

**scAnalyzeR: a comprehensive software package
with graphical user interface for single-cell RNA
sequencing analysis**

User Manual

scAnalyzeR

Version 1.0

04 June 2021

Contents	Page
1. Introduction	1
2. How to Setup	1
3. Upload Dataset	3
4. Pre-processing	6
5. Normalization	12
6. Dimensionality Reduction	15
7. Clustering	18
8. Differential Expression Analysis	19
9. Plots	35
10. Pathway Analysis	41
11. Trajectory Analysis	43
12. References	45

1. Introduction

scAnalyzeR is a comprehensive platform for analysing and visualizing single-cell RNA sequencing (scRNA-seq) data with an interactive graphical user interface. The scRNA-seq technology is becoming popular for the investigation of heterogeneous cell populations at single cell level, such as identifying cell diversity [1] and revealing cell subtypes [2]. Besides, this technology is applied for discovering disease surface markers [3] and constructing the path in progression over time for individual cells [4, 5]. Here is a brief description of the scAnalyzeR tool how it can be run successfully. This user manual covers procedures for uploading, pre-processing (discarding low-quality samples, outlier detection and removal, etc.), gene-expression normalization, highly variable genes(features) identification, cell clustering, differential gene expression analysis (including cluster-oriented marker genes finding), five basic plotting functions (violin plot, feature plot, heatmap, dynamic pca, and correlation plot), pathway enrichment analysis, and pseudo-time construction (trajectories development).

2. How to setup

There are two different ways to setting up the pipeline in your own machine:

Way 1: using from Docker image (strongly recommended**)**

1. Download and install Docker (<https://www.docker.com/products/docker-desktop>)
2. Pull the docker image from Docker Hub by running the following command:

```
docker pull gscdocker/scanalyzer:latest
```

3. Run the docker image locally on your computer and access the link:

To run the docker image, execute the following command:

```
docker run -d --rm -p 3838:3838 gscdocker/scanalyzer:latest
```

After running the docker image successfully, open the following link on a web browser (e.g., Firefox, Google Chrome) to access the pipeline:

```
http://localhost:3838/
```

Way 2: using from source

Firstly, you need to download and install following softwares (install R then RStudio):

Download and install R and RStudio on your machine,

- i. Download and install R (v-3.6.2 or above): <https://cran.r-project.org/>

- ii. Download and install RStudio: RStudio (v-1.1.456 or above):

<https://rstudio.com/products/rstudio/download/>

After installing the R and RStudio on your machine successfully, then, you need to clone this (<https://github.com/sarwarchy20/scAnalyzeR/archive/master.zip>) repository as well as unzip it.

Now, please run the following script on RStudio to install the **renv** R package:

```
install.packages("renv")
```

Next, run the following scripts on RStudio to install all the dependent R packages. Please replace ~ with the location of your scAnalyzeR-master unzipped folder:

```
renv::consent(provided=TRUE)
setwd("~/scAnalyzeR-master")
renv::restore()
```

Finally, run the app using RStudio by running the script below:

```
shiny::runApp ('~/scAnalyzeR-master/')
```

After successfully running the scAnalyzeR, the GUI will be displayed automatically. You can customize (minimize or maximize) window size via mouse clicking. If you would like to open the scAnalyzeR in a browser, click the “Open in Browser” or copy the link (top left corner) and paste it in a browser taskbar (**Fig. 1**).

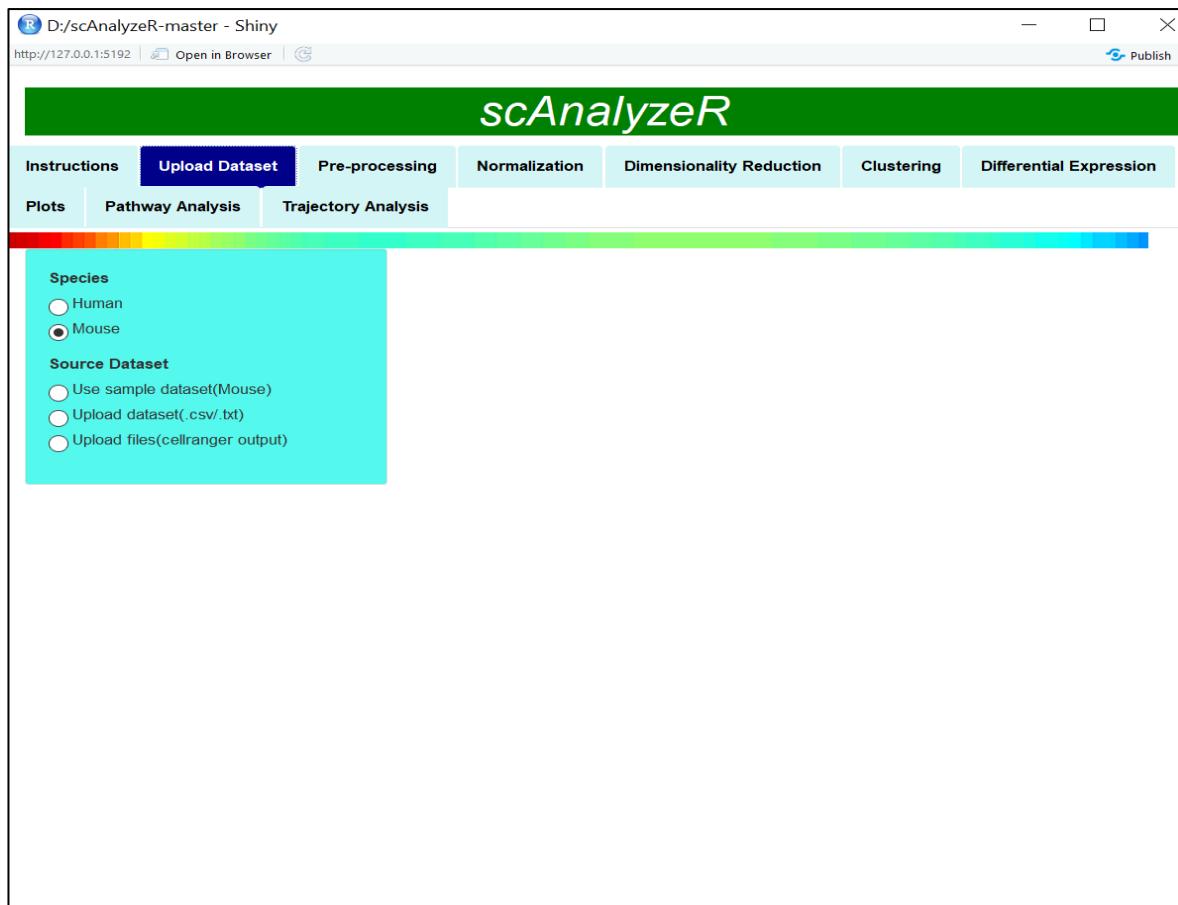


Figure 1. Interface of the scAnalyzeR

3. Upload dataset

This module is the first step for single cell RNA sequencing (scRNA-seq) data analysis using our developed tool, scAnalyzeR. Before uploading the data, first, choose the right species (Human or Mouse) option from where the data was sequenced. After selecting the appropriate species option, then select the perfect data source option which has the best match with your data file format. A sample dataset (mouse ovarian surface epithelium, 14,169 genes across 329 control cells) [6] is integrated with the tool, only for demonstration purpose. If you wish to use it, select the “Upload sample dataset (Mouse),” then the data file is loading automatically, and a progress bar will be shown on the bottom right corner of the interface. A text notification also will be displayed after loading the data file successfully (**Fig. 3**).

If you select the “Upload dataset(.csv/.txt)” option, it will show the data uploading facilities (interface) with some parameters, i.e., Browse, Header, Separator, and Quote. The comma (as a separator), double quote (as a quote), and header are selected by default. You can set different option(s) that will satisfy with your dataset format, which is going to be uploaded. For this data source option, the dataset must be a single file, and the first row provides the identities of the

cells (e.g., cell barcodes), and the first column contains the gene symbols. By clicking the “Browse” button, select the desired file (from your machine directory where the data file is located). After clicking the “open” menu, the file is uploading along with a progress bar, and it will show a text notification (how many genes and cells does exist in the uploaded data file) on the right side of the interface after completion the uploading.

For cellranger’s output data uploading (**Fig. 2**), first, choose the “Upload files(cellranger output),” then click the “Browse” button, and select the three files- matrix.mtx, genes.tsv and barcodes.tsv, respectively. In our example, uploaded a 3k data set (2,700 single cells with 32,738 genes), a cellranger output of Peripheral Blood Mononuclear Cells (PBMC) from 10X Genomics (<https://support.10xgenomics.com/single-cell-gene-expression/datasets/1.1.0/pbmc3k>).

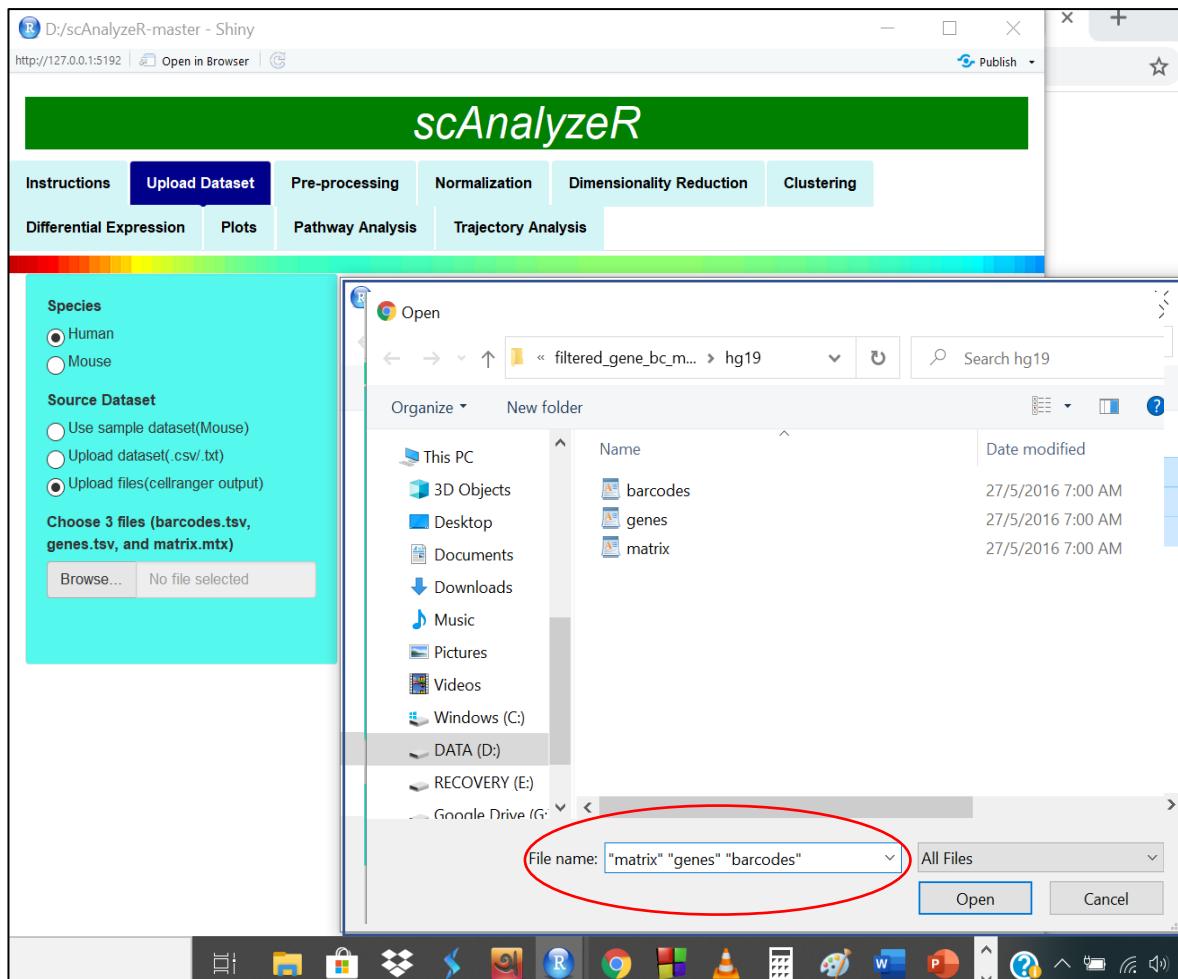


Figure 2. Upload cellranger's files. Follow the file name order (matrix, genes, barcodes) when you are using the Windows OS.



Figure 3. A text notification after successfully uploading the data set

4. Pre-processing

In this step, low-quality samples can be identified and discarded by setting some parameters. There are two-step filtering procedures – create dataset (**Fig. 4-6**), and filtering (**Fig. 7-8**). The dataset can be created by clicking the “Create dataset” button. The created dataset will keep all genes that are expressed in a minimum number of cells (3, by default) and contain those cells

with at least minimum genes (200, by default) detected. The default value(s) can be changed as required. Now, it is ready for showing the list of mitochondrial genes, plot metadata, gene summary plot, and download mitochondrial gene expression. To show a plot or list, select a particular checkbox via mouse clicking. Moreover, it can be copied and saved via right-click on the mouse. Deselect the checkbox to clear plot (s) or list.

Sometimes very few genes present in low-quality cells or empty droplets, and an abnormally high gene count may lie in cell doublets or multiplets. Even in, immense mitochondrial contamination points out in dying cells or low-quality cells. Now our target is discarding cells that have a paranormally high and low gene counts and a significant level of mitochondrial contamination. In the filtering step, check the “Filtering Cells” checkbox, then the filtering interface will be shown with parameters list along default values setting. The default values can be edited by mouse clicking or typing. After clicking the “Submit” button, a text notification will be shown, which indicates how many cells remain in the dataset after filtering. In this instance, discarded cells that have unique gene counts over 2500 or less than 200 and greater than 5% mitochondrial counts. If the notification shows a message “*Error: Cannot find cells provided,*” it means that the dataset is empty (no cell exists) after filtered with these parameters setting. In this situation, it is necessary to change the parameter(s) value and click the “Submit” button again.

