

RESEARCH

# Tensor decomposition- and principal component analysis-based unsupervised feature extraction to select more reasonable differentially expressed genes: Optimization of standard deviation versus state-of-art methods

Y-h. Taguchi<sup>1\*</sup> and Turki Turki<sup>2</sup>

## Abstract

**Background:** 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.

**Results:** 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.

## Conclusions:

Tensor decomposition- and principal component analysis-based unsupervised feature extraction are perhaps better than state-of-art methods in regard to predicting differentially expressed genes because they achieve the desired property that the less expressed differentially expressed genes should be less likely selected or even associated with the same amount of logarithmic fold change, although they assume neither negative binomial distribution nor dispersion relation, which is usually assumed in state-of-art methods.

**Keywords:** tensor decomposition; principal component analysis; feature extraction; standard deviation; differentially expressed genes

## 1 Background

2 Identifying differentially expressed genes (DEGs) on  
 3 the basis of comparative analyses [1, 2] has always  
 4 been difficult. This challenge is attributable to mul-  
 5 tiple reasons; however, the primary reason is it be-  
 6 ing a *large p small n* problem. In a *large p small n*  
 7 problem, it is difficult to select features based on sta-  
 8 tistical criteria because a small number of samples  
 9 ( $= n$ ) have a tendency to lead to low significance;  
 10 in reality, the obtained  $P$ -values must be heavily cor-  
 11 rected by considering a large number of features ( $= p$ ).  
 12 This makes it difficult to find features with signifi-  
 13 cance. To resolve this difficulty, many methods spe-  
 14 cific to gene expression analysis have been proposed.  
 15 For example, significant analysis microarray (SAM) [3]  
 16 adds a small amount of constancy to gene expression,  
 17 thereby avoiding the misidentification of low expressed  
 18 genes as DEGs. Limma [4] applied a Bayesian strategy  
 19 to logarithmic gene expression. After high-throughput  
 20 sequencing (HTS) became popular,  $P$ -values are at-  
 21 tributed to individual genes, assuming that gene ex-  
 22 pression follows a negative binomial (NB) distribu-  
 23 tion [5, 6], which is one of the simplest positively val-  
 24 ued distributions with a tunable mean and variance. In  
 25 addition to this, the so-called dispersion relation [5, 6],

$$28 \quad \frac{\alpha(\mu)}{\mu^2} = \alpha_0 + \frac{\alpha_1}{\mu}, \quad (1)$$

29 has also been assumed, where  $\mu$  and  $\alpha$  are the mean  
 30 and variance, respectively, and  $\alpha_0$  and  $\alpha_1$  are regres-  
 31 sion coefficients; to our knowledge, eq. (1) is purely em-  
 32 pirical and lacks rationalization. Despite these difficul-  
 33 ties, many proposed state-of-art methods [5, 6, 7, 8, 9]  
 34 have been widely employed and used in various stud-  
 35 ies.

36 Contrary to these empirical methods, we proposed  
 37 tensor decomposition (TD)- and principal component  
 38 analysis (PCA)-based unsupervised feature extraction  
 39 (FE) [10] that only assumes that principal component  
 40 (PC) and singular value vectors (SVVs) obey Gaussian  
 41 distribution. Despite this simplicity, TD- and PCA-  
 42 based unsupervised FE have been successfully applied  
 43 to a wide range of genomic analyses. However, there  
 44 have been two problems: 1. The histogram of the  $P$ -  
 45 values is not fully coincident with the null hypothesis  
 46 that PC and SVV obey Gaussian distribution and 2.  
 47 The number of genes selected is too small to have no  
 48 false negatives. In this paper, we have shown that the

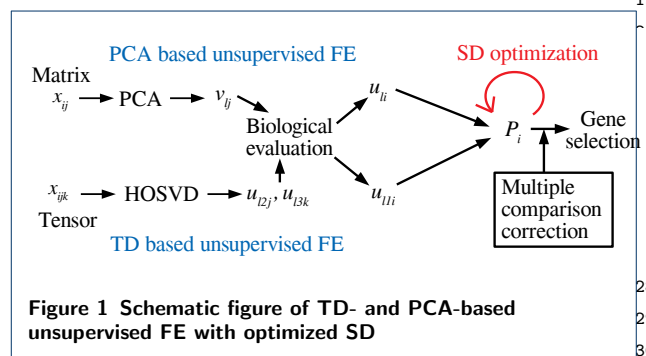
1 optimization of standard deviation (SD) in Gaussian  
 2 distribution can resolve these problems.

3 We tried optimizing SD for PCA-based unsuper-  
 4 vised FE and applied this to two highly curated data  
 5 sets—MAQC and SEQC. Then, we tested the opti-  
 6 mization of SD for TD-based unsupervised FE and  
 7 applied it to two more realistic problems: 1. drug repo-  
 8 sitioning for SARS-CoV-2 and 2. the analysis of gene  
 9 expression of multiple organs treated with multiple  
 10 drugs, to which TD-based unsupervised FE without  
 11 SD optimization was already applied.

## 12 Results

### 13 Outlines of TD and PCA based unsupervised FE

14 In this section, we have briefly explained the algorithm  
 15 of PCA- and TD-based unsupervised FE (Fig. 1) be-  
 16 fore explaining how we could improve them. When



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

$$36 \quad \sum_i x_{ij} = 0 \quad (2)$$

$$37 \quad \sum_i x_{ij}^2 = N, \quad (3)$$

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

$$40 \quad \sum_{i'} \left( \sum_j x_{ij}x_{i'j} \right) u_{\ell i'} = \lambda_{\ell} u_{\ell i} \quad (4)$$

41 where  $u_{\ell i} \in \mathbb{R}^{N \times N}$  is the  $\ell$ th PC score attributed to  
 42 gene  $i$ . The  $\ell$ th PC loading attributed to the  $j$ th sam-  
 43 ple can be computed as

$$44 \quad v_{\ell j} = \sum_i x_{ij} u_{\ell i} \in \mathbb{R}^{M \times M}. \quad (5)$$

51 \*Correspondence: tag@granular.com

52 <sup>1</sup>Department of Physics, Chuo University, 1-13-27 Kasuga, Bunkyo-ku,  
 53 112-8551 Tokyo, JAPAN

54 <sup>2</sup>Full list of author information is available at the end of the article

55 <sup>†</sup>Equal contributor

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

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

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,

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

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

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) [10]. 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 [10]; 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:

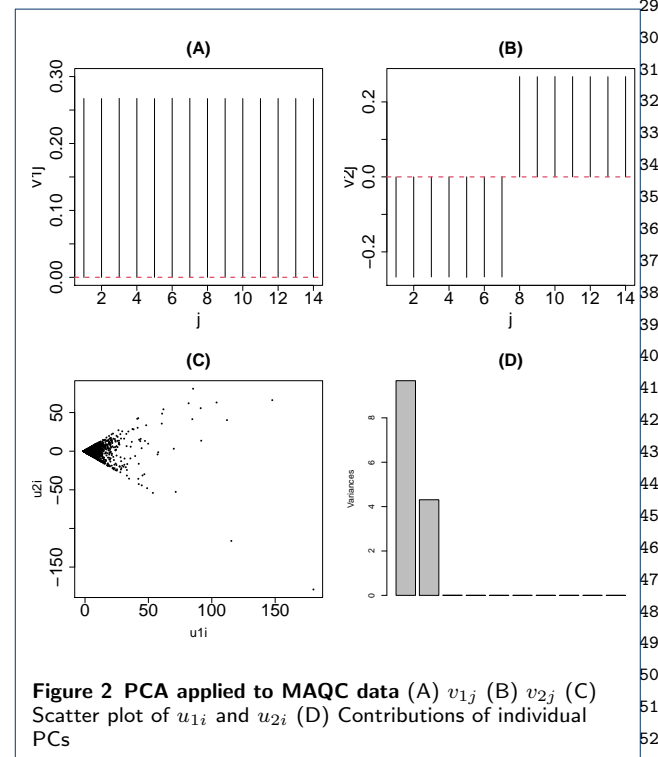
- Too small a number of genes were selected to have no false negatives.
- The histogram of  $P_i$  did not fully obey the null assumption that  $u_{\ell i}$  and  $u_{\ell_1 i}$  obey the Gaussian distribution.

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 [11] data set, which has been carefully curated and frequently used for benchmark studies. Figure 2C shows



**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

a scatter plot of genes using  $u_{1i}$  and  $u_{2i}$ . Figure 2A

1 and B show the PC loading  $v_{1j}$  and  $v_{2j}$ ;  $v_{1j}$  represents  
 2 the mean gene expression and  $v_{2j}$  represents the dif-  
 3 ferential expression between universal human reference  
 4 (UHR) and brain. Occasionally, this reminds us of the  
 5 horizontal and vertical axes of an MAPlot; the horizon-  
 6 tal axis of an MAPlot represents the mean expression  
 7 of individual genes, typically the mean logarithmic ex-  
 8 pression,  
 9

$$10 \quad \frac{1}{M} \sum_{j=1}^M \log_2 x_{ij} \quad (15)$$

14 whereas the vertical axis of an MAPlot represents the  
 15 differential expression between the two classes, typi-  
 16 cally the mean logarithmic fold change (LFC),  
 17

$$18 \quad \frac{1}{M_A} \sum_{j \in A} \log_2 x_{ij} - \frac{1}{M_B} \sum_{j \in B} \log_2 x_{ij} \quad (16)$$

