H., Safavi, E., Rezapour, A., Amoodizaj, F. F. & asl Iranifam, R. Drug repurposing for coronavirus (SARS-CoV-2) based on gene co-expression network analysis. *Sci. Rep.* **11**, 21872. <https://doi.org/10.1038/s41598-021-01410-3> (2021).
22. Mansouri, A., Kowsar, R., Zakariazadeh, M., Hakimi, H. & Miyamoto, A. The impact of calcitriol and estradiol on the SARS-CoV-2 biological activity: A molecular modeling approach. *Sci. Rep.* **12**, 717. <https://doi.org/10.1038/s41598-022-04778-y> (2022).
23. Zu, S. *et al.* Tamoxifen and clomiphene inhibit SARS-CoV-2 infection by suppressing viral entry. *Signal Transduct. Targeted Therapy* **6**, 435. <https://doi.org/10.1038/s41392-021-00853-4> (2021).
24. Zhu, W. *et al.* Identification of SARS-CoV-2 3cl protease inhibitors by a quantitative high-throughput screening. *ACS Pharmacol. Transl. Sci.* **3**, 1008–1016. <https://doi.org/10.1021/acspsci.0c00108> (2020).
25. Paragas, J., Blatt, L. M., Hartmann, C., Huggins, J. W. & Endy, T. P. Interferon alfacon1 is an inhibitor of SARS-corona virus in cell-based models. *Antiviral Res.* **66**, 99–102. <https://doi.org/10.1016/j.antiviral.2005.01.002> (2005).
26. Ripon, M. A. R., Bhowmik, D. R., Amin, M. T. & Hossain, M. S. Role of arachidonic cascade in covid-19 infection: A review. *Prostaglandins Other Lipid Mediators* **154**, 106539. <https://doi.org/10.1016/j.prostaglandins.2021.106539> (2021).
27. Chowdhury, T., Roymahapatra, G. & Mandal, S. M. In silico identification of a potent arsenic based approved drug darinaparsin against sars-cov-2: Inhibitor of RNA dependent RNA polymerase (RdRp) and necessary proteases. *ChemRxiv*. <https://doi.org/10.26434/chemrxiv.12200495.v1> (2020).
28. Clemente-Moragón, A. *et al.* Metoprolol in critically ill patients with COVID-19. *J. Am. Coll. Cardiol.* **78**, 1001–1011. <https://doi.org/10.1016/j.jacc.2021.07.003> (2021).
29. The RECOVERY Collaborative Group, Dexamethasone in hospitalized patients with covid-19. *N. Engl. J. Med.* **384**, 693–704. <https://doi.org/10.1056/nejmoa2021436> (2021).
30. Parthasarathy, H., Tandel, D. & Harshan, K. H. Metformin suppresses SARS-CoV-2 in cell culture. *bioRxiv*. <https://doi.org/10.1101/2021.11.18.469078> (2021).
31. Salesi, M., Shojaie, B., Farajzadegan, Z., Salesi, N. & Mohammadi, E. TNF- $\alpha$  blockers showed prophylactic effects in preventing COVID-19 in patients with rheumatoid arthritis and seronegative spondyloarthropathies: A case-control study. *Rheumatol. Therapy* **8**, 1355–1370. <https://doi.org/10.1007/s40744-021-00342-8> (2021).
32. Petruk, G. *et al.* SARS-CoV-2 spike protein binds to bacterial lipopolysaccharide and boosts proinflammatory activity. *J. Mol. Cell Biol.* **12**, 916–932. <https://doi.org/10.1093/jmcb/mjaa067> (2020).
33. Pasquereau, S. *et al.* Resveratrol inhibits HCoV-229E and SARS-CoV-2 coronavirus replication in vitro. *Viruses* **13**, 354. <https://doi.org/10.3390/v13020354> (2021).
34. Taguchi, Y.-h. & Turki, T. Universal nature of drug treatment responses in drug-tissue-wide model-animal experiments using tensor decomposition-based unsupervised feature extraction. *Front. Genet.* **11**, 695. <https://doi.org/10.3389/fgene.2020.00695> (2020).
35. Zhou, Y. *et al.* Metascape provides a biologist-oriented resource for the analysis of systems-level datasets. *Nat. Commun.* **10**, 1523. <https://doi.org/10.1038/s41467-019-09234-6> (2019).
36. Xie, Z. *et al.* Gene set knowledge discovery with Enrichr. *Curr. Protocols* **1**, e90. <https://doi.org/10.1002/cpz1.90> (2021).
37. Stupnikov, A. *et al.* Robustness of differential gene expression analysis of RNA-seq. *Comput. Struct. Biotechnol. J.* **19**, 3470–3481. <https://doi.org/10.1016/j.csbj.2021.05.040> (2021).
38. Leinonen, R., Sugawara, H. & Shumway, M. On behalf of the international nucleotide sequence database collaboration, the sequence read archive. *Nucleic Acids Res.* **39**, D19–D21. <https://doi.org/10.1093/nar/gkq1019> (2010).
39. Srivastava, A., Sarkar, H., Gupta, N. & Patro, R. RapMap: A rapid, sensitive and accurate tool for mapping RNA-seq reads to transcriptomes. *Bioinformatics* **32**, i192–i200. <https://doi.org/10.1093/bioinformatics/btw277> (2016).
40. Putri, G. H., Anders, S., Pyl, P. T., Pimanda, J. E. & Zanini, F. Analysing high-throughput sequencing data in python with htseq 2.0. *Bioinformatics* **38**, 2943–2945. <https://doi.org/10.1093/bioinformatics/btac166> (2022).
41. Huber, W. *et al.* Orchestrating high-throughput genomic analysis with bioconductor. *Nat. Methods* **12**, 115–121. <https://doi.org/10.1038/nmeth.3252> (2015).
42. Huang, D. W., Sherman, B. T. & Lempicki, R. A. Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. *Nat. Protoc.* **4**, 44–57. <https://doi.org/10.1038/nprot.2008.211> (2008).
43. Huang, D. W., Sherman, B. T. & Lempicki, R. A. Bioinformatics enrichment tools: Paths toward the comprehensive functional analysis of large gene lists. *Nucleic Acids Res.* **37**, 1–13. <https://doi.org/10.1093/nar/gkn923> (2008).

## Acknowledgements

This work was supported by the Japan Society for the Promotion of Science, KAKENHI [Grant numbers 19H05270, 20K12067, and 20H04848] to YHT.

## Author contributions

Y.H.T. planned the research and performed analyses. Y.H.T. and T.T. evaluated the results, discussions, and outcomes and drafted and reviewed the manuscript.

## Competing interests

The authors declare no competing interests.

### Additional information

**Supplementary Information** The online version contains supplementary material available at <https://doi.org/10.1038/s41598-022-21474-z>.

**Correspondence** and requests for materials should be addressed to Y.T.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022