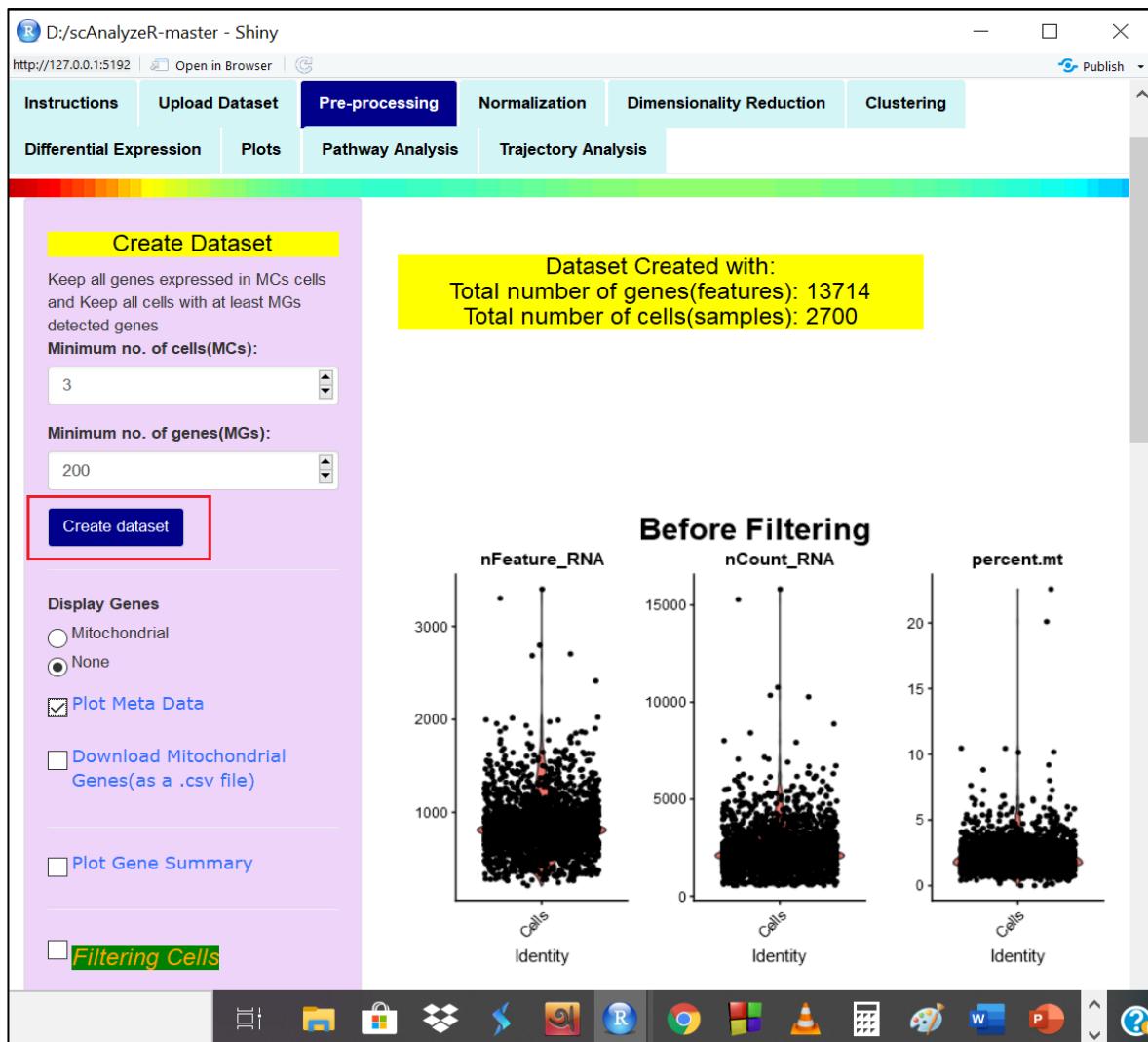


Figure 4. Create data: Meta data plot after creating the data set

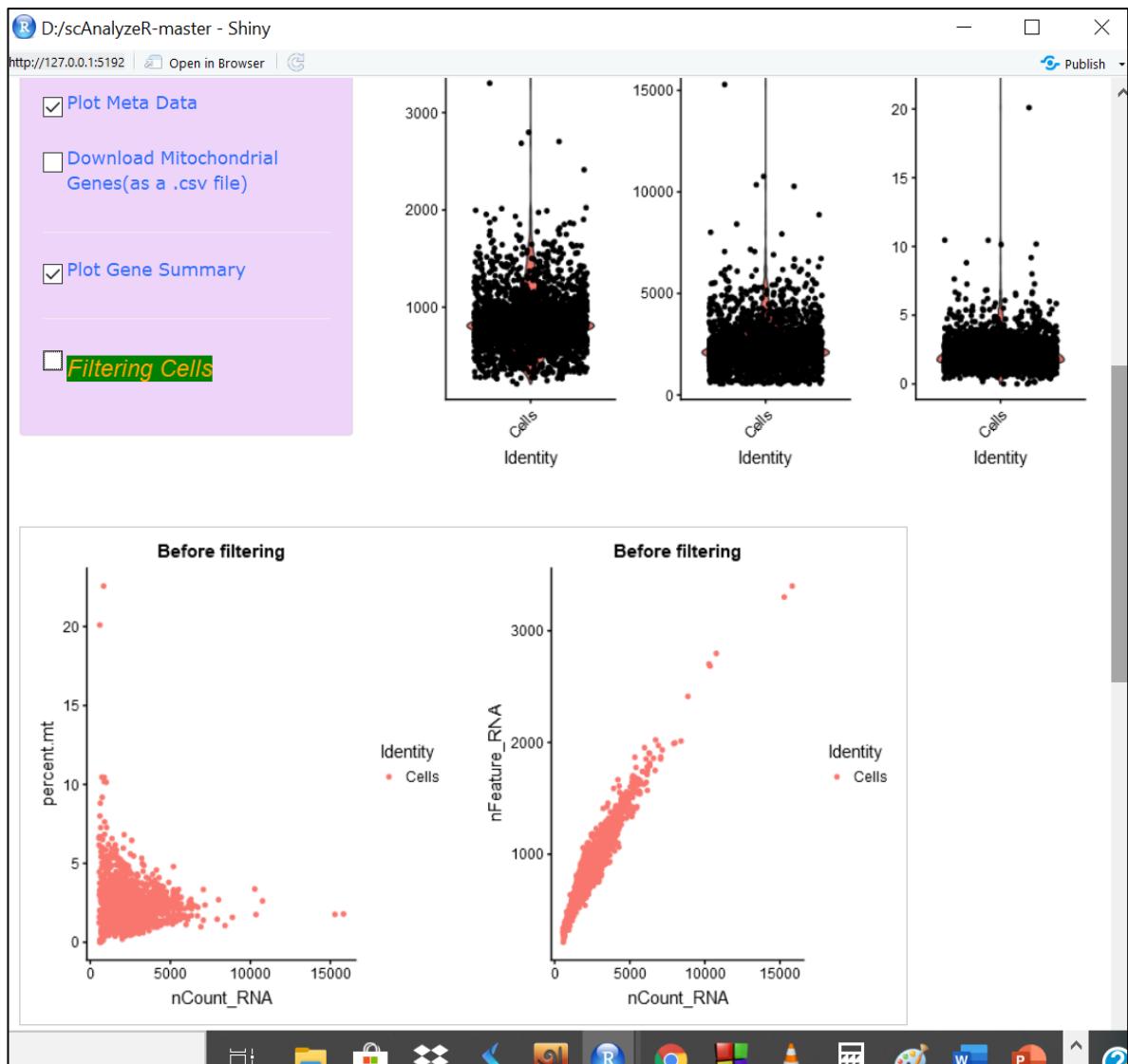


Figure 5. Create data: Gene summary plot (bottom) after creating the data set

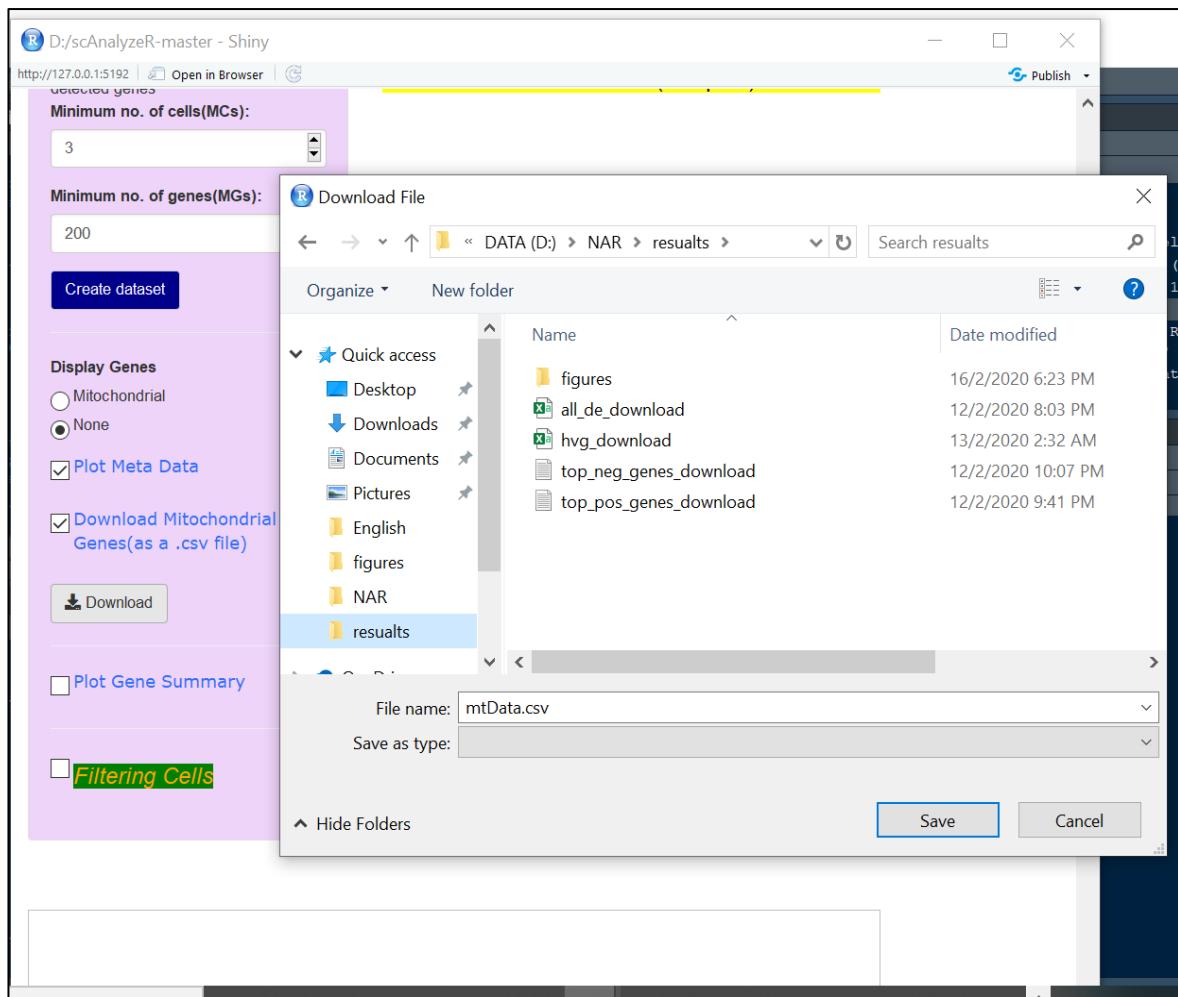


Figure 6. Create data: Download mitochondrial gene expressions after creating the data set

D:/scAnalyzeR-master - Shiny

http://127.0.0.1:5192 | Open in Browser |

None

Plot Meta Data

Download Mitochondrial Genes(as a .csv file)

Plot Gene Summary

Filtering Cells

Filter(discard) cells that have unique Gene counts over nFeature_RNA(>) or less than nFeature_RNA(<), and mitochondrial counts above percent.mt(>)

nFeature_RNA(<):
200

nFeature_RNA(>):
2500

percent.mt(>):
5

Plot Meta Data

Plot Gene Summary

Filtered Dataset:

Total number of genes(features): 13714
Total number of cells(samples): 2638



Figure 7. Filtering

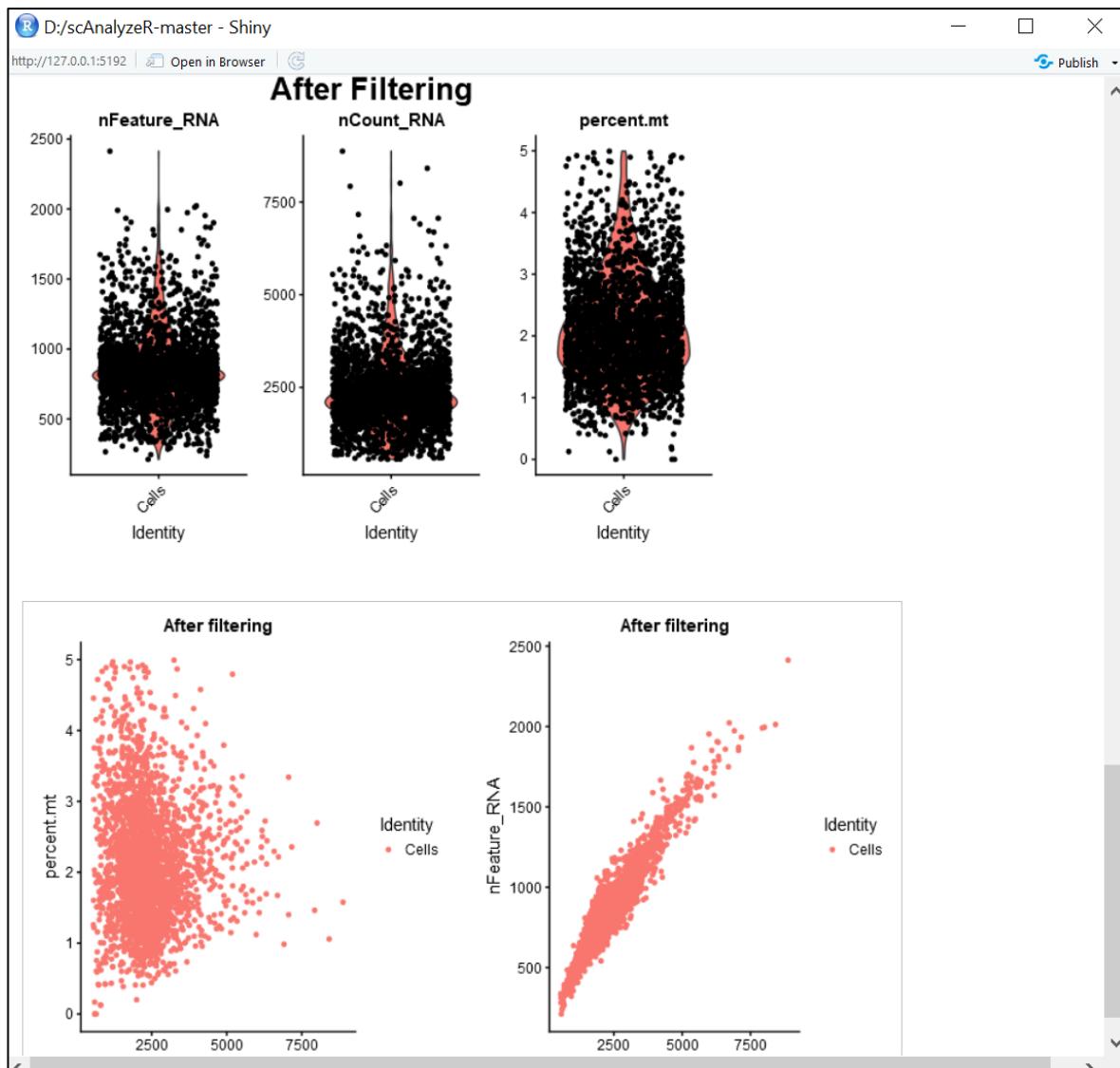


Figure 8. Filtering: Meta data plot(top), Gene summary plot(bottom) after filtering

5. Normalization

After discarding unwanted cells, in this module, the normalization method will be applied, and highly variable Genes will be identified from the normalized data. To normalize, choose a normalization method from the “Normalization Method” drop-down list (**Fig. 9**). The ‘Log Normalization’ method selected by default. The user can also change the scaling factor value. Now, press the “Submit” button, then a message will be shown after successfully normalized.