21 where  $M_A$  and  $M_B (= M - M_A)$  are sample numbers  
 22 within one of the two classes, A and B, respectively,  
 23 and summations are taken within individual classes.  
 24 As can be seen in Fig. 2D, which represents the contri-  
 25 bution of PC loading,  $x_{ij}$  can be expressed almost fully  
 26 in the 2-dimensional space spanned by the first two  
 27 PCs. Thus, PCA can derive, in a fully unsupervised  
 28 manner, something that qualitatively corresponds to  
 29 an MAPlot (Fig.8), which is usually drawn artificially.  
 30 In spite of that, unfortunately, the genes selected by  
 31 the adjusted  $P_i$  are too small to have no false negatives  
 32 (Table 3) and an histogram of  $P_i$  is hardly regarded to  
 33 obey the null hypothesis; the left panel of Fig. 3 shows  
 34 the histogram of  $1 - P_i$ , where  $P_i$ s were computed from  
 35  $u_{2i}$  by eq. (6) using  $\sigma_2$  defined as  
 36

$$37 \quad \sigma_2 = \sqrt{\frac{1}{N} \sum_i (u_{2i} - \langle u_{2i} \rangle)^2} \quad (17)$$

$$38 \quad \langle u_{2i} \rangle = \frac{1}{N} \sum_i u_{2i}. \quad (18)$$

43 If  $1 - P_i$  is coincident with the null hypothesis; the  
 44 histogram of  $1 - P_i < 1$  should have a flat distribution  
 45 and that of  $1 - P_i \sim 1$  should have a sharp peak.  
 46

47 *Top ranked genes are coincident with DESeq2*

48 To understand the problem of  $P_i$ s computed by PCA-  
 49 based unsupervised FE, we compared  $P_i$ s computed  
 50 by PCA-based unsupervised FE with those computed  
 51 by DESeq2, a state-of-art method. At first, AUC was  
 52 computed to predict the top 1000 genes based on  $P_i$   
 53 derived with DESeq2 using  $P_i$ s computed by PCA-based  
 54 unsupervised FE; the area under the curve (AUC) was  
 55

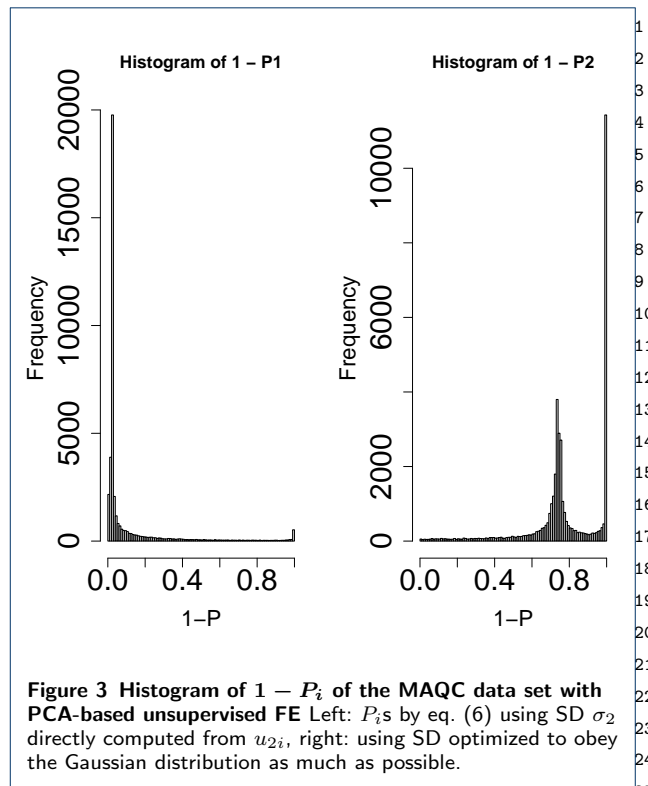


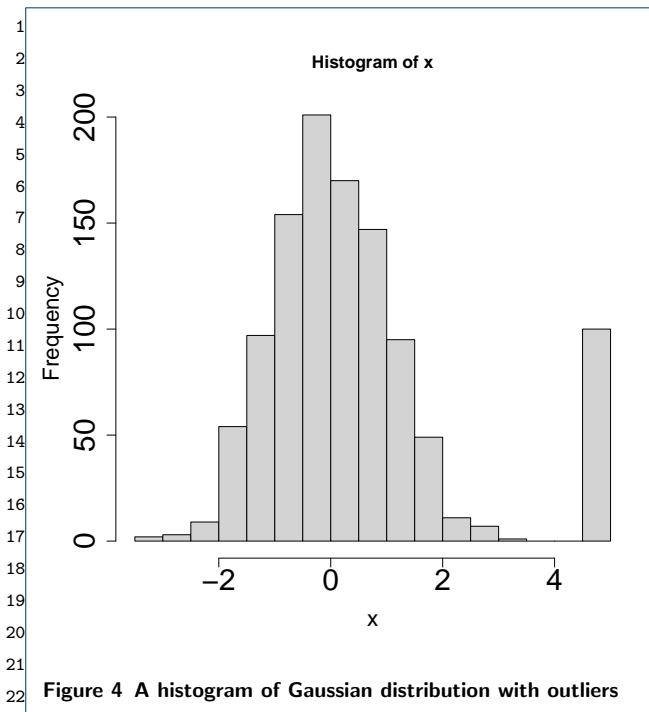
Figure 3 Histogram of  $1 - P_i$  of the MAQC data set with PCA-based unsupervised FE Left:  $P_i$ s by eq. (6) using SD  $\sigma_2$  directly computed from  $u_{2i}$ , right: using SD optimized to obey the Gaussian distribution as much as possible.

0.97. Next, in contrast, the AUC was computed to pre-  
 28 dict the top 1000 genes based on  $P_i$  derived with PCA-  
 29 based unsupervised FE using  $P_i$ s computed using DE-  
 30 Seq2; the AUC was 0.98. This indicated that the top-  
 31 ranked genes were suitably shared between PCA-based  
 32 unsupervised FE and DESeq2. Thus, the problem of  
 33 PCA-based unsupervised FE is not the genes' ranking  
 34 but the absolute value of  $P_i$ s.  
 35

36 *Optimization of SD*

37 Based on the observations at the end of the subsub-  
 38 section, we arrived at optimizing  $\sigma_\ell$  such that  $u_{\ell i}$  and  
 39  $u_{\ell_1 i}$  obeyed the Gaussian distribution. Generally, opti-  
 40 mizing SD to be fitted to the null hypothesis is not  
 41 easy. For example, Mudge et al [12] had to assume the  
 42 equivalence between Type I and II errors, which we  
 43 cannot assume because of an imbalance of numbers  
 44 between DEGs and the other genes; typically, DEGs  
 45 are expected to be minorities. Next, we decided to em-  
 46 ploy an alternative and more empirical approach. To  
 47 visualize the idea, we have shown some illustrative ex-  
 48 amples. Figure 4 shows a histogram of the variable  $x_i$   
 49 derived from the Gaussian distribution and outliers. If  
 50 we attribute the  $P$ -values to the  $i$ th variable with  $x_i$   
 51  
 52

$$53 \quad P_i = P_{\chi^2} \left[ > \left( \frac{x_i}{\sigma} \right)^2 \right] \quad (19)$$



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

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

$$\langle x_i \rangle = \frac{1}{N} \sum_{i=1}^N x_i \quad (21)$$

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.

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

	True	not outliers	outliers
predicted adjusted $P$ -values $> 0.01$		1000	100
predicted adjusted $P$ -values $\leq 0.01$		0	0

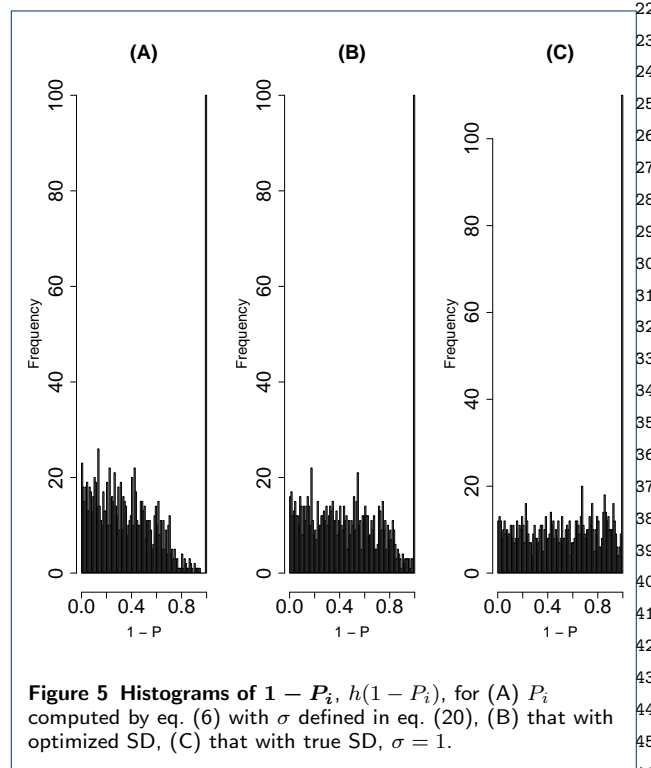
We computed the histogram of  $1 - P_i$ , Fig. 5A, which is far being idealized, Fig. 5C, that should have a constant histogram  $h(1 - P_i)$  up to  $1 - P_i$  very close to 1 and has one with a narrow peak near  $1 - P_i \sim 1$ . To optimize the SD, we tried to find an optimal SD such that the histogram for those not recognized as outliers was as flat as possible, i.e., obeying the null

hypothesis of the Gaussian distribution; we decided<sup>1</sup> to find the optimal SD that results in the most flat<sup>2</sup>  $h(1 - P_i)$  for  $1 - \text{adjusted } P_i$  less than threshold value<sup>3</sup>  $1 - \text{adjusted } P_0$  (adjusted  $P_0$  should be small enough).<sup>4</sup> To minimize the SD of binned  $h_i = h(1 - P_i)$ ,  $\sigma_h$ ,<sup>5</sup>