To find highly variable Genes (**Fig. 9**), set the number (2000 by default) for how many Genes will be considered as highly variable to all remained modules where it (*‘High Variable Genes Only’*) is indicated (e.g., Dimensionality Reduction, and Trajectory Analysis). It is necessary to press the “Find” button to detect highly variable Genes. To see the computed list of highly

variable Genes, check the “Show Highly Variable Genes” checkbox. The gene list is sorted descending order according to standardized variance, and 10 entities are showing by default. To increase the number of genes for showing, select a number from the drop-down menu, and write a specific gene symbol in the search box for searching individual highly variable gene properties. Top (10 by default) highly variable Genes’ plot will be shown through checking on the “Variable Genes Plot” checkbox (**Fig. 10**). To download the expression matrix of detected high variable Genes, check the “Download Highly Variable Genes(as a .csv file)” box, then press the download button.

D:/scAnalyzeR-master - Shiny
http://127.0.0.1:5192 | Open in Browser |

Normalization

Normalized Dataset:
Total number of genes/features): 13714
Total number of cells(samples): 2638

Set Normalization Parameters

Normalization method: Log Normalization

Scaling factor: 10000

Find Highly Variable Genes

Number of Genes: 2000

Show Highly Variable Genes (checkbox checked)

Variable Genes Plot (checkbox unchecked)

Download Highly Variable Genes expression value (normalized)(as a .csv file) (checkbox unchecked)

Submit (button) First press the 'Submit' button before pressing the 'Find' button

Highly Variable Genes Detected: 2000

Show 10 entries Search:

	mean	variance	Standardized.variance
PPBP	0.248673237300986	9.79555189849471	11.2721864976797
LYZ	10.4082638362396	575.359613139306	8.42747141249425
S100A9	6.14632297194845	284.275130577484	8.41802141460753
IGLL5	0.276724791508719	9.00947702017393	8.11569480212293
GNLY	1.59931766489765	46.2023051846024	7.82994548132357
FTL	27.9996209249431	2042.63519136173	7.67322220326353
PF4	0.111827141774071	2.78271179111742	7.18497629148837
FTH1	21.535633055345	1224.96630659568	6.6361658659417
GNG11	0.0625473843821077	0.742009738937032	5.59538329864408
S100A8	3.17475360121304	79.3380502805618	5.58046242634516

Showing 1 to 10 of 2,000 entries

Figure 9. Normalization

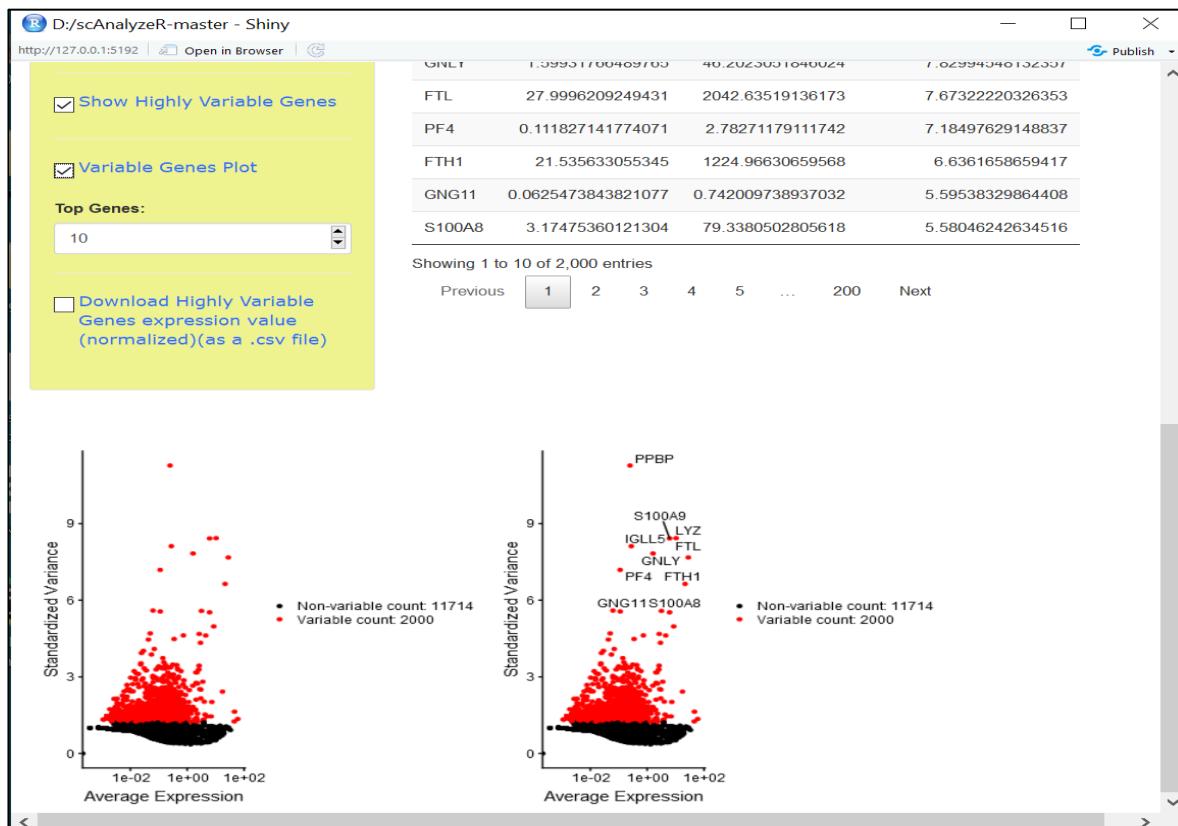


Figure 9. Normalization: Variable genes plot (top 10)

6. Dimensionality Reduction

Now we perform the Principal Component Analysis (PCA) on normalized data. To compute PCA, either use previously (on Normalization module) identified highly variable Genes or all Genes as input. The number of PCs (20 by default) will be computed via clicking the “Compute PCA” button (**Fig. 10**). Various plots can be shown if check the individual checkbox, e.g., PCA plot, t-SNE plot, Elbow plot, etc. However, the ‘PCA plot’ draws a 2D scatter plot with PC1 and PC2, where each dot represents a cell. The ‘t-SNE plot’ graphs the output from applying the nonlinear dimensionality reduction technique (**Fig. 11**), and the ‘Print PCA’ presents all genes in an individual PC as well as remarks which genes are belonging to the positive and negative PC scores. It is recommended to choose the right number of PCs that will be processed to the clustering module. The ‘Elbow plot’ and ‘Jack Straw plot’ will help you to identify the correct number of PC(s). The ‘Elbow plot’ ranking PCs according to percentages of variance, each black dot represents a single PC. In our example, we can see there is no significant drop after the PC10; it is concluding that the first 10 PCs (PC1-10) are covering the maximum number of true signals (**Fig. 12**). On the other hand, the ‘Jack Straw plot’ shows a uniform distribution for each PC with their p-values. Curves (above the black dashed line) with low p-values indicate significant PCs. In our demonstration, there is a significant drop after the first 10-12 PCs. We recommend that repeat the PCA procedure with different number of PCs, and try to compute PCA with 10 or more number of PCs.

The PC heatmap helps to explore the primary sources of heterogeneity in the dataset. To draw the heatmap plot, check the “PC Heatmap” box, then the interface panel will be visible with the default parameters setting. The user can set different values, if necessary. After setting arguments, click the “Show plot” button to display the heatmap. The heatmap shows a correlation for an equal number of genes from +ve and -ve PC’s sets. In this case, it shows 30 genes (15 with positive PC scores and another 15 with negative PC scores) in the first two principle components (PC1 and PC2) (**Fig. 13**).

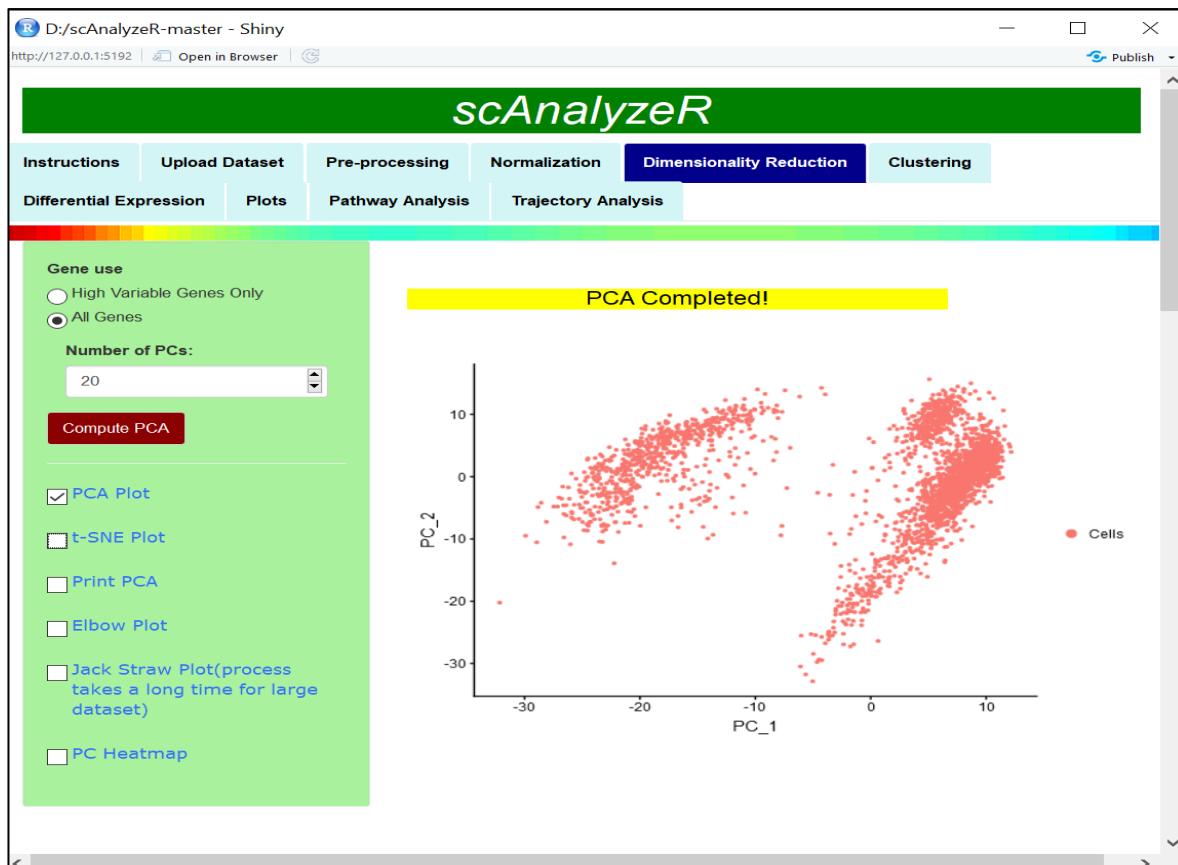


Figure 10. Dimensionality Reduction: Principal component analysis (PCA) for 20 PCs

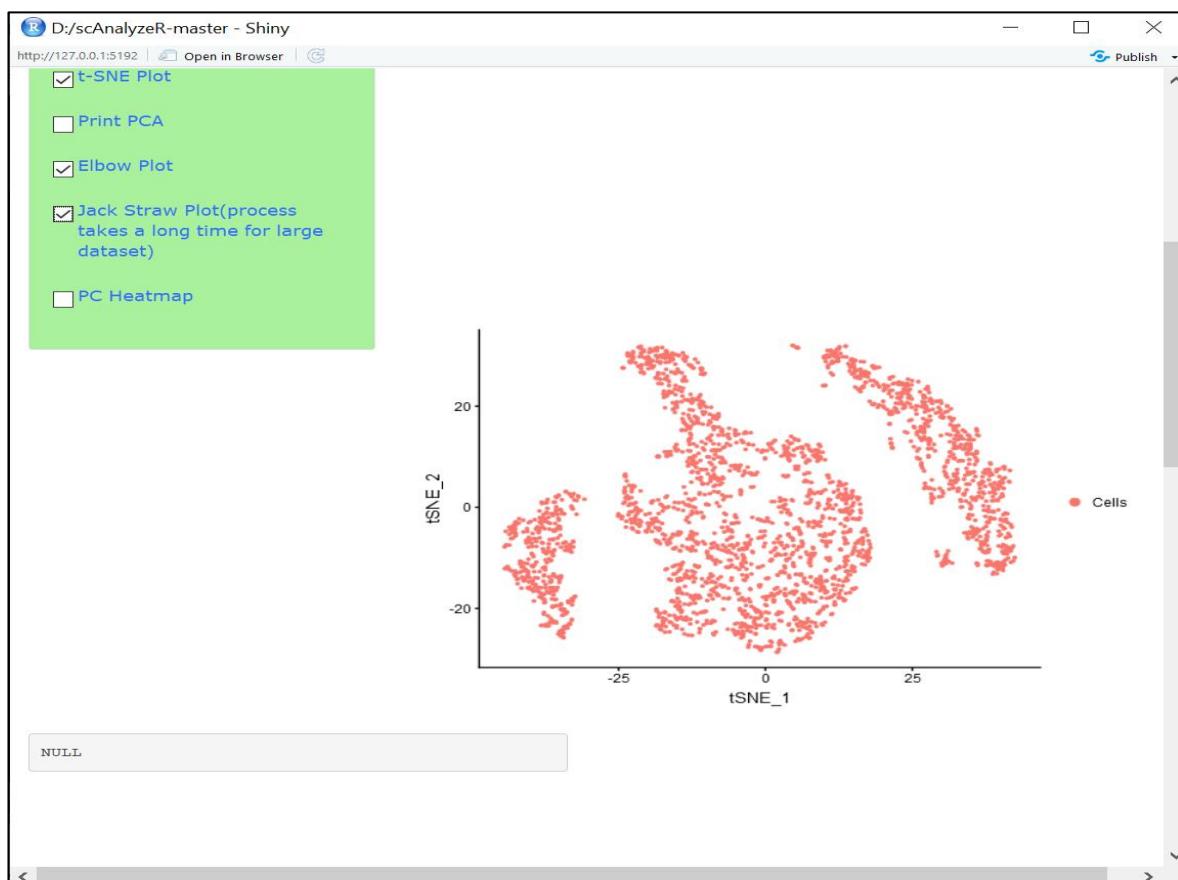


Figure 11. Dimensionality Reduction: t-SNE plot

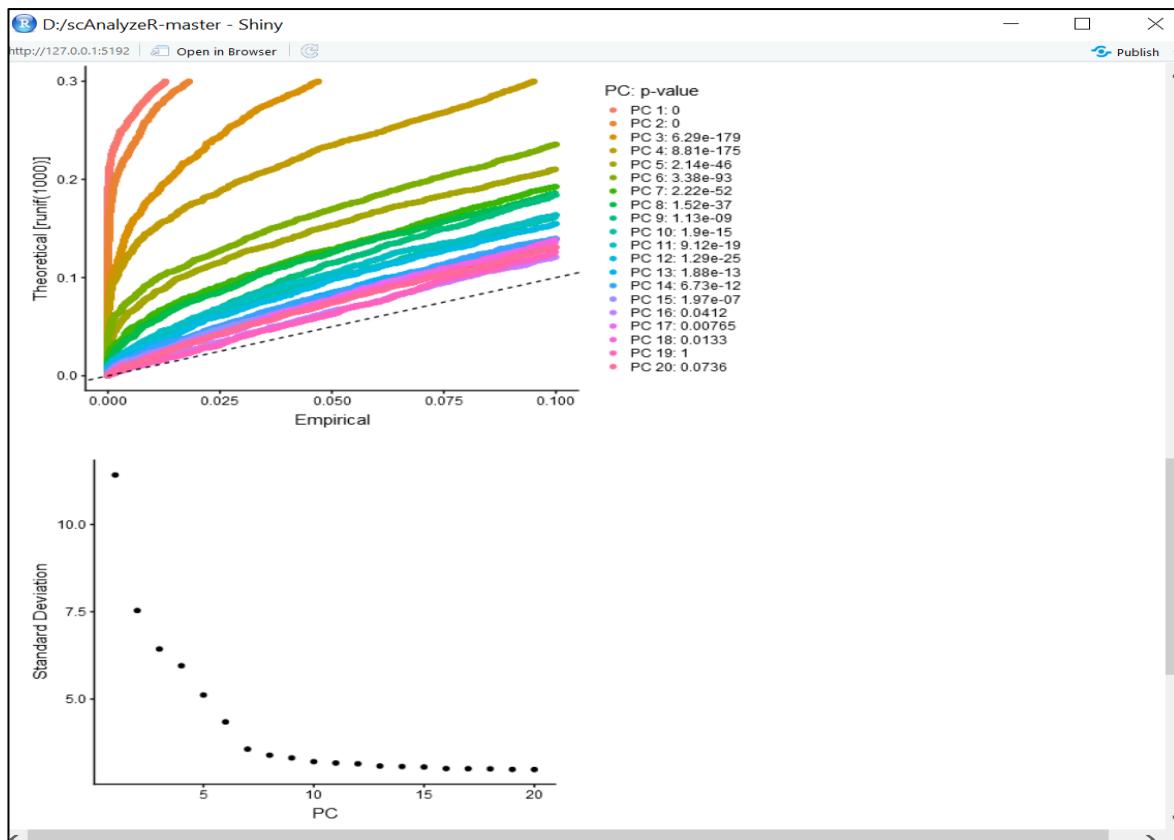


Figure 12. Dimensionality Reduction: Jack Straw plot (top), and Elbow plot (bottom) for 20 PCs

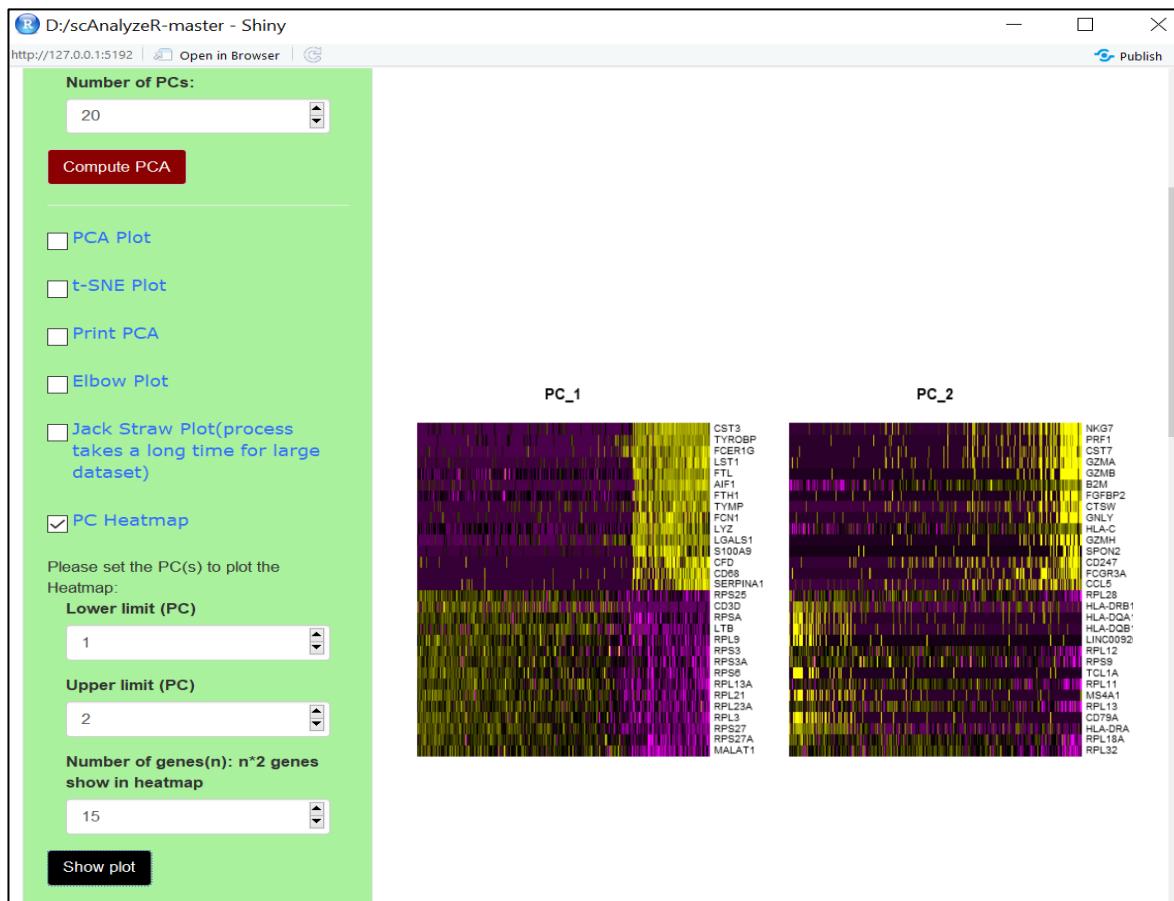


Figure 13. Dimensionality Reduction: Heatmap for top two PCs (PC1, PC2)

7. Clustering

In this module, cells are clustered based on pre-calculated PCs (PC1-PC10). The user can make clusters using either default parameter values or different values. It is highly recommended that use only the significant PCs for better results, for instance, we use first 10 PCs for performing clustering cells. The default value 0.5 is provided for the resolution field, and the minimum value of the resolution is 0 (zero), if you increase the resolution, it will also increase the number of clusters. Generally, setting this value between 0.4 -1.2 gives good results where the dataset contains about 3k cells. In this case, we set the resolution to 0.4, and the default clustering method (*Louvain algorithm*) was selected from three available clustering algorithms (Fig. 14). However, it is necessary to click the “Do cluster” button to perform the clustering. After pressing the button, then click checkboxes to show the clusters’ t-SNE plot and bar plot (Fig. 15). It is mentioned that this clustering results will be applied to later modules.

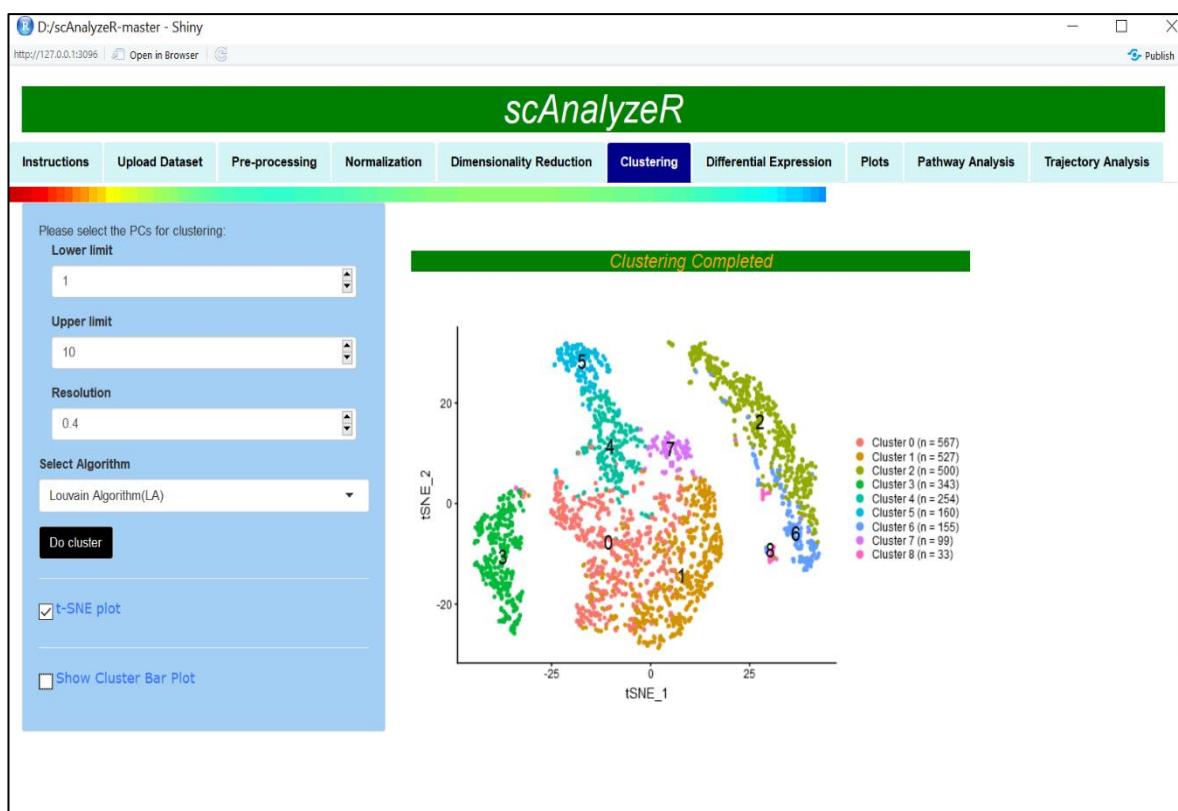


Figure 14. Clustering: t-SNE plot



Figure 15. Clustering: Clustering Bar plot

8. Differential Expression Analysis

In differential expression analysis, you can find markers (differentially expressed genes) that define clusters. The differential expression module is divided into three submodules, namely, ‘find all markers,’ ‘find markers by cluster,’ and ‘find markers by cluster vs other clusters’.

To find differentially expressed genes for all clusters, click the “Find All markers” button (**Fig. 16-27**). You can choose a different test method from the “Select Test Method” drop-down menu, and then press the button to find markers. The “Save all markers as csv” button and markers bar plot will be shown when the marker finding process is completed. You can also save the full list of markers as a csv file. To show the full table of marker list, check the “Show All Markers” checkbox, and check the “Filtering Markers” checkbox for filtering the generated marker list. After checking the filtering checkbox, the filtering interface panel will be visible with necessary default setting parameters. If you choose “Positive only” from “Markers Selection,” it means that both the “Show Filtered list” and “Show Top genes” also have the

same selection from the filtered list. Similar action will be done if the “Negative Only” is selected. You can see the top differentially expressed genes (10 by default) from each cluster clicking via “Show Top genes” as well as press a particular download button (‘Save All Top genes as a text file,’ ‘Save +ve Top genes as a text file’ or ‘Save -ve Top genes as a text file’) to save the top markers list.

In the ‘find marker by cluster’ submodule (**Fig. 28- 35**), a list of markers will be identified in a particular cluster (cluster ‘0’ is selected by default), compared to all remained clusters. For our example, we used cluster ‘0’, compared to all other clusters. The marker list will be shown and saved, which is similar to the ‘Find all markers’ submodule. Besides, a heatmap plot will be displayed clicking via the “Show Heatmap” checkbox for that particular cluster’s top markers as well as you can enlarge the heatmap plot for that specific cluster clicking on the “Zoom the cluster” checkbox.

The ‘find markers by cluster vs other clusters’ submodule (**Fig. 36-41**) can detect markers by cluster(s) versus cluster(s). First, write the cluster(s) number in both textboxes, i.e., “Select cluster/s” and “Select complementary cluster/s”. The cluster number must be separated by a comma, and no common cluster number is allowed as an input for both textboxes. If the same cluster number is written in both input boxes, it will be shown an error message, “*Error: No features pass logfc.threshold threshold*”. Other functionalities of this module are identical to the ‘find marker by cluster’ submodule. In this case, we used clusters ‘0’, and ‘1’, compared with clusters ‘2’, ‘6’, and ‘8’.