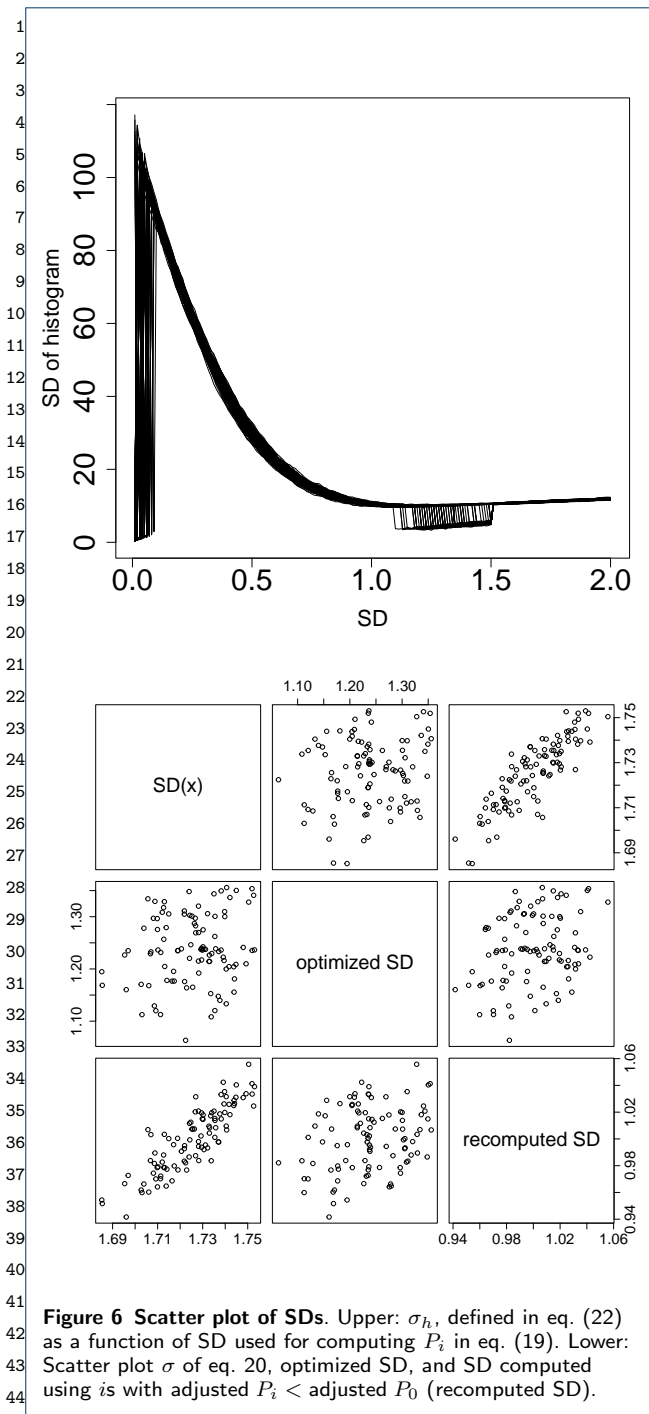
$$\sigma_h = \sqrt{\frac{\sum_{\text{adjusted } P_i < \text{adjusted } P_0} (h_i - \langle h_i \rangle)^2}{N(\text{adjusted } P_0)}} \quad (22)$$

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

with respect to  $\sigma$ , where  $N(\text{adjusted } P_0)$  is the number<sup>14</sup> of  $i$ s associated with adjusted  $P_i > \text{adjusted } P_0$ , i.e.,<sup>15</sup> not recognized as outliers and recognized as a part<sup>16</sup> of the Gaussian distribution. After optimizing  $\sigma$ , we<sup>17</sup> recomputed  $P_i$ . Fig. 5A and 5B show the histogram of<sup>18</sup>  $1 - P_i$  using  $\sigma = 1.75$  and optimized SD, respectively;<sup>19</sup> the latter is closer to an idealized histogram of  $P_i$ , Fig.<sup>20</sup> 5C, than the former.<sup>21</sup>



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  $i$ s 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



**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  $P_i < adjusted P_0$  (recomputed SD).

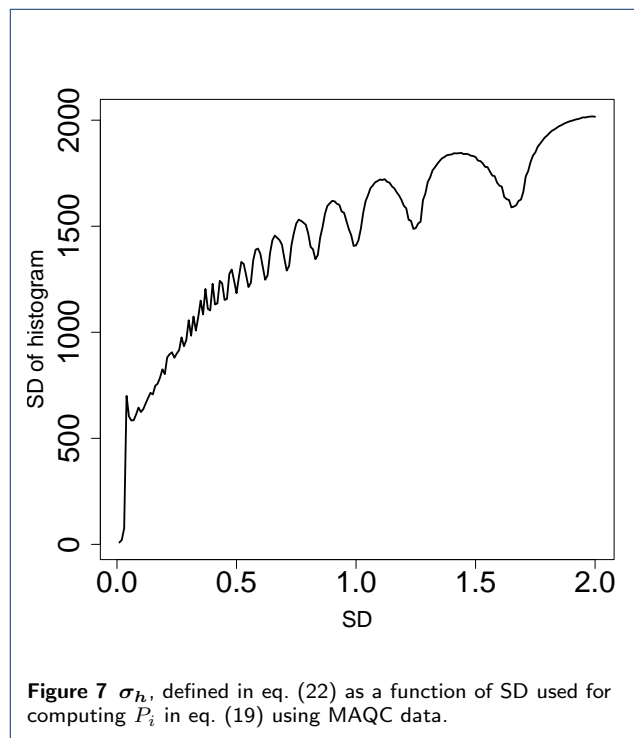
fact that SD computed using  $is$  for adjusted  $P_i < adjusted P_0$ , which is expected to correspond to the Gaussian distribution part in Fig. 4, is almost 1 helps justify our optimization procedure (Fig. 6, lower panel). The reason why  $SD = 0$  with  $\sigma_h = 0$  in the upper panel of Fig. 6 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).

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

	True	not outliers	outliers
predict adjusted $P$ -values $> 0.01$	1000	0	0
predict adjusted $P$ -values $\leq 0.01$	0	100	0

Next, we applied this strategy to the MAQC data set. Figure 7 shows  $\sigma_h$ , defined in eq. (22), as a func-



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

tion 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 recomputed using  $is$  with adjusted  $P_i < adjusted P_0$ , 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 Additional file 1), the original PCA-based unsupervised FE, than by using optimal SD (list of genes available as Additional file 2). 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.

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

	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

*Less expressed genes are less likely to be DEGs*

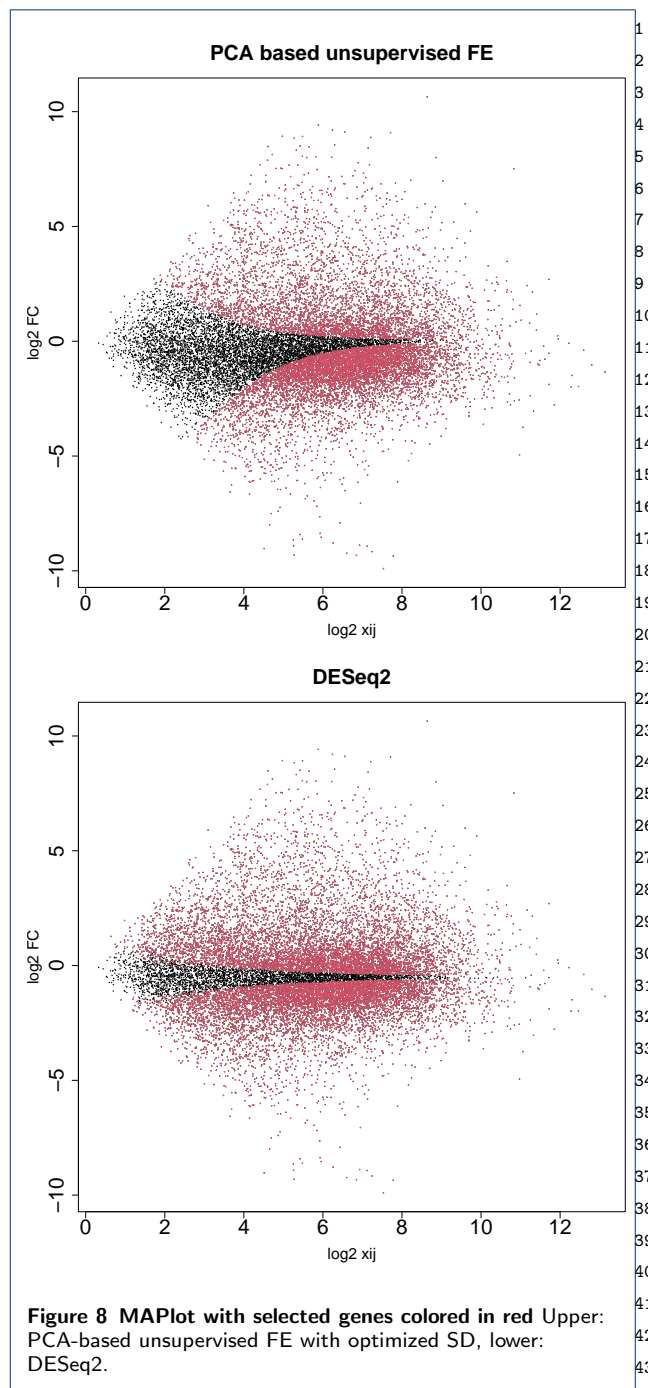
Figure 8 shows the selected genes in MAPlot. Although we assumed neither NB distribution nor dispersion relation, eq. (1), the distribution of selected genes in the MAPlot is reasonable; genes with the same LFC (vertical axis) are less likely selected when associated with smaller mean expression (horizontal 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 [13], 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*

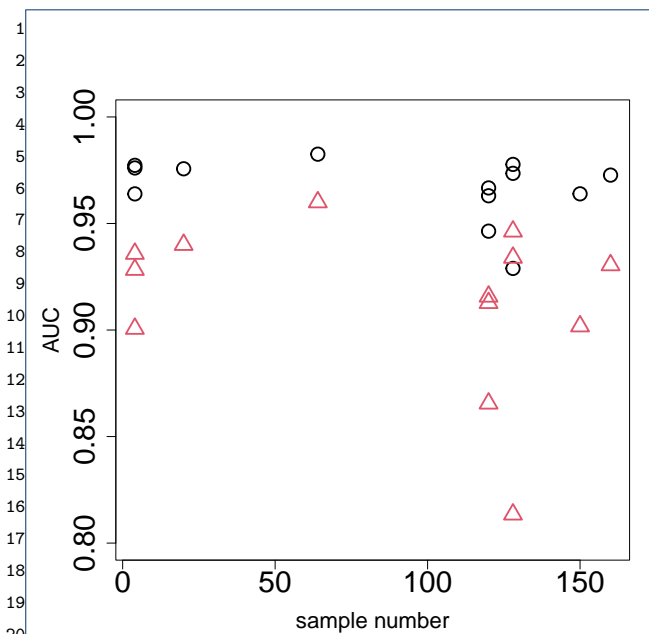
Based on the above results, PCA-based unsupervised FE is seemingly better than DESeq2. Nonetheless, PCA-based unsupervised FE can select a reasonable number of genes regardless of sample numbers (Fig. 10), and less expressed genes are unlikely to be DEGs when genes are selected by PCA-based unsupervised FE with optimized SD (Figs. 8 and 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 by PCA-based unsupervised FE