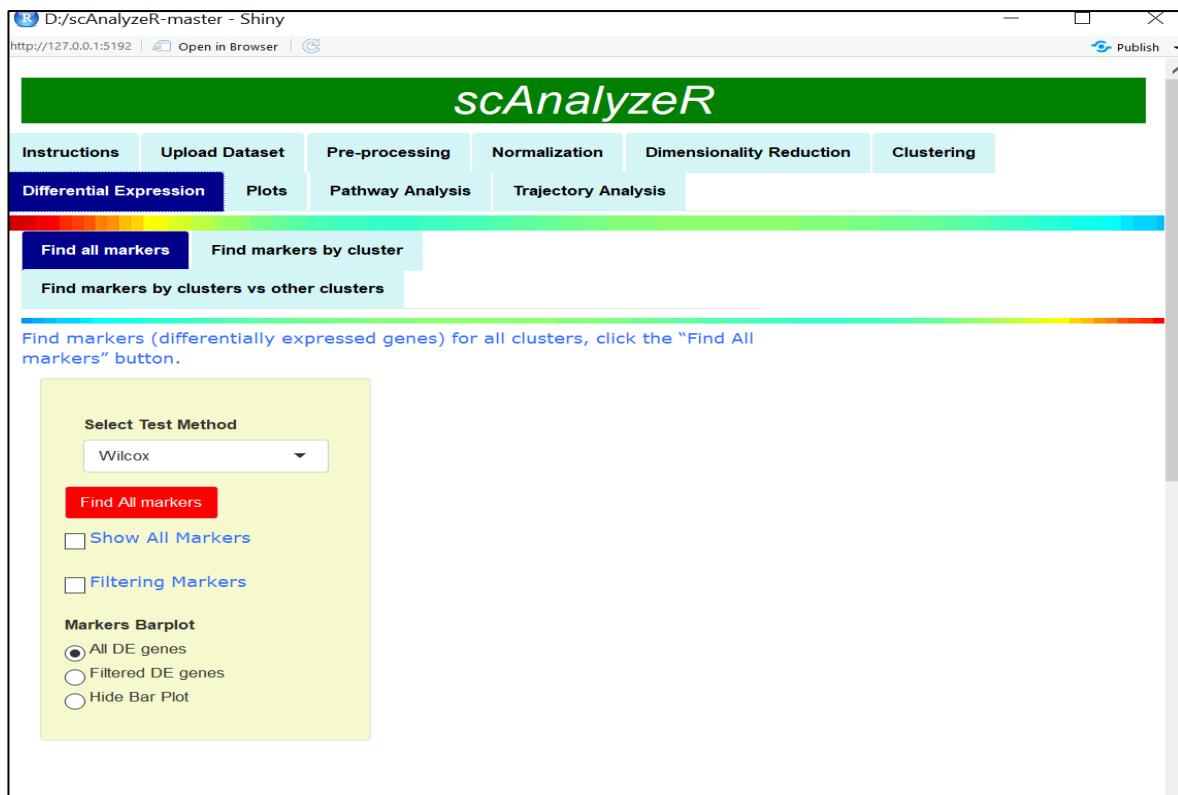


Figure 16. Differential Expression (Find all markers): the user interface

	p_val	avg_logFC	pct.1	pct.2	p_val_adj	cluster
IL32	6.3432747883065e-112	0.899295497529341	0.937	0.447	8.69916704468353e-108	0
LTB	4.02632739336073e-104	0.925974614479914	0.974	0.63	5.52170538725491e-100	0
IL7R	2.49122850646073e-87	0.857734806576934	0.751	0.307	3.41647077376025e-83	0
CD3D	4.26263800013313e-81	0.645704626962905	0.91	0.413	5.84578175338257e-77	0
LDHB	5.0841209962221e-76	0.643469073230949	0.935	0.604	6.97236353421914e-72	0
TNFRSF4	5.04253805656252e-72	0.704146064769952	0.215	0.015	6.91533669076984e-68	0
HLA-DRA	1.09662727290198e-68	-2.44013653865999	0.33	0.623	1.50391464205777e-64	0
AQP3	2.08570008109926e-64	0.866048540188016	0.404	0.102	2.86032909121952e-60	0
CD2	2.37931130617026e-64	0.860269070964188	0.621	0.236	3.2629875252819e-60	0
HLA-DRB1	7.73596834515611e-64	-1.78437131912177	0.183	0.535	1.06091069885471e-59	0

Figure 17. Differential Expression (Find all markers): The table shows all DE genes (up & down-regulated) before filtering

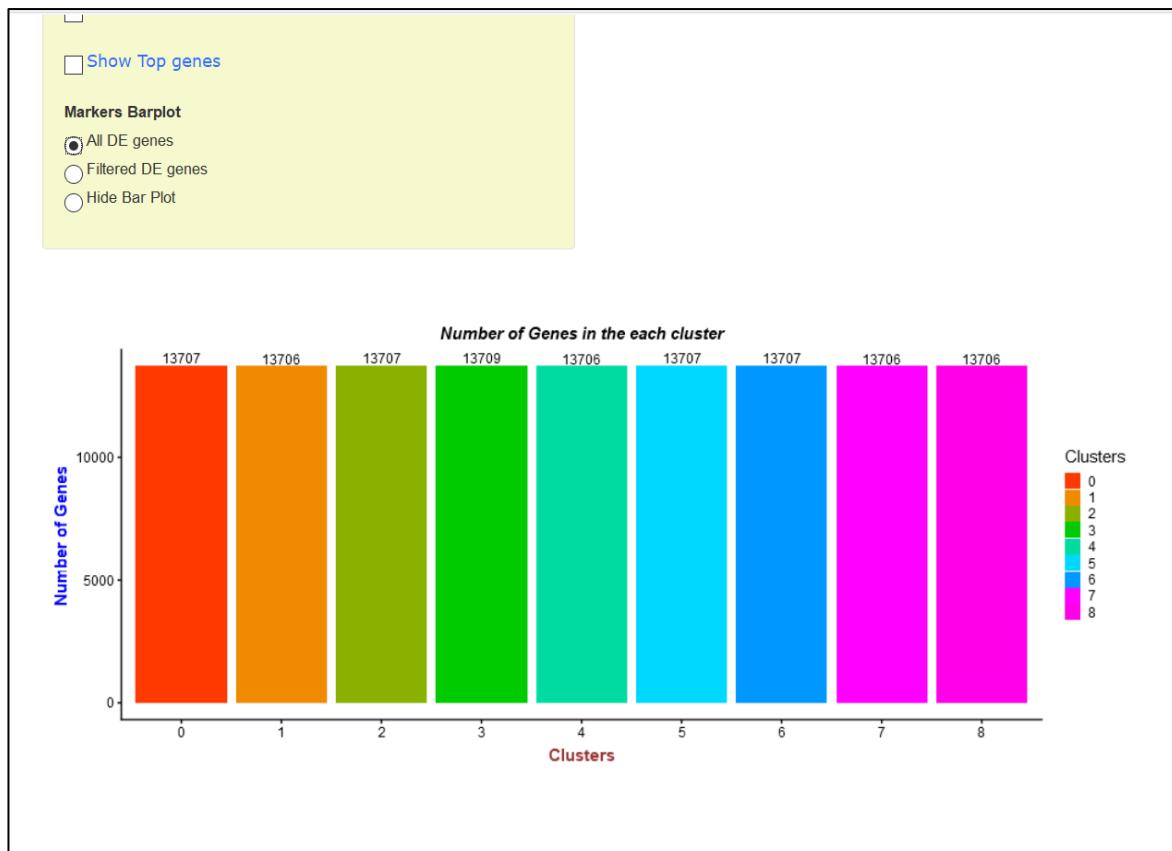


Figure 18. Differential Expression (Find all markers): The Bar plot shows number of DE genes identified in each cluster before filtering

The figure displays a table of filtered differential expression markers. The left sidebar contains filtering options: "Find All markers" (red button), "Save All markers as a csv" (button), "Show All Markers" (checkbox), "Filtering Markers" (checkbox checked), "Avg.logFC threshold" (input 0.25), "Min % (min.pct)" (input 0.25), "Adjusted p-value(p_val_adj)<input value" (input 0.05), "Markers Selection" (radio buttons for Positive only and Negative only, with Negative only selected), "Show filtered list" (checkbox checked), and "Show Top genes" (checkbox). The main area shows a table with columns: p_val, avg_logFC, pct.1, pct.2, p_val_adj, cluster, and gene. The table lists 10 rows of data, each with a unique ID (1-10), p-values, logFC values, and cluster/gene names. The table includes navigation buttons for "Previous", "Next", and page numbers 1 through 62.

	p_val	avg_logFC	pct.1	pct.2	p_val_adj	cluster	gene
1	1.09662727290198e-68	-2.44013653865999	0.33	0.623	1.50391464205777e-64	0	HLA-DRA
2	4.50233337988743e-63	-1.6879509312323	0.808	0.854	6.17449999717762e-59	0	CD74
3	7.5263071937534e-58	-1.70634580319977	0.254	0.559	1.03215776855134e-53	0	HLA-DPA1
4	1.7654793183345e-52	-1.69831911825197	0.309	0.576	2.42117833716393e-48	0	HLA-DPB1
5	3.35196935606943e-48	-0.779374769703536	0.831	0.852	4.59689077491361e-44	0	CYBA
6	4.89322728738516e-42	-1.57658411608557	0.989	0.99	6.71057190192001e-38	0	FTL
7	7.03136724290577e-34	-1.2497820309464	0.39	0.559	9.64281703692097e-30	0	CTSS
8	7.21957419610659e-34	-0.686692716157707	0.921	0.91	9.90092405254057e-30	0	OAZ1
9	2.58252342942545e-33	-0.2881408686262775	0.998	0.994	3.54167263111407e-29	0	MT-CO1
10	3.83589126438756e-30	-2.50839989746729	0.513	0.63	5.26054127998111e-26	0	LYZ

Figure 19. Differential Expression (Find all markers): The table shows down-regulated DE genes after filtering

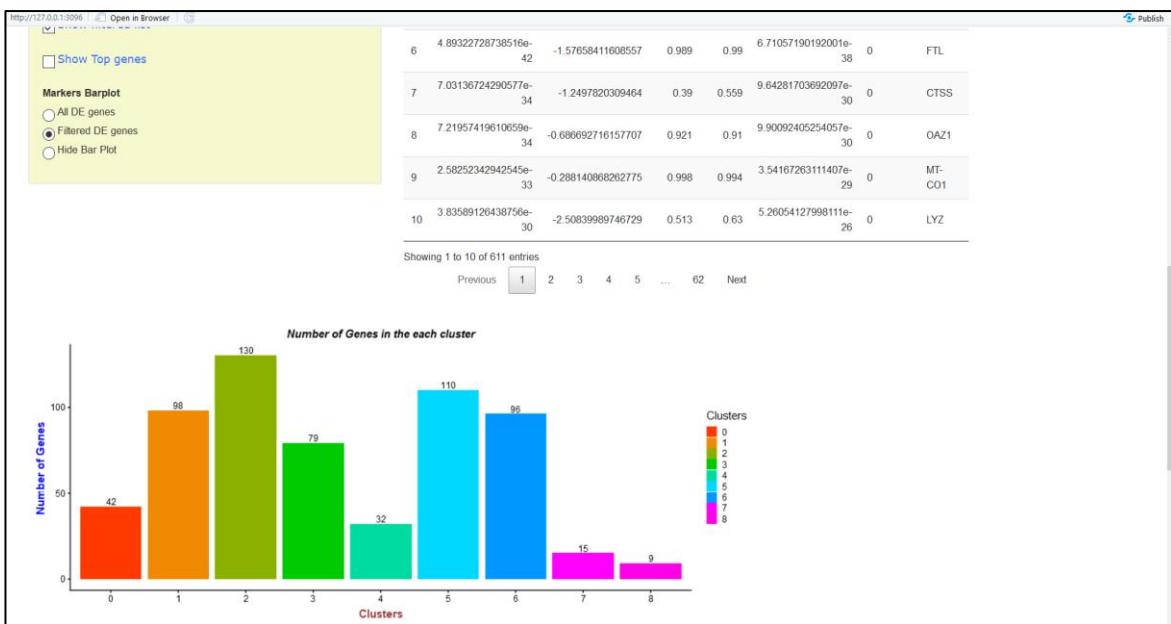


Figure 20. Differential Expression (Find all markers): The bar plot shows number of down-regulated DE genes in each cluster after filtering

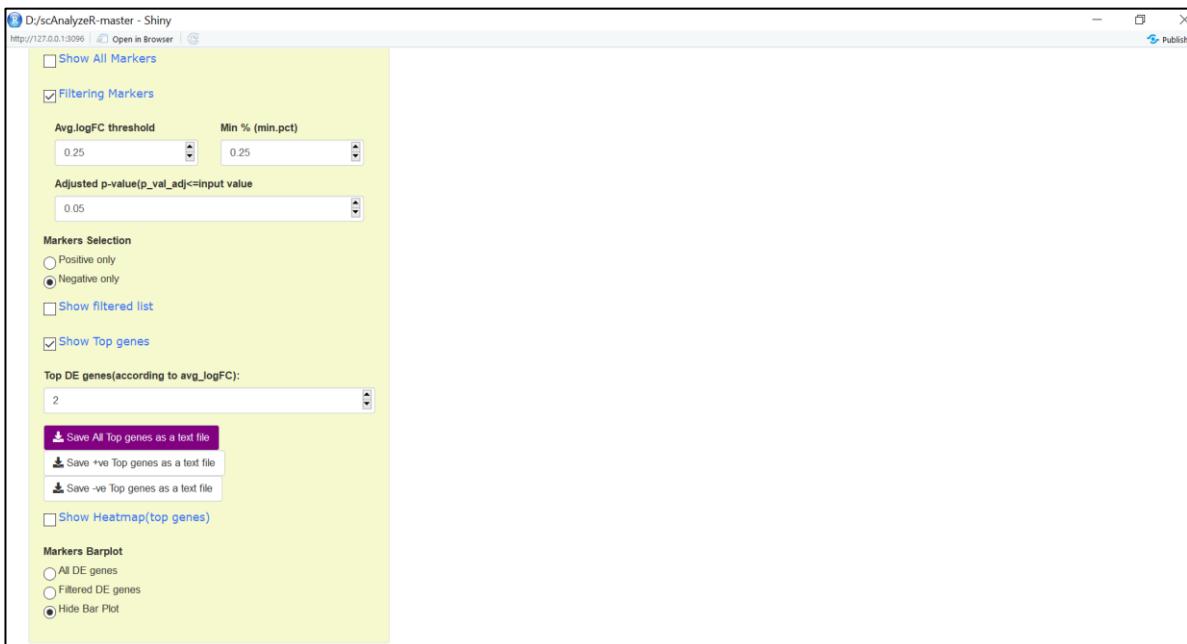


Figure 21. Differential Expression (Find all markers): DE genes downloading options after filtering