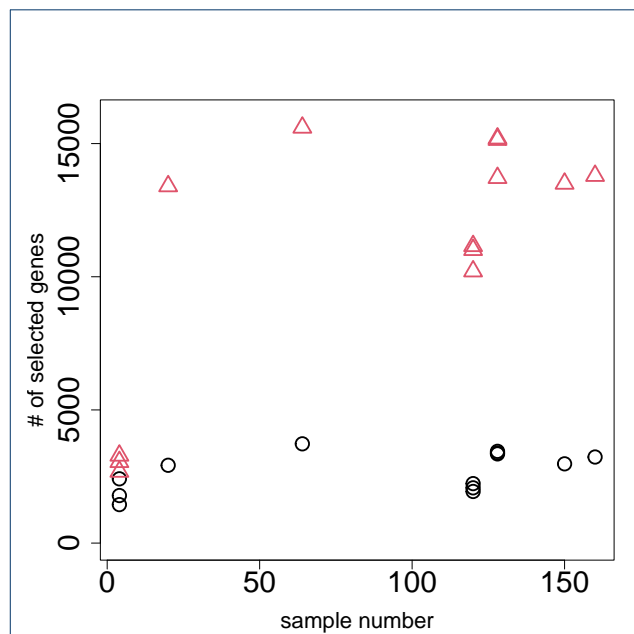
**Figure 8** MAPlot with selected genes colored in red Upper: PCA-based unsupervised FE with optimized SD, lower: DESeq2.

were better than those selected by DESeq2 (Full list of enrichment analysis is available in Additional files 1 and 2).

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 [6] (full list of enrichment analysis available in Additional file 3), voom [8] (full list of enrichment analysis available in



**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.

Additional file 4), and NOISEq [9] (full list of enrichment analysis available in Additional file 5); 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 outperform DESeq2 in more realistic data sets, we considered the drug repositioning of SARS-CoV-2, to which we applied TD-based unsupervised FE [14] and its kernelized version [15].

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  $u_{li}$  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 Additional file 6),

which is a much higher number of genes than 163 genes than that selected in previous studies [14, 15].

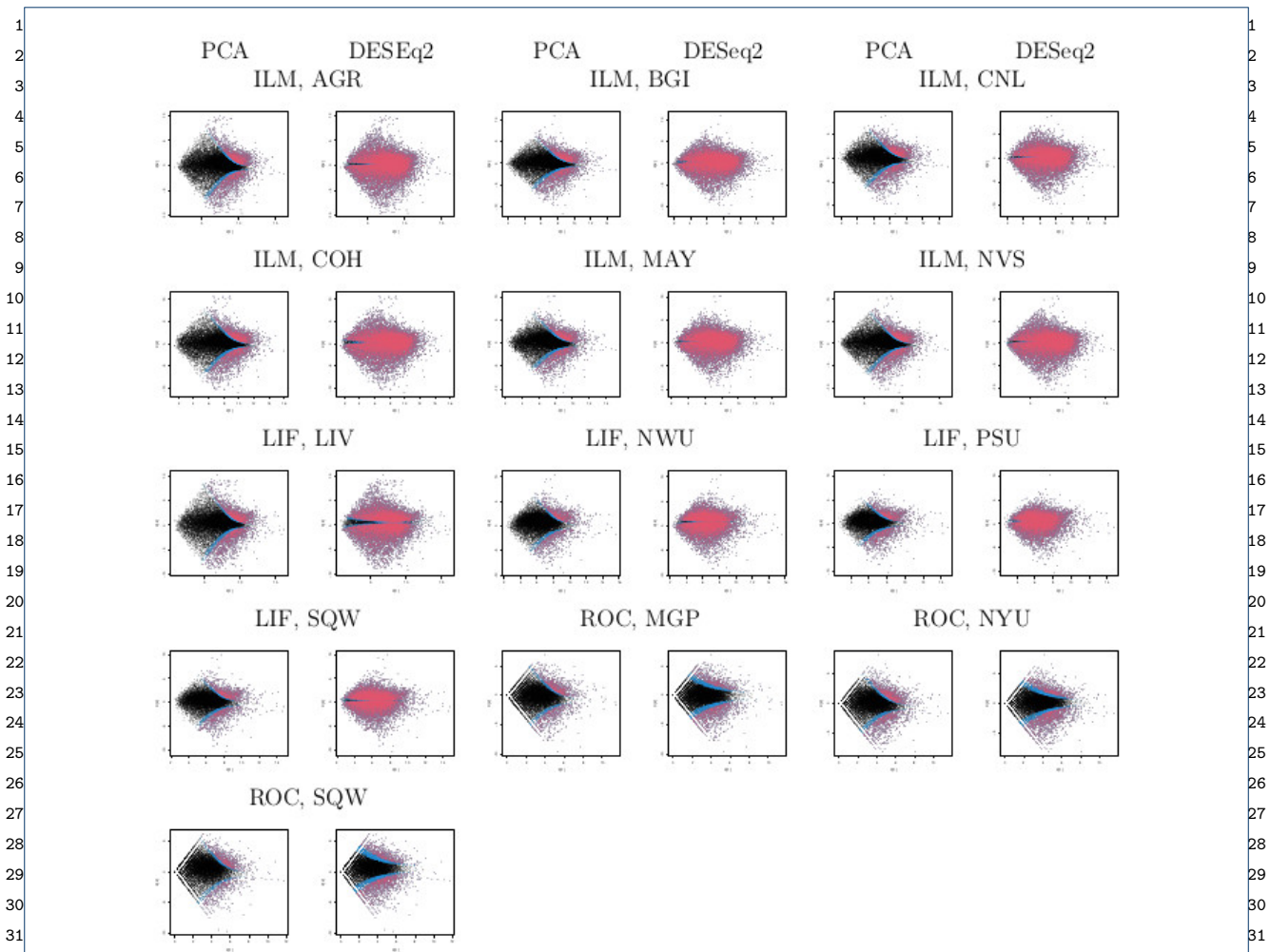
### Overlap with human genes known to interact with SARS-CoV-2 protein

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 [14, 15] (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 Additional File 7). 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 Additional file 6). The first one, imatinib, was once identified as a promising drug toward COVID-19, although it was rejected later [16]. The second one, apratoxin A, was reported to be a promising compound based on its protein binding affinity [17]. The third and fourth one, doxycycline, was supposed to be a promising drug toward COVID-19 [18]. The seventh one, trovafloxacin, was reported to be a promising compound based on its





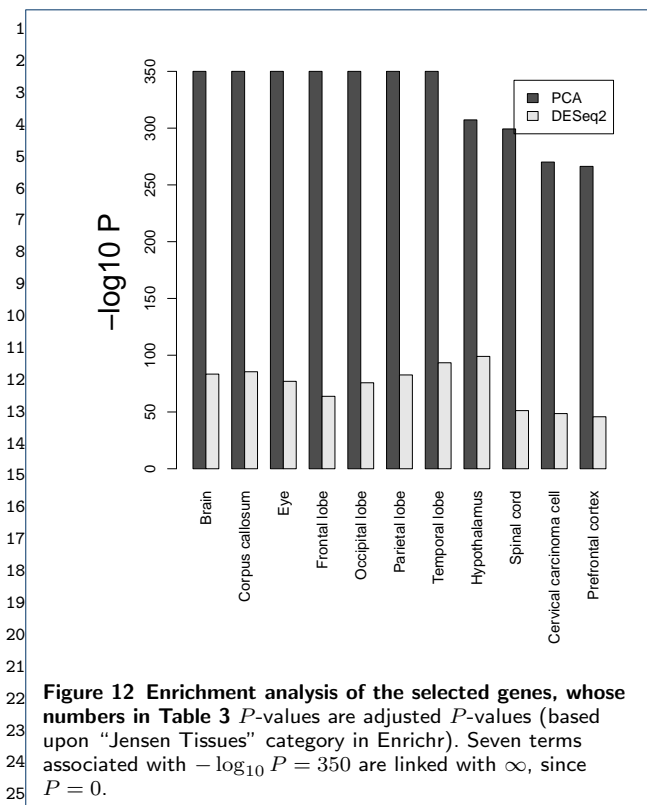
**Figure 11** MAPlot for SEQC PCA-based unsupervised FE with optimized SD: the first, third, and fifth columns, DESeq2: the second, fourth, and sixth columns. Three character IDs represent platform and sites. Blue: genes associated with adjusted  $P$ -values less than 0.1 but greater than 0.01. Red: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

protein binding affinity [19]. The eighth one, doxorubicin, was also reported to be a promising compound based on its protein binding affinity [20]. The ninth one, cisplatin, and the tenth one, carboplatin, were proposed as a result of drug repositioning [21]. 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 [22]. The second one, tamoxifen, was reported to inhibit SARS-CoV-2 infection by suppressing viral entry [23]. 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 [24], although its efficiency was limited to 40 %. The sixth one, IFN-alphacon1, was reported to be an inhibitor of SARS-CoV [25] but not

for SARS-CoV-2. The seventh one, arachidonic acid, was generally expected to inhibit SARS-CoV-2 infection [26]. The eighth one, arsenic, was also generally expected to act against the RdRp of coronavirus [27]. The ninth one, metoprolol, was reported to be a promising drug toward COVID-19 [28]. 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 Additional file 8), 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. [29], The seventh one, metformin, suppressed SARS-CoV-2 in cell culture [30]. The eighth one, etanercept, significantly decreased the risk of developing COVID-19 in patients with rheuma-



28 toid arthritis or spondyloarthropathies [31]. The tenth  
 29 one, lipopolysaccharide, is not a compound but a bac-  
 30 terial protein reported to bind to the SARS-CoV-2  
 31 spike protein [32].