D:/scAnalyzeR-master - Shiny
http://127.0.0.1:3096 | Open in Browser |

Show 25 entries Search:

	p_val	avg_logFC	pct.1	pct.2	p_val_adj	cluster	gene
1	3.83589126438756e-30	-2.50839989746729	0.513	0.63	5.26054127998111e-26	0	LYZ
2	1.09662727290198e-68	-2.44013653865999	0.33	0.623	1.50391464205777e-64	0	HLA-DRA
3	1.89531700178944e-74	-2.53385795064068	0.266	0.634	2.59923773625404e-70	1	HLA-DRA
4	6.11349718759364e-28	-2.49881216732877	0.512	0.628	8.38405004306592e-24	1	LYZ
5	5.09488632846149e-80	-1.38524404657505	0.254	0.682	6.98712711085209e-76	2	CXCR4
6	3.12398799449261e-219	-1.17019359055668	0.998	1	4.28423713564717e-215	2	MALAT1
7	8.50149899306664e-20	-2.38353621052836	0.423	0.632	1.16589557190916e-15	3	LYZ
8	2.05409797385055e-112	-2.16700044089328	0.359	0.881	2.81698996133864e-108	3	S100A4
9	8.5564008766643e-15	-2.42197344937327	0.445	0.622	1.17342481622574e-10	4	LYZ
10	5.20061109956522e-23	-2.0738014768374	0.343	0.583	7.13211806194374e-19	4	HLA-DRA
11	1.40046088867604e-9	-2.39601242903279	0.444	0.615	0.0000192059206273032	5	LYZ
12	6.25389216534243e-14	-2.02287760696784	0.344	0.574	8.5765877155506e-10	5	HLA-DRA
13	7.30307263399988e-19	-1.42774688114116	0.658	0.707	1.00154338102674e-14	6	LTB
14	1.32495625626882e-29	-1.38676628667677	0.381	0.694	1.81704500984706e-25	6	LDHB
15	2.92739588178109e-9	-2.33766097855824	0.273	0.618	0.0000401463071227459	7	LYZ
16	1.40477504530937e-8	-1.3760807527351	0.899	0.993	0.000192650849713727	7	FTL
17	0.00000205842199188706	-1.56102874632868	0.515	0.706	0.0282291991967391	8	LTB
18	8.14567874747102e-7	-1.34384751083004	0.273	0.646	0.0111709838342818	8	PTPRCAP

Showing 1 to 18 of 18 entries Previous Next

Figure 22. Differential Expression (Find all markers): The table shows top-two down-regulated DE genes in each cluster after filtering

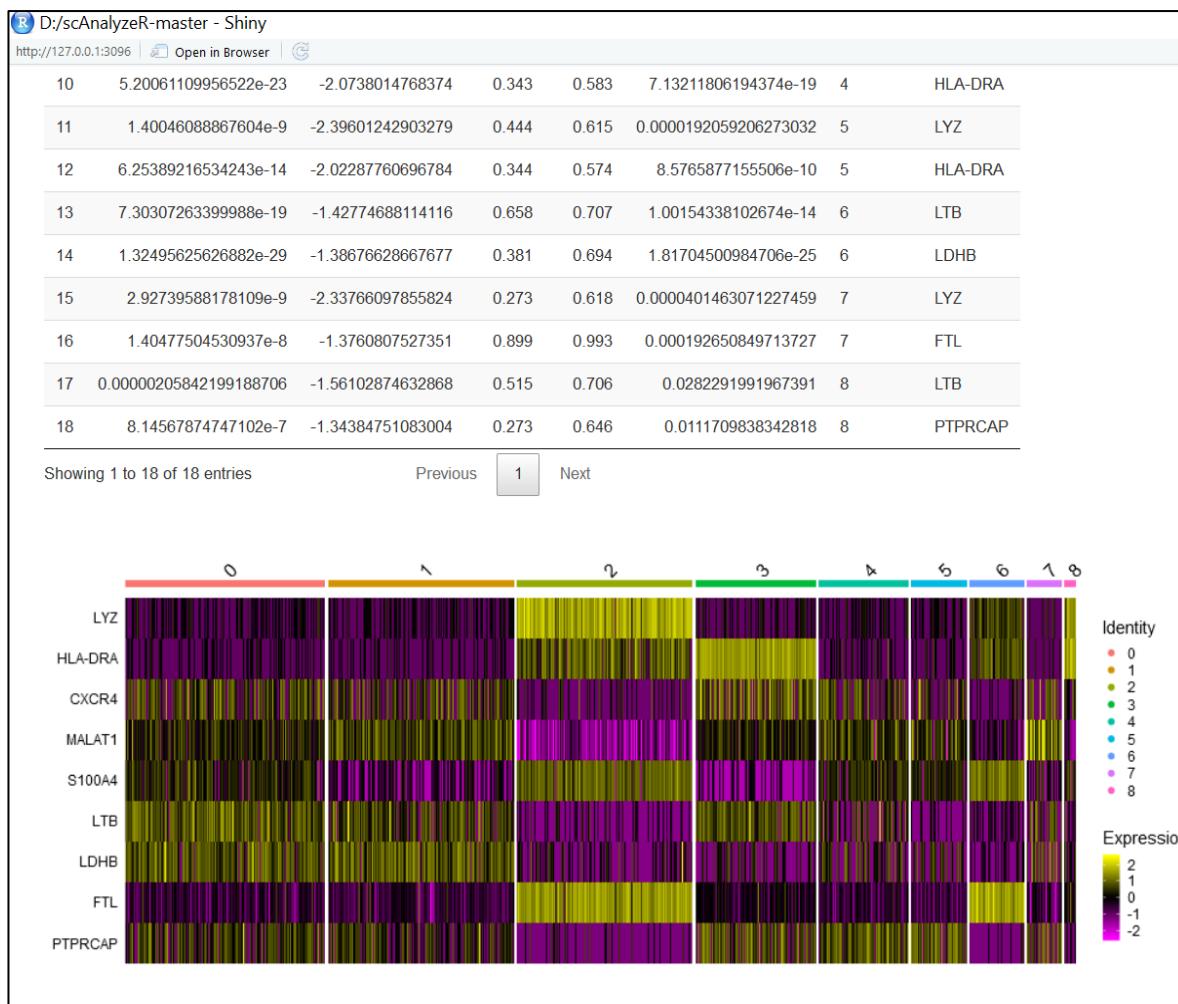


Figure 23. Differential Expression (Find all markers): The heatmap illustrates expressions for top-two down-regulated DE genes (unique genes only) in each cluster after filtering

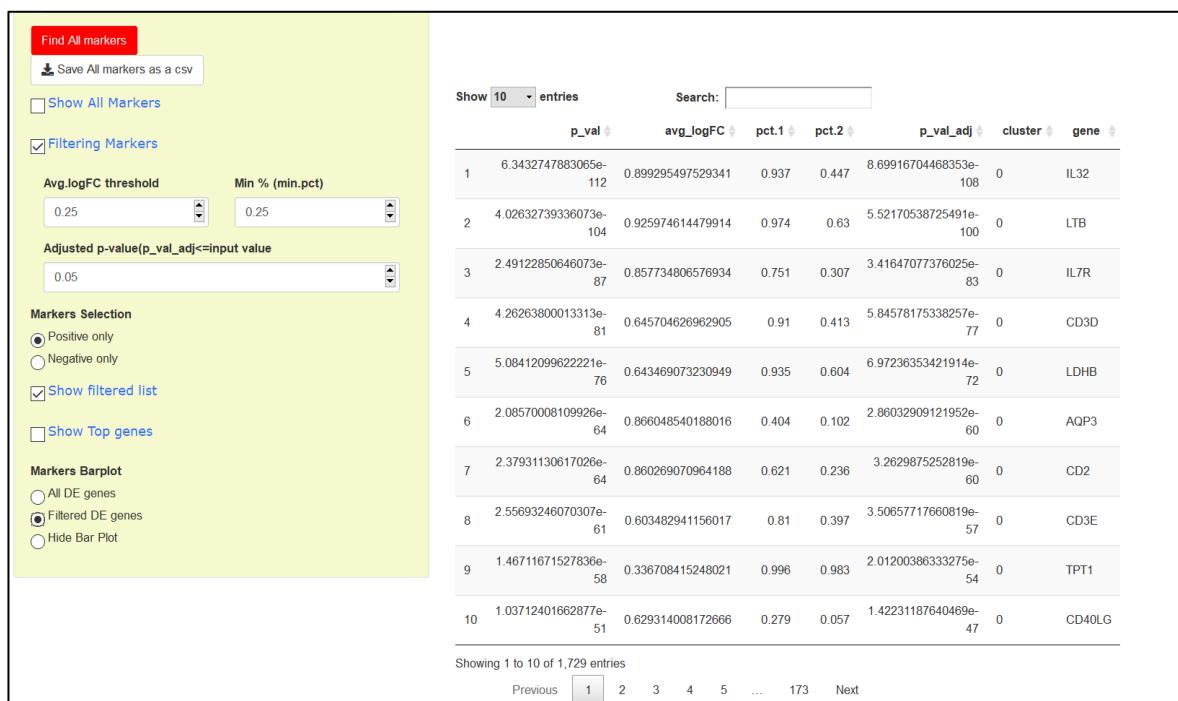


Figure 24. Differential Expression (Find all markers): The table shows up-regulated DE genes after filtering

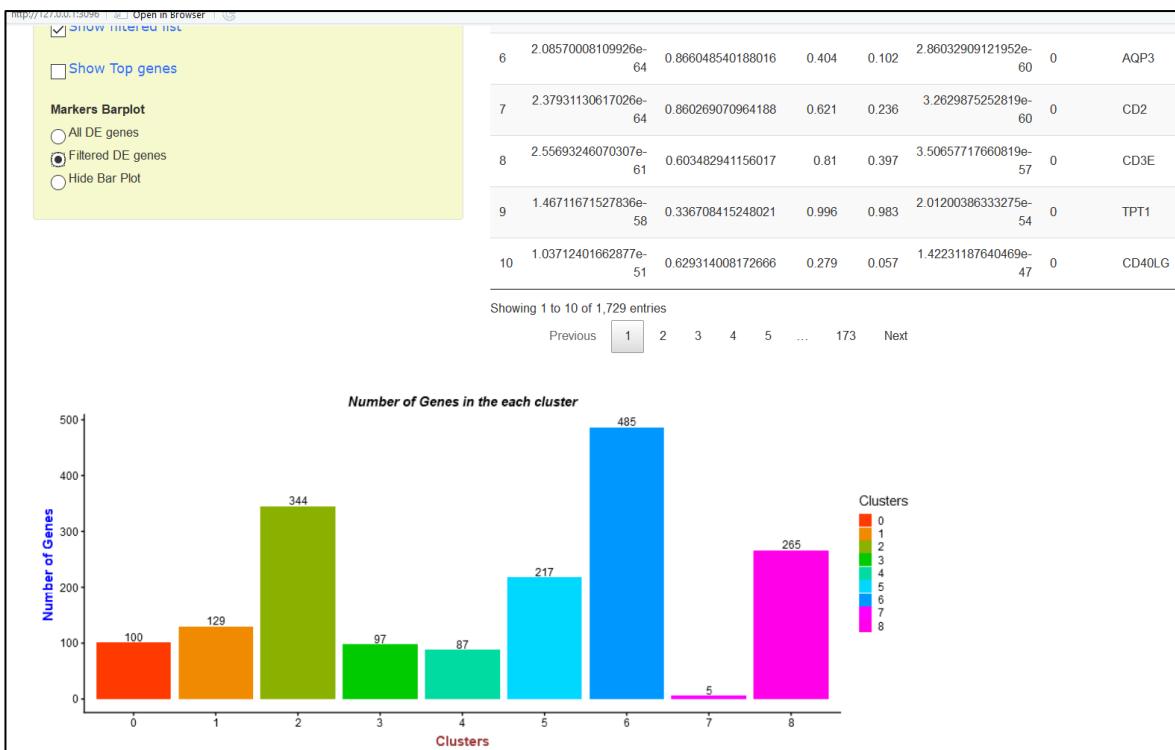


Figure 25. Differential Expression (Find all markers): The bar plot shows number of up-regulated DE genes in each cluster after filtering

	p_val	avg_logFC	pct.1	pct.2	p_val_adj	cluster	gene
1	4.02632739336073e-104	0.925974614479914	0.974	0.63	5.52170538725491e-100	0	LTB
2	6.3432747883065e-112	0.899295497529341	0.937	0.447	8.69916704468353e-108	0	IL32
3	2.08570008109926e-64	0.866048540188016	0.404	0.102	2.86032909121952e-60	0	AQP3
4	2.37931130617026e-64	0.860269070964188	0.621	0.236	3.2629875252819e-60	0	CD2
5	2.49122850646073e-87	0.857734806576934	0.751	0.307	3.41647077376025e-83	0	IL7R
6	6.1637788434081e-118	1.08595287004588	0.545	0.109	8.45300630584987e-114	1	CCR7
7	3.18621333071222e-59	0.780073922571657	0.395	0.108	4.36957296173874e-55	1	LEF1
8	6.25003656699994e-64	0.761293675885978	0.26	0.037	8.57130014798371e-60	1	FHIT
9	1.72036901029043e-108	0.743469321521585	0.966	0.603	2.35931406071229e-104	1	LDHB
10	9.15755184532473e-55	0.741896070153267	0.391	0.113	1.25586666006783e-50	1	PRKCQ-AS1

Showing 1 to 10 of 45 entries

Figure 26. Differential Expression (Find all markers): The table shows top-five up-regulated DE genes in each cluster after filtering