32 See Table 7. The first and fourth one, resveratrol, in-  
 33 hibits HCoV-229E and SARS-CoV-2 coronavirus repli-  
 34 cation in vitro [33]. The second, third, and fifth one,  
 35 carboplatin, was proposed as a result of drug repo-  
 36 sitioning [21]. The seventh one, lipopolysaccharide, is  
 37 listed in Table 6, too.

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

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

#### 51 Comparison of methods using multi-organ 52 measurements with multiple drug treatments

53 One might wonder if the proposed methods, TD- and  
 54 PCA-based unsupervised FE with optimized SD, are

1 applicable to a more complicated set-up. To investigate  
 2 this point, we checked the case where multiple drugs  
 3 are applied to mice whose gene expression of multiple  
 4 tissues are measured, to which we applied TD-based  
 5 unsupervised FE [34].

#### 6 *Enrichment of tissue-specific genes*

7 In the previous study [34], although we applied TD-  
 8 based unsupervised FE to gene expression profiles,  
 9 there existed some problems. First of all, the number of  
 10 genes selected was too small to have no false negatives.  
 11 Using the optimized SD, the number of selected genes  
 12 increased (Table 8; for more details, e.g., the defini-  
 13 tion of the four gene sets, neurons and testis, muscle,  
 14 gastrointestinal 1 and 2, see the previous study [34].  
 15 This topic has not been discussed herein as it is not  
 16 directly related to the comparison of the performance,  
 17 between the original TD-based unsupervised FE and  
 18 that with the optimised SD. The full list of the se-  
 19 lected genes is available in Additional file 9). Although  
 20 an increased number of genes is meaningless if the bi-  
 21 ological reliability is less, the biological reliability of  
 22 selected genes is also improved (lower panel of Fig. 15,  
 23 which corresponds to a present study and is associated  
 24 with a greater number of cell lines and tissue specificity  
 25 than that in the upper panel of Fig. 15, which corre-  
 26 sponds to a previous study). Thus, the employment  
 27 of optimized SD is also effective to a more complicated  
 28 data set than simple pairwise comparisons between the  
 29 treated and control samples investigated in the previ-  
 30 ous sections.

#### 32 *Coincidence with drug treatment*

33 We have also performed additional validation of the  
 34 genes selected by TD-based unsupervised FE with  
 35 optimized SD associated with adjusted  $P$ -values less  
 36 than 0.1 (Table 8, full list is available in Additional  
 37 files 10–13). We have uploaded selected genes to En-  
 38 richr [36] and evaluated the overlaps between the genes  
 39 selected and those whose expression was altered with  
 40 the treatment of the 15 drugs used in this study.  
 41 Then, we found that all four gene sets in Table 8  
 42 had a significant overlap with the genes whose expres-  
 43 sion was altered with the treatment of 5 of the drugs  
 44 (acetaminophen, cisplatin, clozapine, doxycycline, and  
 45 olanzapine) in DrugMatrix, which does not include  
 46 other drug treatments (Supplementary material). This  
 47 suggests that TD-based unsupervised FE with optimal  
 48 SD can correctly recognize drug treatments based on  
 49 gene expression; this was impossible in the previous  
 50 study [34] because of the very small number of genes  
 51 selected (Table 8). Thus, considering the optimization  
 52 of SD enables TD-based unsupervised FE to recognize  
 53 a greater number of biologically reliable genes than the  
 54 original TD-based unsupervised FE, which did not in-  
 55 clude the optimization of SD.

1 1  
2 2

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

Rank	Term	Overlap	P-value	Adjusted P-value	Odds Ratio
4 1	imatinib (glivec) 123596 human GSE12211 sample 2518	316/442	7.81E-137	7.06E-134	12.3
5 2	apratoxin A 6326668 human GSE2742 sample 3071	279/389	3.77E-121	1.57E-118	12.3
6 3	doxycycline DB00254 human GSE2624 sample 3074	294/425	5.22E-121	1.57E-118	10.9
7 4	doxycycline DB00254 human GSE2624 sample 3077	278/391	3.83E-119	8.64E-117	11.9
8 5	grepafloxacin 72474 human GSE9166 sample 2627	320/495	5.62E-119	1.02E-116	8.96
9 6	clinafloxacin 60063 human GSE9166 sample 2625	309/470	8.04E-118	1.21E-115	9.38
10 7	trovafloxacin 62959 human GSE9166 sample 2629	302/459	3.05E-115	3.94E-113	9.38
11 8	doxorubicin, 2xEC50, 5 d 31703 human GSE6930 sample 3265	314/493	4.76E-114	5.37E-112	8.57
12 9	cisplatin DB00515 human GSE6410 sample 2532	239/315	1.06E-112	1.07E-110	15.1
13 10	carboplatin DB00958 human GSE7035 sample 3060	284/422	4.57E-112	4.13E-110	9.99

14 14  
15 15  
16 16

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

Rank	Term	Overlap	P-value	Adjusted P-value	Odds Ratio
18 1	estradiol 5757 human GSE4668 sample 3063	276/367	1.26E-128	1.14E-125	14.74
19 2	tamoxifen DB00675 human GSE4025 sample 2820	271/361	6.30E-126	2.85E-123	14.61
20 3	apratoxin A 6326668 human GSE2742 sample 3068	278/389	4.61E-120	1.12E-117	12.16
21 4	estradiol DB00783 human GSE4668 sample 2727	261/350	4.96E-120	1.12E-117	14.19
22 5	MK-886 CID 3651377 human GSE3202 sample 3193	268/368	5.29E-119	9.59E-117	12.98
23 6	IFN-alphacon1 DB05258 human GSE5542 sample 2474	242/313	2.21E-117	3.34E-115	16.41
24 7	Arachidonic acid DB04557 human GSE3737 sample 3171	277/395	2.80E-116	3.63E-114	11.39
25 8	ARSENIC 5359596 human GSE6907 sample 3529	276/394	1.15E-115	1.30E-113	11.35
26 9	metoprolol DB00264 human GSE3356 sample 2786	306/469	2.67E-115	2.68E-113	9.16
27 10	estradiol 5757 human GSE4668 sample 3062	245/325	1.92E-114	1.74E-112	14.75

28 28  
29 29  
30 30

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

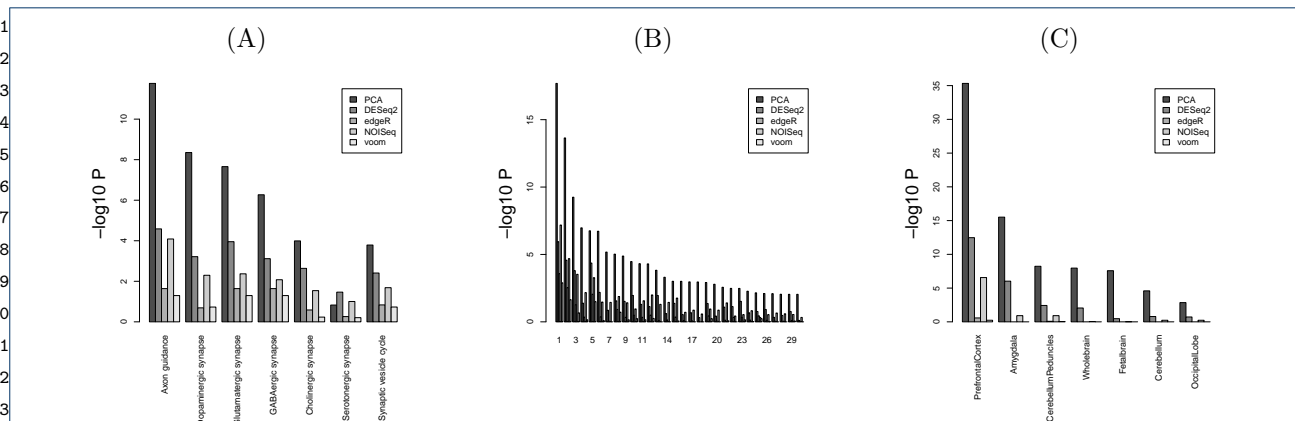
Rank	Term	Overlap	P-value	Adjusted P-value	Odds Ratio
32 1	PLX4032 DB05238 human GSE24862 sample 2568	65/318	1.59E-29	1.42E-26	7.06
33 2	dexamethasone DB01234 human GSE34313 sample 2714	51/297	7.68E-20	3.44E-17	5.59
34 3	dexamethasone DB01234 human GSE54608 sample 3093	52/322	5.45E-19	1.63E-16	5.19
35 4	VX 39793 human GSE33606 sample 3376	54/367	8.17E-18	1.58E-15	4.65
36 5	PLX4032 DB05238 human GSE24862 sample 2570	56/393	8.78E-18	1.58E-15	4.49
37 6	formoterol DB00983 human GSE30242 sample 2631	49/315	2.83E-17	4.23E-15	4.94
38 7	metformin DB00331 human GSE33612 sample 2483	50/343	2.07E-16	2.65E-14	4.58
39 8	etanercept DB00005 human GSE41663 sample 2605	45/322	3.29E-14	3.69E-12	4.33
40 9	cisplatin DB00515 human GSE47856 sample 3145	40/267	8.93E-14	8.91E-12	4.68
41 10	Lipopolysaccharide 11970143 human GSE5504 sample 3486	35/224	9.25E-13	8.30E-11	4.89

42 42  
43 43  
44 44

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

Rank	Term	Overlap	P-value	Adjusted P-value	Odds Ratio
46 1	resveratrol DB02709 human GSE25412 sample 3500	70/250	2.90E-41	2.63E-38	10.81
47 2	carboplatin (30 h) 10339178 human GSE13525 sample 3031	85/423	7.47E-38	3.38E-35	7.09
48 3	carboplatin (36 h) 10339178 human GSE13525 sample 3032	74/392	3.93E-31	1.19E-28	6.46
49 4	resveratrol DB02709 human GSE25412 sample 3501	51/194	7.59E-29	1.72E-26	9.66
50 5	Carboplatin DB00958 human GSE13525 sample 3089	65/357	1.69E-26	3.07E-24	6.11
51 6	NSC319726 5351307 human GSE35972 sample 2479	59/309	2.99E-25	4.52E-23	6.43
52 7	Lipopolysaccharide 11970143 human GSE5504 sample 3483	72/468	1.29E-24	1.67E-22	5.01
53 8	dasatinib DB01254 human GSE59357 sample 3306	57/298	1.81E-24	1.98E-22	6.43
54 9	thapsigargin 446378 human GSE19519 sample 3236	66/399	1.97E-24	1.98E-22	5.43
55 10	Y15 23627197 human GSE43452 sample 2554	64/390	1.59E-23	1.44E-21	5.37

56 56  
57 57



**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 axonextension,” 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.

adjusted <i>P</i> -values	TD-based unsupervised FE [34]		TD-based unsupervised FE with optimized SD	
	$\leq 0.01$	$\leq 0.01$	$\leq 0.1$	$\leq 0.1$
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 [34] and optimal SD with multi-organ data sets

## 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

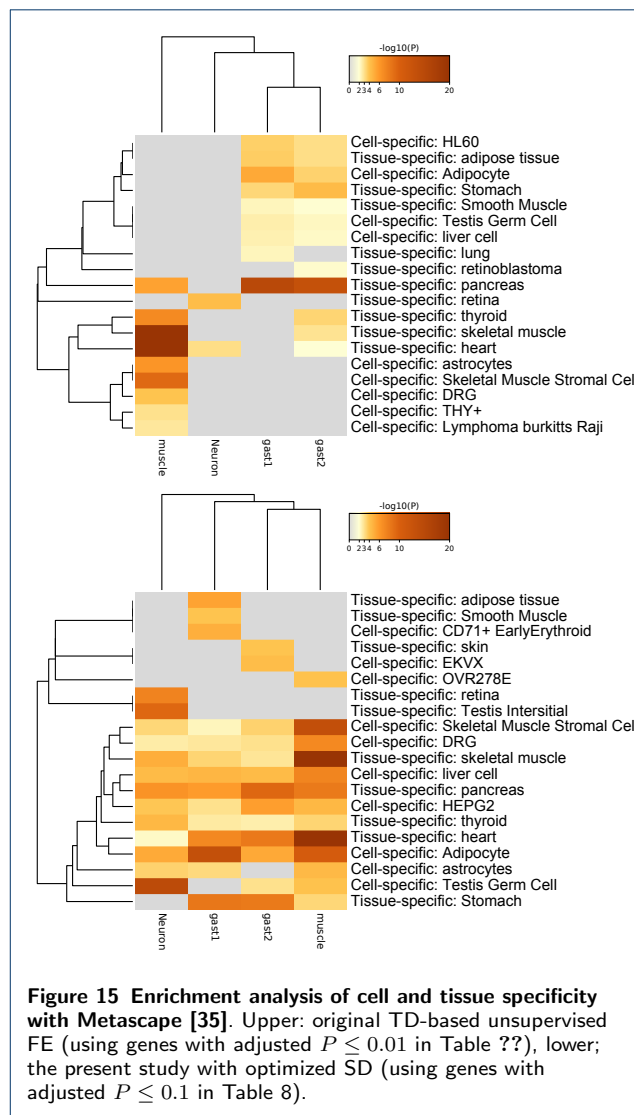
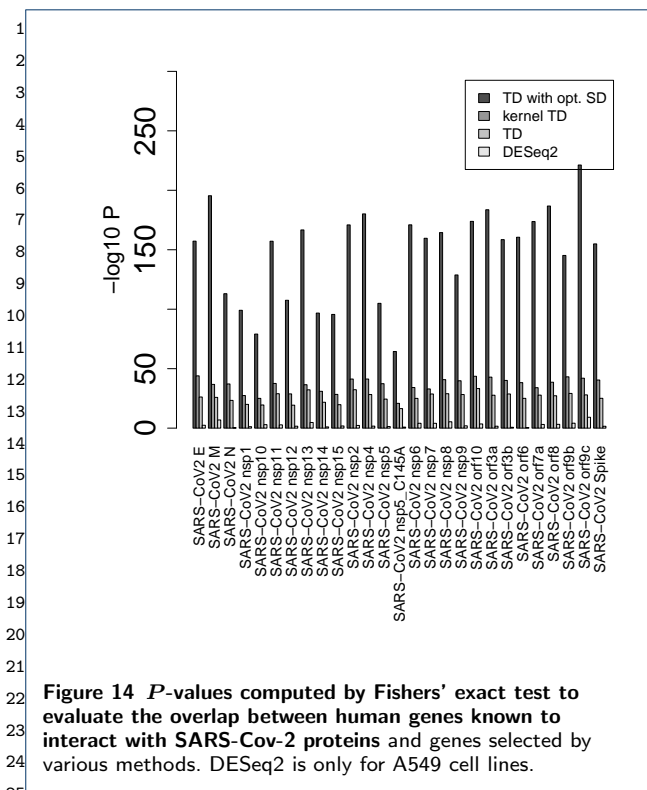
$$\Delta \equiv x - y \quad (24)$$

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

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.

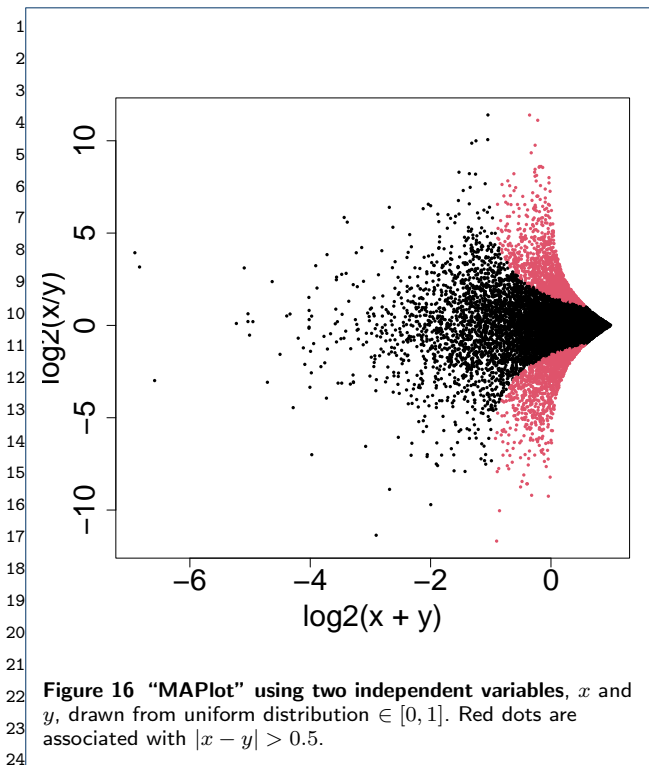
Another remarkable point of TD- and PCA-based unsupervised FE with optimized SD is that it does not have to screen for selected genes by LFC after the genes are selected using  $P$ -values. As can be seen in Fig. 10, state-of-art methods, including DESeq2, often identify too many DEGs. In these circumstances, LFC is often used to reduce the number of DEGs. Nevertheless, Stupnikov et al [37] 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.

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.

## Conclusions

In this study, we optimized SD to improve TD- and PCA-based unsupervised FE. As a result, not only the



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

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

### Gene expression profiles

#### MAQC

Seven human brain expression profiles were downloaded from SRA [38] (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 [39]. htseq-count [40] 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.primary\\_assemblies.gtf.gz](ftp://ftp.ensembl.org/pub/release-105/gtf/homo_sapiens/Homo_sapiens.GRCh38.primary_assemblies.gtf.gz)

#### SEQC

SEQC [13] were obtained from bioconductor [41] 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 [14]. 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 [34]. After getting  $u_{\ell i}$ , the SD was optimized as described in the main text.

#### 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 additional file 14 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 checked by using MAQC (Fig. 9) as follows.

<sup>1</sup>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.

### <sup>13</sup>Enrichment analyses

<sup>14</sup>Enrichment analyses were performed using either Metascape [35] or Enrichr [36] 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) [42, 43] was used for conversion.

### <sup>21</sup>DEG identification of SARS-CoV-2 data by DESeq2

<sup>22</sup>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.

### <sup>34</sup>Declarations

<sup>35</sup>Ethics approval and consent to participate  
<sup>36</sup>Not applicable.

<sup>37</sup>  
<sup>38</sup>Consent for publication  
<sup>39</sup>Not applicable.

### <sup>40</sup>Availability of data and materials

<sup>41</sup>The MAQC data set can be downloaded from SRA with ID SRX016359 and SRX016367. The SEQC data is a part of the bioconductor seqc package. SARS-CoV-2 data can be downloaded from Gene Expression Omnibus (GEO) with the GEO ID GSE147507. Multi-organ data can be downloaded from GEO with the GEO ID GSE142068.

<sup>46</sup>Competing interests  
<sup>47</sup>The authors declare that they have no competing interests.

### <sup>48</sup>Funding

<sup>49</sup>This work was supported by the Japan Society for the Promotion of Science, KAKENHI [Grant numbers 19H05270, 20K12067, and 20H04848] to YHT.

### <sup>52</sup>Author's contributions

<sup>53</sup>YHT planned the research and performed analyses. YHT and TT evaluated the results, discussions, and outcomes and drafted and reviewed the manuscript.

### Acknowledgements

Not applicable.

### Author details

<sup>1</sup>Department of Physics, Chuo University, 1-13-27 Kasuga, Bunkyo-ku, 112-8551 Tokyo, JAPAN. <sup>2</sup>Department of Computer Science, King Abdulaziz University, 21589 Jeddah, Saudi Arabia.