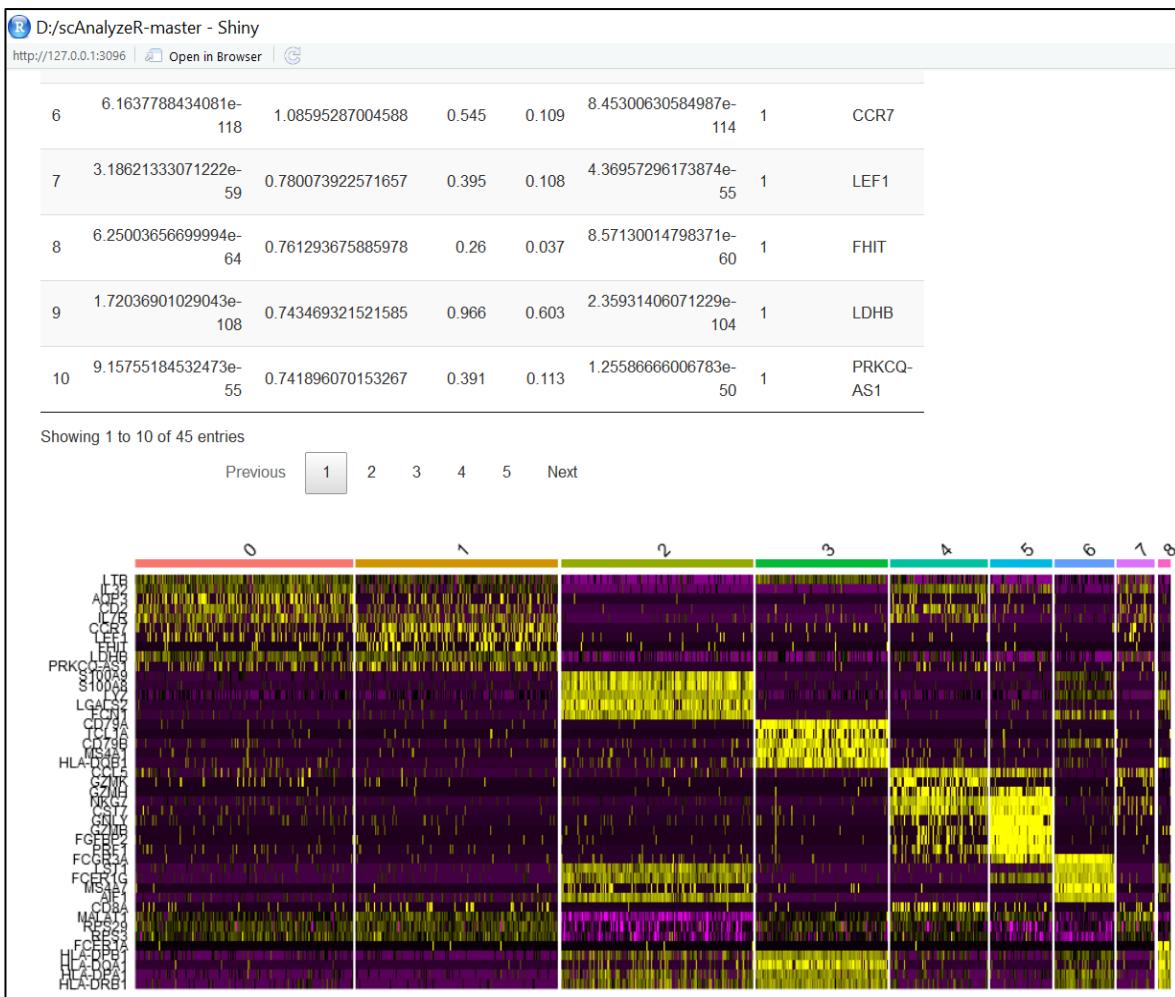


Figure 27. Differential Expression (Find all markers): The heatmap illustrates expressions for top-five up-regulated DE genes (unique genes only) in each cluster after filtering

Please select parameters:

Select cluster logfc threshold

0 0.25

Min % (min.pct)

0.25

Select Test Method

Wilcox

Show All Markers(both positive and negative)

Filtering Markes

A list of markers will be identified in a particular cluster (cluster 0 is selected by default), compared to all remained clusters.

Figure 28. Differential Expression (Find markers by cluster): the user interface

	Show 10 entries	Search:	p_val	avg_logFC	pct.1	pct.2	p_val_adj
IL32	6.3432747883065e-112	0.899295497529341	0.937	0.447	8.69916704468353e-108		
LTB	4.02632739336073e-104	0.925974614479914	0.974	0.63	5.52170538725491e-100		
IL7R	2.49122850646073e-87	0.857734806576934	0.751	0.307	3.41647077376025e-83		
CD3D	4.26263800013131e-81	0.645704626962905	0.91	0.413	5.84578175338257e-77		
LDHB	5.08412099622221e-76	0.643469073230949	0.935	0.604	6.97236353421914e-72		
HLA-DRA	1.09662727290198e-68	-2.44013653865999	0.33	0.623	1.50391464205777e-64		
AQP3	2.08570008109926e-64	0.866048540188016	0.404	0.102	2.86032909121952e-60		
CD2	2.37931130617026e-64	0.860269070964188	0.621	0.236	3.2629875252819e-60		
HLA-DRB1	7.73596834515611e-64	-1.78437131912177	0.183	0.535	1.06091069885471e-59		
CD74	4.5023337988743e-63	-1.8879509312323	0.808	0.854	6.17449999717762e-59		

Showing 1 to 10 of 218 entries

Figure 29. Differential Expression (Find markers by cluster): The table shows DE genes (both up & down-regulated) in cluster '0'

D:/scAnalyzeR-master - Shiny
http://127.0.0.1:3096 | Open in Browser |

Please select parameters:

Select cluster	logfc threshold
0	0.25

Min % (min.pct)

0.25

Select Test Method

Wilcox

Find Markers

Save All markers as a csv

Show All Markers(both positive and negative)

Filtering Markes

Adjusted p-value(p_val_adj<=input value)

0.05

Markers Selection

Positive only

Negative only

Show Heatmap

Show Filtered Markers

Show 10 entries Search:

Gene	p_val	avg_logFC	pct.1	pct.2	p_val_adj
2 LTB	4.02632739336073e-104	0.925974614479914	0.974	0.63	5.52170538725491e-100
1 IL32	6.3432747883065e-112	0.899295497529341	0.937	0.447	8.69916704468353e-108
6 AQP3	2.08570008109926e-64	0.866048540188016	0.404	0.102	2.86032909121952e-60
7 CD2	2.37931130617026e-64	0.860269070964188	0.621	0.236	3.2629875252819e-60
3 IL7R	2.49122850646073e-87	0.857734806576934	0.751	0.307	3.41647077376025e-83
4 CD3D	4.26263800013313e-81	0.645704626962905	0.91	0.413	5.84578175338257e-77
5 LDHB	5.08412099622221e-76	0.643469073230949	0.935	0.604	6.97236353421914e-72
10 CD40LG	1.03712401662877e-51	0.629314008172666	0.279	0.057	1.42231187640469e-47
15 SPOCK2	1.2740400724193e-43	0.614542846456007	0.423	0.151	1.74721855531583e-39
8 CD3E	2.55693246070307e-61	0.603482941156017	0.81	0.397	3.50657717660819e-57

Showing 1 to 10 of 100 entries

Previous 1 2 3 4 5 ... 10 Next

Figure 30. Differential Expression (Find markers by cluster): The table shows up-regulated DE genes (in cluster '0') after filtering

D:/scAnalyzeR-master - Shiny
http://127.0.0.1:3096 | Open in Browser |

Please select parameters:

Select cluster	logfc threshold
0	0.25

Min % (min.pct)

0.25

Select Test Method

Wilcox

Find Markers

Save All markers as a csv

Show All Markers(both positive and negative)

Filtering Markes

Adjusted p-value(p_val_adj<=input value)

0.05

Markers Selection

Positive only

Negative only

Show Heatmap

Show Filtered Markers

Show 10 entries Search:

Gene	p_val	avg_logFC	pct.1	pct.2	p_val_adj
25 S100A9	1.017055043023e-29	-2.90912742632973	0.198	0.401	1.39478928600174e-25
26 S100A8	1.68892342306758e-29	-2.60624443816181	0.106	0.324	2.31618958239488e-25
24 LYZ	3.83589126438756e-30	-2.50839989746729	0.513	0.63	5.26054127998111e-26
1 HLA-DRA	1.09662727290198e-68	-2.44013653865999	0.33	0.623	1.50391464205777e-64
6 TYROBP	4.55374355483474e-55	-2.40854503171091	0.15	0.465	6.24500391110036e-51
13 CST3	1.17103532028823e-38	-2.3502295922136	0.213	0.449	1.60595783824328e-34
44 NKG7	3.4430903360272e-18	-2.19057333620457	0.175	0.333	4.7218540868277e-14
7 FCER1G	1.18760188099493e-53	-1.93252240218982	0.104	0.428	1.62867721959645e-49
11 AIF1	3.06098338657606e-40	-1.82794166344406	0.173	0.43	4.19783261635041e-36
2 HLA-DRB1	7.73596834515611e-64	-1.78437131912177	0.183	0.535	1.06091069885471e-59

Showing 1 to 10 of 83 entries

Previous 1 2 3 4 5 ... 9 Next

Figure 31. Differential Expression (Find markers by cluster): The table shows down-regulated DE genes (in cluster '0') after filtering

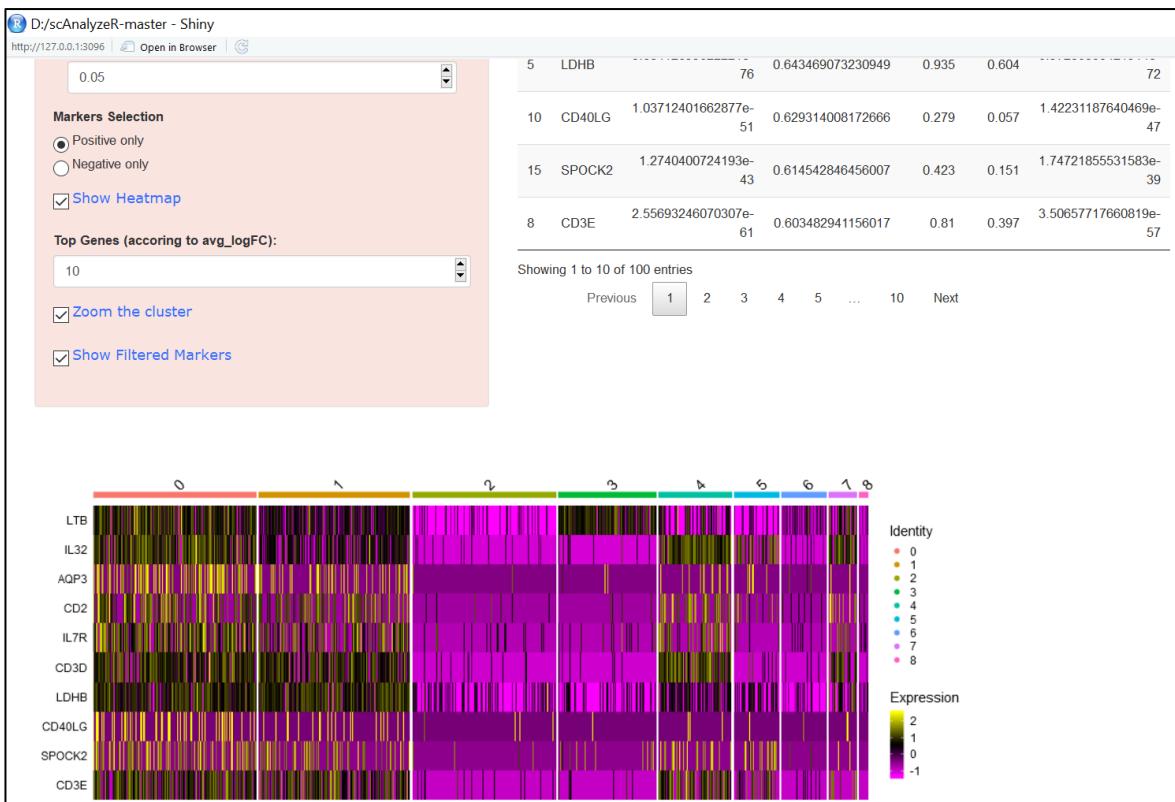


Figure 32. Differential Expression (Find markers by cluster): The heatmap illustrates expressions for top-ten up-regulated DE genes after filtering

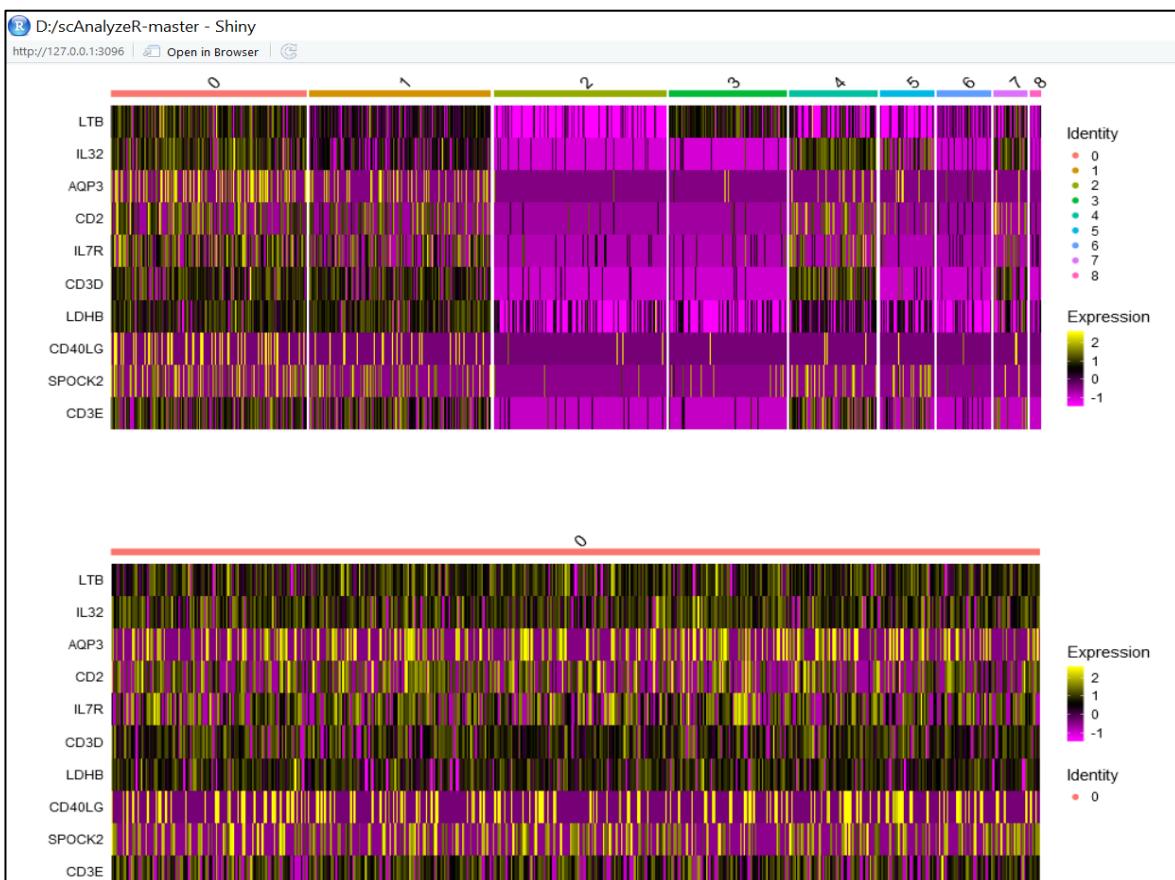


Figure 33. Differential Expression (Find markers by cluster): The heatmap(bottom) illustrates expressions for top-ten up-regulated DE genes in the cluster '0' only (zooming the cluster '0'), after filtering