### References

1. Taguchi, Y.-h.: Comparative transcriptomics analysis. In: Ranganathan, S., Gribkov, M., Nakai, K., Schönbach, C. (eds.) Encyclopedia of Bioinformatics and Computational Biology, pp. 814–818. Academic Press, Oxford (2019). doi:10.1016/B978-0-12-809633-8.20163-5. <https://www.sciencedirect.com/science/article/pii/B9780128096338201635>
2. Rapaport, F., Khanin, R., Liang, Y., Pirun, M., Krek, A., Zumbo, P., Mason, C.E., Succi, N.D., Betel, D.: Erratum to: Comprehensive evaluation of differential gene expression analysis methods for RNA-seq data. *Genome Biology* **16**(1) (2015). doi:10.1186/s13059-015-0813-z
3. Tusher, V.G., Tibshirani, R., Chu, G.: Significance analysis of microarrays applied to the ionizing radiation response. *Proceedings of the National Academy of Sciences* **98**(9), 5116–5121 (2001). doi:10.1073/pnas.091062498. <https://www.pnas.org/content/98/9/5116.full.pdf>
4. Ritchie, M.E., Phipson, B., Wu, D., Hu, Y., Law, C.W., Shi, W., Smyth, G.K.: limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Research* **43**(7), 47–47 (2015). doi:10.1093/nar/gkv007. <https://academic.oup.com/nar/article-pdf/43/7/e47/7207289/gkv007.pdf>
5. Love, M.I., Huber, W., Anders, S.: Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biology* **15**(12) (2014). doi:10.1186/s13059-014-0550-8
6. Robinson, M.D., McCarthy, D.J., Smyth, G.K.: edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* **26**(1), 139–140 (2009). doi:10.1093/bioinformatics/btp616. <https://academic.oup.com/bioinformatics/article-pdf/26/1/139/443156/btp616.pdf>
7. McCarthy, D.J., Chen, Y., Smyth, G.K.: Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research* **40**(10), 4288–4297 (2012). doi:10.1093/nar/gks042. <https://academic.oup.com/nar/article-pdf/40/10/4288/25335174/gks042.pdf>
8. Law, C.W., Chen, Y., Shi, W., Smyth, G.K.: voom: precision weights unlock linear model analysis tools for RNA-seq read counts. *Genome Biology* **15**(2), 29 (2014). doi:10.1186/gb-2014-15-2-r29
9. 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**(B), 18–19 (2012). doi:10.14806/ej.17.B.265
10. Taguchi, Y.-h.: Unsupervised Feature Extraction Applied to Bioinformatics. Springer, Singapore (2020). doi:10.1007/978-3-030-22456-1. <https://doi.org/10.1007/978-3-030-22456-1>
11. Shi, L., Shi, L., Reid, L.H., Jones, W.D., Shippy, R., Warrington, J.A., Baker, S.C., Collins, P.J., de Longueville, F., Kawasaki, E.S., Lee, K.Y., Luo, Y., Sun, Y.A., Willey, J.C., Setterquist, R.A., Fischer, G.M., Tong, W., Dragan, Y.P., Dix, D.J., Frueh, F.W., Goodsaid, F.M., Herman, D., Jensen, R.V., Johnson, C.D., Lobenhofer, E.K., Puri, R.K., Scherf, U., Thierry-Mieg, J., Wang, C., Wilson, M., Wolber, P.K., Zhang, L., Amur, S., Bao, W., Barbacioru, C.C., Lucas, A.B., Bertholet, V., Boysen, C., Bromley, B., Brown, D., Brunner, A., Canales, R., Cao, X.M., Cebula, T.A., Chen, J.J., Cheng, J., Chu, T.-M., Chudin, E., Corson, J., Corton, J.C., Croner, L.J., Davies, C., Davison, T.S., Delenstarr, G., Deng, X., Dorris, D., Eklund, A.C., Fan, X.-h., Fang, H., Fulmer-Smentek, S., Fuscoe, J.C., Gallagher, K., Ge, W., Guo, L., Guo, X., Hager, J., Haje, P.K., Han, J., Han, T., Harbottle, H.C., Harris, S.C., Hatchwell, E., Hauser, C.A., Hester, S., Hong, H., Hurban, P., Jackson, S.A., Ji, H., Knight, C.R., Kuo, W.P., LeClerc, J.E., Levy, S., Li, Q.-Z., Liu, C., Liu, Y., Lombardi, M.J., Ma, Y., Magnuson, S.R., Maqsoodi, B., McDaniel, T., Mei, N., Myklebost, 55

	Cell lines	adjusted $P$ -values $\leq 0.01$	alternative conditions	the number of DEG2	
1	Calu3	16432	adjusted $P$ -value $\leq 0.05$ , LFC $> 2.0$	340	1
2	NHBE	327	adjusted $P$ -value $\leq 0.05$ , LFC $> 0.5$	171	2
3	A549				3
4	MOI 0.2	15852	adjusted $P$ -value $\leq 0.05$ , LFC $> 2.0$	176	4
5	MOI 2.0	7431	adjusted $P$ -value $\leq 0.05$ , LFC $> 2.0$	547	5
5	ACE2 expressed	7509	adjusted $P$ -value $\leq 0.05$ , LFC $> 1.0$	756	5

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

O., Ning, B., Novorodovskaya, N., Orr, M.S., Osborn, T.W., Papallo, A., Patterson, T.A., Perkins, R.G., Peters, E.H., Peterson, R., Phillips, K.L., Pine, P.S., Pusztai, L., Qian, F., Ren, H., Rosen, M., Rosenzweig, B.A., Samaha, R.R., Schena, M., Schroth, G.P., Shchegrova, S., Smith, D.D., Staedtler, F., Su, Z., Sun, H., Szallasi, Z., Tezak, Z., Thierry-Mieg, D., Thompson, K.L., Tikhonova, I., Turpaz, Y., Vallanat, B., Van, C., Walker, S.J., Wang, S.J., Wang, Y., Wolfinger, R., Wong, A., Wu, J., Xiao, C., Xie, Q., Xu, J., Yang, W., Zhang, L., Zhong, S., Zong, Y., Slikker, W.: The MicroArray quality control (MAQC) project shows inter- and intraplatform reproducibility of gene expression measurements. *Nature Biotechnology* **24**(9), 1151–1161 (2006). doi:10.1038/nbt1239

12. Mudge, J.F., Baker, L.F., Edge, C.B., Houlihan, J.E.: Setting an optimal  $\alpha$  that minimizes errors in null hypothesis significance tests. *PLOS ONE* **7**(2), 1–7 (2012). doi:10.1371/journal.pone.0032734

13. A comprehensive assessment of RNA-seq accuracy, reproducibility and information content by the sequencing quality control consortium. *Nature Biotechnology* **32**(9), 903–914 (2014). doi:10.1038/nbt.2957

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**(9), 1–16 (2020). doi:10.1371/journal.pone.0238907

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 Journal of Selected Topics in Signal Processing* **15**(3), 746–758 (2021). doi:10.1109/JSTSP.2021.3061251

16. Zhao, H., Mendenhall, M., Deininger, M.W.: Imatinib is not a potent anti-SARS-CoV-2 drug. *Leukemia* **34**(11), 3085–3087 (2020). doi:10.1038/s41375-020-01045-9

17. Naidoo, D., Roy, A., Kar, P., Mutanda, T., Anandraj, A.: Cyanobacterial metabolites as promising drug leads against the mpro and plpro of sars-cov-2: an in silico analysis. *Journal of Biomolecular Structure and Dynamics* **39**(16), 6218–6230 (2021). doi:10.1080/07391102.2020.1794972. PMID: 32691680. <https://doi.org/10.1080/07391102.2020.1794972>

18. Dorobisz, K., Dorobisz, T., Janczaki, D., Zatoński, T.: Doxycycline in the coronavirus disease 2019 therapy. *Therapeutics and Clinical Risk Management* **Volume 17**, 1023–1026 (2021). doi:10.2147/tcrm.s314923

19. Gimeno, A., Mestres-Truyol, J., Ojeda-Montes, M.J., Macip, G., Saldivar-Espinoza, B., Cereto-Massagué, A., Pujadas, G., Garcia-Vallé, S.: Prediction of novel inhibitors of the main protease (M-pro) of SARS-CoV-2 through consensus docking and drug reposition. *International Journal of Molecular Sciences* **21**(11) (2020). doi:10.3390/ijms21113793

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. *Journal of Biomolecular Structure and Dynamics* **0**(0), 1–15 (2021). doi:10.1080/07391102.2021.1905551. PMID: 33826483. <https://doi.org/10.1080/07391102.2021.1905551>

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. *Scientific Reports* **11**(1) (2021). doi:10.1038/s41598-021-01410-3

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. *Scientific Reports* **12**(1) (2022). doi:10.1038/s41598-022-04778-y

23. Zu, S., Luo, D., Li, L., Ye, Q., Li, R.-T., Wang, Y., Gao, M., Yang, H., Deng, Y.-Q., Cheng, G.: Tamoxifen and clomiphene inhibit SARS-CoV-2 infection by suppressing viral entry. *Signal Transduction and Targeted Therapy* **6**(1) (2021). doi:10.1038/s41392-021-00853-4

24. Zhu, W., Xu, M., Chen, C.Z., Guo, H., Shen, M., Hu, X., Shinn, P., Klumpp-Thomas, C., Michael, S.G., Zheng, W.: Identification of SARS-CoV-2 3cl protease inhibitors by a quantitative high-throughput screening. *ACS Pharmacology & Translational Science* **3**(5), 1008–1016 (2020). doi:10.1021/acspsci.0c00108

25. Paragas, J., Blatt, L.M., Hartmann, C., Mandala, S.M., Endy, T.P.: Interferon alfacon1 is an inhibitor of SARS-corona virus in cell-based models. *Antiviral Research* **66**(2), 99–102 (2005). doi:10.1016/j.antiviral.2005.01.002

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 (2021). doi:10.1016/j.prostaglandins.2021.106539

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 (2020). doi:10.26434/chemrxiv.12200495.v1

28. Clemente-Moragón, A., Martínez-Milla, J., Oliver, E., Santos, A., Flandes, J., Fernández, I., Rodríguez-González, L., del Castillo, C.S., Ioan, A.-M., López-Álvarez, M., Gómez-Talavera, S., Galán-Arriola, C., Fuster, V., Pérez-Calvo, C., Ibáñez, B.: Metoprolol in critically ill patients with COVID-19. *Journal of the American College of Cardiology* **78**(10), 1001–1011 (2021). doi:10.1016/j.jacc.2021.07.003

29. Dexamethasone in hospitalized patients with covid-19. *New England Journal of Medicine* **384**(8), 693–704 (2021). doi:10.1056/nejmoa2021436

30. Parthasarathy, H., Tandel, D., Harshan, K.H.: Metformin suppresses SARS-CoV-2 in cell culture. *bioRxiv* (2021). doi:10.1101/2021.11.18.469078. <https://www.biorxiv.org/content/early/2021/11/22/2021.11.18.469078.full.pdf>

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. *Rheumatology and Therapy* **8**(3), 1355–1370 (2021). doi:10.1007/s40744-021-00342-8

32. Petruk, G., Puthia, M., Petrolva, J., Samsudin, F., Strömdahl, A.-C., Cerps, S., Uller, L., Kjellström, S., Bond, P.J., Schmidtchen, A.: SARS-CoV-2 spike protein binds to bacterial lipopolysaccharide and boosts proinflammatory activity. *Journal of Molecular Cell Biology* **12**(12), 916–932 (2020). doi:10.1093/jmcb/mjaa067. <https://academic.oup.com/jmcb/article-pdf/12/12/916/36546065/mjaa067.pdf>

33. Pasquereau, S., Nehme, Z., Haidar Ahmad, S., Daouad, F., Van Assche, J., Wallet, C., Schwartz, C., Rohr, O., Morot-Bizot, S., Herbein, G.: Resveratrol inhibits HCoV-229E and SARS-CoV-2 coronavirus replication in vitro. *Viruses* **13**(2) (2021). doi:10.3390/v13020354

34. Taguchi, Y., Turki, T.: Universal nature of drug treatment responses in drug-tissue-wide model-animal experiments using tensor decomposition-based unsupervised feature extraction. *Frontiers in Genetics* **11** (2020). doi:10.3389/fgene.2020.00695

35. Zhou, Y., Zhou, B., Pache, L., Chang, M., Khodabakhshi, A.H., Tanaseichuk, O., Benner, C., Chanda, S.K.: Metascape provides a biologist-oriented resource for the analysis of systems-level datasets. *Nature Communications* **10**(1) (2019). doi:10.1038/s41467-019-09234-6



1	36. Xie, Z., Bailey, A., Kuleshov, M.V., Clarke, D.J.B., Evangelista, J.E.,	Additional file 7 — Overlap with human genes known to interact with	1
2	Jenkins, S.L., Lachmann, A., Wojciechowicz, M.L., Kropiwnicki, E.,	SARS-CoV-2 protein by DESeq2	2
3	Jagodnik, K.M., Jeon, M., Ma'ayan, A.: Gene set knowledge discovery	Overlap with human genes known to interact with SARS-CoV-2 proteins by	3
4	with Enrichr. <i>Current Protocols</i> <b>1</b> (3), 90 (2021). doi:10.1002/cpz1.90.	DESeq2.	4
5	<a href="https://currentprotocols.onlinelibrary.wiley.com/doi/pdf/10.1002/cpz1.90">https://currentprotocols.onlinelibrary.wiley.com/doi/pdf/10.1002/cpz1.90</a>		5
6	37. Stupnikov, A., McInerney, C.E., Savage, K.I., McIntosh, S.A.,	Additional file 8 — Genes selected for SARS-CoV-2 infected A549 cell lines	6
7	Emmert-Streib, F., Kennedy, R., Salto-Tellez, M., Prise, K.M., McArt,	by DESeq2	6
8	D.G.: Robustness of differential gene expression analysis of RNA-seq.	Genes selected by DESeq2, for A549 cell lines, shown in Table 14, and drug	7
9	<i>Computational and Structural Biotechnology Journal</i> <b>19</b> , 3470–3481	repositioning.	8
10	(2021). doi:10.1016/j.csbj.2021.05.040		9
11	38. Leinonen, R., Sugawara, H., Shumway, o.b.o.t.I.N.S.D.C. Martin: The	Additional file 9 — Genes selected by TD-based unsupervised FE with	10
12	Sequence Read Archive. <i>Nucleic Acids Research</i> <b>39</b> (suppl.1), 19–21	optimized SD for multi-organ study	10
13	(2010). doi:10.1093/nar/gkq1019.	Genes selected by TD-based unsupervised FE with optimized SD for a	11
14	<a href="https://academic.oup.com/nar/article-pdf/39/suppl.1/D19/7624335/gkq1019.pdf">https://academic.oup.com/nar/article-</a>	multi-organ study.	12
15	<a href="https://academic.oup.com/nar/article-pdf/39/suppl.1/D19/7624335/gkq1019.pdf">pdf/39/suppl.1/D19/7624335/gkq1019.pdf</a>		13
16	39. Srivastava, A., Sarkar, H., Gupta, N., Patro, R.: RapMap: a rapid,	Additional file 10 — Drug repositioning for neuron and tesis gene sets	14
17	sensitive and accurate tool for mapping RNA-seq reads to	Drug repositioning for neuron and tesis gene sets.	14
18	transcriptomes. <i>Bioinformatics</i> <b>32</b> (12), 192–200 (2016).		15
19	doi:10.1093/bioinformatics/btw277.	Additional file 11 — Drug repositioning for muscle gene sets	16
20	<a href="https://academic.oup.com/bioinformatics/article-pdf/32/12/i192/17130476/btw277.pdf">https://academic.oup.com/bioinformatics/article-</a>	Drug repositioning for muscle gene sets.	17
21	<a href="https://academic.oup.com/bioinformatics/article-pdf/32/12/i192/17130476/btw277.pdf">pdf/32/12/i192/17130476/btw277.pdf</a>		18
22	40. Putri, G.H., Anders, S., Pyl, P.T., Pimanda, J.E., Zanini, F.: Analysing	Additional file 12 — Drug repositioning for gast 1 gene sets	19
23	high-throughput sequencing data in Python with HTSeq 2.0 (2021).	Drug repositioning for gast 1 gene sets.	19
24	2112.00939. <a href="https://arxiv.org/abs/2112.00939">https://arxiv.org/abs/2112.00939</a>		20
25	41. Huber, W., Carey, V.J., Gentleman, R., Anders, S., Carlson, M.,	Additional file 13 — Drug repositioning for gast 2 gene sets	21
26	Carvalho, B.S., Bravo, H.C., Davis, S., Gatto, L., Girke, T., Gottardo,	Drug repositioning for gast 2 gene sets.	22
27	R., Hahne, F., Hansen, K.D., Irizarry, R.A., Lawrence, M., Love, M.I.,		23
28	MacDonald, J., Obenchain, V., Oleś, A.K., Pagès, H., Reyes, A.,	Additional file 14 — Source code	24
29	Shannon, P., Smyth, G.K., Tenenbaum, D., Waldron, L., Morgan, M.:	R source code to perform PCA- and TD-based unsupervised FE with	24
30	Orchestrating high-throughput genomic analysis with bioconductor.	optimized SD.	25
31	<i>Nature Methods</i> <b>12</b> (2), 115–121 (2015). doi:10.1038/nmeth.3252		26
32	42. Huang, D.W., Sherman, B.T., Lempicki, R.A.: Systematic and		27
33	integrative analysis of large gene lists using DAVID bioinformatics		28
34	resources. <i>Nature Protocols</i> <b>4</b> (1), 44–57 (2008).		29
35	doi:10.1038/nprot.2008.211		30
36	43. Huang, D.W., Sherman, B.T., Lempicki, R.A.: Bioinformatics		31
37	enrichment tools: paths toward the comprehensive functional analysis		32
38	of large gene lists. <i>Nucleic Acids Research</i> <b>37</b> (1), 1–13 (2008).		33
39	doi:10.1093/nar/gkn923. <a href="https://academic.oup.com/nar/article-pdf/37/1/1/17059338/gkn923.pdf">https://academic.oup.com/nar/article-</a>		34
40	<a href="https://academic.oup.com/nar/article-pdf/37/1/1/17059338/gkn923.pdf">pdf/37/1/1/17059338/gkn923.pdf</a>		35
41			36
42	<b>Figures</b>		37
43	<b>Tables</b>		38
44	<b>Additional Files</b>		39
45	Additional file 1 — Genes selected by DESeq2 for MAQC		40
46	Genes associated with adjusted <i>P</i> -values less than 0.1 using DESeq2 and		41
47	enrichment analysis associated with them.		42
48			43
49	Additional file 2 — Genes selected by PCA-based unsupervised FE with		44
50	optimized SD for MAQC		45
51	Genes associated with adjusted <i>P</i> -values less than 0.1 using PCA-based		46
52	unsupervised FE with optimized SD and enrichment analysis associated		47
53	with them.		48
54			49
55	Additional file 3 — Genes selected by EdgeR for MAQC		50
	Genes associated with adjusted <i>P</i> -values less than 0.1 using EdgeR and		51
	enrichment analysis associated with them.		52
			53
	Additional file 4 — Genes selected by voom for MAQC		54
	Genes associated with adjusted <i>P</i> -values less than 0.1 using voom and		55
	enrichment analysis associated with them.		
	Additional file 5 — Genes selected by NOISeq for MAQC		
	Genes associated with adjusted <i>P</i> -values less than 0.1 using NOISeq and		
	enrichment analysis associated with them.		
	Additional file 6 — Genes selected by TD-based unsupervised FE with		
	optimized SD for SARS-CoV-2		
	Genes associated with adjusted <i>P</i> -values less than 0.1 by TD-based		
	unsupervised FE with optimized SD for SARS-CoV-2 and drug		
	repositioning associated with the genes.		