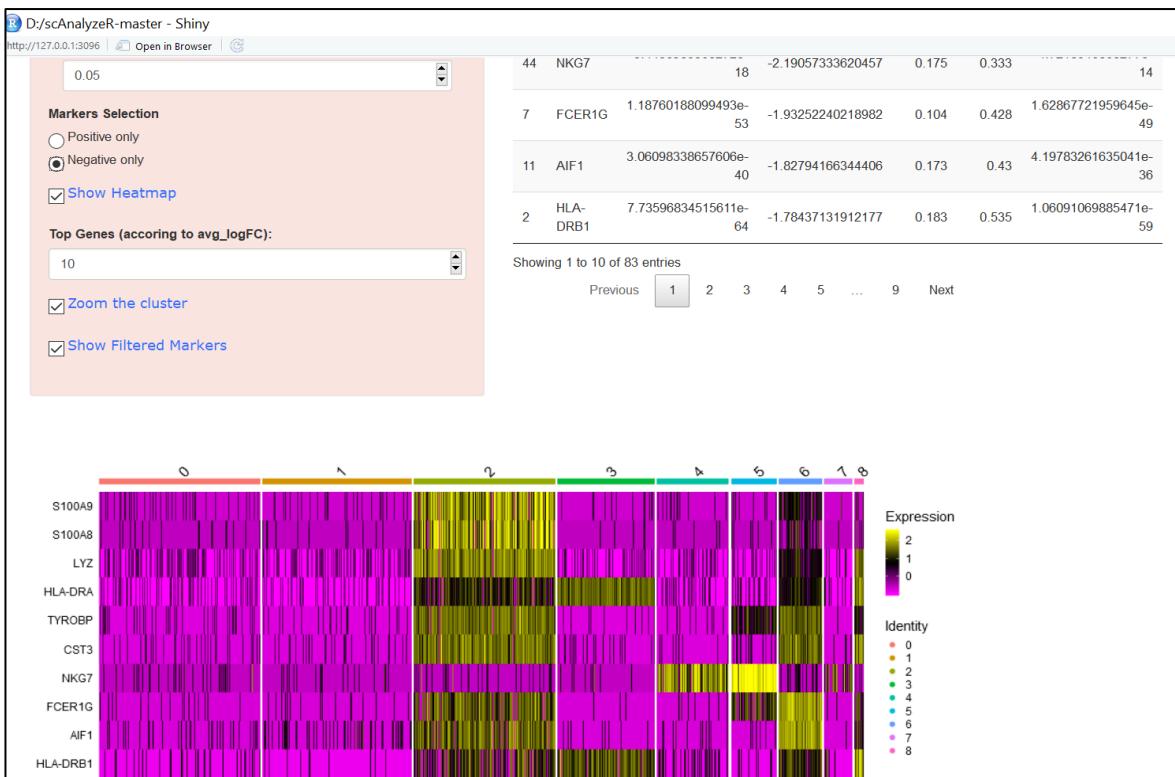


Figure 34. Differential Expression (Find markers by cluster): The heatmap illustrates expressions for top-ten down-regulated DE genes after filtering



Figure 35. Differential Expression (Find markers by cluster): The heatmap(bottom) illustrates expressions for top-ten down-regulated DE genes in the cluster '0' only (zooming the cluster '0'), after filtering

D:/scAnalyzeR-master - Shiny
http://127.0.0.1:5192 | Open in Browser | Publish

Find markers by clusters vs other clusters

Find markers by cluster(s) versus cluster(s). The cluster number must be separated by a comma, and no common cluster number is allowed as an input for both textboxes. If the same cluster number is written in both input boxes, it will be shown an error message, "Error: No features pass logfc.threshold threshold".

Please select parameters:

Select Cluster/s
The cluster number must be separated by a ,

Select complementary cluster/s
The cluster number must be separated by a ,

logfc threshold **Min % (min.pct)**
0.25 0.25

Select Test Method
Wilcox

Find markers

Show All Markers(Both Positive and Negative)
 Filtering Markers

Figure 36. Differential Expression (Find markers by clusters vs other clusters): the user interface

D:/scAnalyzeR-master - Shiny
http://127.0.0.1:5192 | Open in Browser | Publish

Find all markers **Find markers by cluster** **Find markers by clusters vs other clusters**

Find markers by cluster(s) versus cluster(s). The cluster number must be separated by a comma, and no common cluster number is allowed as an input for both textboxes. If the same cluster number is written in both input boxes, it will be shown an error message, "Error: No features pass logfc.threshold threshold".

Please select parameters:

Select Cluster/s
0,1

Select complementary cluster/s
2,6,8

logfc threshold **Min % (min.pct)**
0.25 0.25

Select Test Method
Wilcox

Find markers

Show All Markers(Both Positive and Negative)
 Filtering Markers

Show 10 entries							Search:
	p_val	avg_logFC	pct.1	pct.2	p_val_adj	Gene	
TYROBP	0	-3.37987372919936	0.137	0.978	0	TYROBP	
CST3	1.36054376443493e-305	-3.39305665473068	0.2	0.987	1.86584971854606e-301	CST3	
FCER1G	5.18309071268409e-296	-2.82390057565644	0.1	0.942	7.10809060337495e-292	FCER1G	
LST1	8.83303299598355e-283	-2.75759494753316	0.171	0.949	1.21136214506918e-278	LST1	
FTH1	1.56636726227736e-270	-2.11632961734668	0.991	1	2.14811606348717e-266	FTH1	
S100A9	3.20712460712218e-270	-4.03443466468004	0.169	0.935	4.39825068620735e-266	S100A9	
FTL	5.95213597149063e-268	-2.49824736283444	0.987	1	8.16275927130225e-264	FTL	
AIF1	1.74505546026807e-267	-2.48018862542981	0.235	0.951	2.39316905821163e-263	AIF1	
LGALS1	1.14137005858933e-265	-2.41177942750266	0.286	0.967	1.56527489834941e-261	LGALS1	
LYZ	2.23396095550141e-263	-3.5588118657203	0.513	0.985	3.06365405437463e-259	LYZ	

Showing 1 to 10 of 703 entries

Figure 37. Differential Expression (Find markers by clusters vs other clusters): The table shows both up & down-regulated DE genes

The screenshot shows the scAnalyzeR Shiny application interface. On the left, a sidebar contains input fields for 'Select Cluster/s' (0,1), 'Select complementary cluster/s' (2,6,8), 'logfc threshold' (0.25), 'Min % (min.pct)' (0.25), 'Select Test Method' (Wilcox), and checkboxes for 'Show All Markers(Both Positive and Negative)', 'Filtering Markers' (checked), 'Adjust p-value' (0.05), 'Markers Selection' (Positive only), and buttons for 'Save Filtered Markers as csv' and 'Show Heatmap'. On the right, a main panel displays a table of differential expression results with columns: Gene, p_val, avg_logFC, pct.1, pct.2, and p_val_adj. The table lists 242 entries, showing genes like IL32, LTB, CD3D, etc., with their respective statistics. A navigation bar at the bottom shows pages 1 through 25.

Gene	p_val	avg_logFC	pct.1	pct.2	p_val_adj
18 IL32	7.75165047207468e-200	2.25234713405081	0.854	0.128	1.06306134574032e-195
10 LTB	6.56485680696429e-218	2.21297148164323	0.947	0.353	9.00304462507083e-214
6 CD3D	5.21204281738511e-221	2.14887783529107	0.888	0.093	7.14779551976194e-217
29 CD3E	9.32562487067535e-172	1.76964318770807	0.796	0.11	1.27891619476442e-167
30 PTPRCAP	2.26643808327139e-168	1.76741660083204	0.799	0.131	3.10819318739839e-164
46 IL7R	7.10031460249768e-124	1.65880382565359	0.682	0.137	9.73737144586532e-120
16 LDHB	1.24305757942696e-204	1.51373142535212	0.95	0.422	1.70472916442613e-200
57 CD7	5.83534293311262e-106	1.47168928996993	0.567	0.055	8.00258929847065e-102
63 CD2	3.53154654265231e-96	1.44055107086799	0.536	0.061	4.84316292859338e-92
50 AES	5.8082706075369e-115	1.34609849142842	0.705	0.185	7.9654623111761e-111

Figure 38. Differential Expression (Find markers by clusters vs other clusters): The table shows only up-regulated DE genes after filtering

The screenshot shows the scAnalyzeR Shiny application interface. The layout is identical to Figure 38, with a sidebar for input parameters and a main panel for displaying differential expression results. The main panel table lists 439 entries, showing genes like S100A9, S100A8, LYZ, CST3, TYROBP, FCER1G, HLA-DRA, LST1, FCN1, and LGALS2, along with their statistics. Navigation is provided for pages 1 through 44.

Gene	p_val	avg_logFC	pct.1	pct.2	p_val_adj
6 S100A9	3.20712460712218e-270	-4.03443466468004	0.169	0.935	4.39825068620735e-266
20 S100A8	3.44456993696565e-227	-3.67360535942173	0.101	0.82	4.72388321155469e-223
10 LYZ	2.23396095550141e-263	-3.5588118657203	0.513	0.985	3.06365405437463e-259
2 CST3	1.36054376443493e-305	-3.39305665473068	0.2	0.987	1.86584971854606e-301
1 TYROBP	0	-3.37987372919936	0.137	0.978	0
3 FCER1G	5.18309071268409e-296	-2.82390057565644	0.1	0.942	7.10809060337495e-292
13 HLA-DRA	2.43434112698301e-251	-2.76471973908818	0.299	0.942	3.3384554215445e-247
4 LST1	8.83303299598355e-283	-2.75759494753316	0.171	0.949	1.21136214506918e-278
12 FCN1	1.59741755406685e-252	-2.70796394706965	0.11	0.871	2.19069843364728e-248
25 LGALS2	2.78514625992767e-217	-2.52750384046246	0.036	0.735	3.8195495808648e-213

Figure 39. Differential Expression (Find markers by clusters vs other clusters): The table shows only down-regulated DE genes after filtering

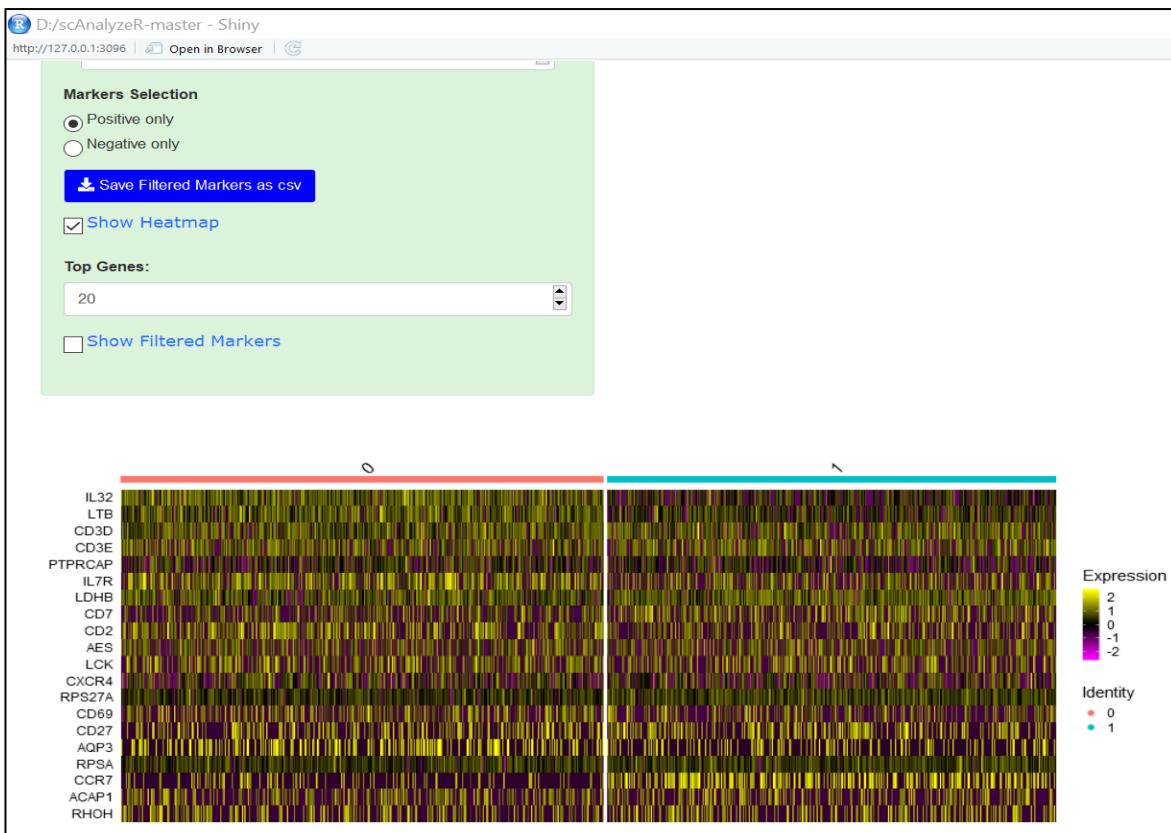


Figure 40. Differential Expression (Find markers by clusters vs other clusters): The heatmap illustrates expressions for top-twenty up-regulated DE genes after filtering

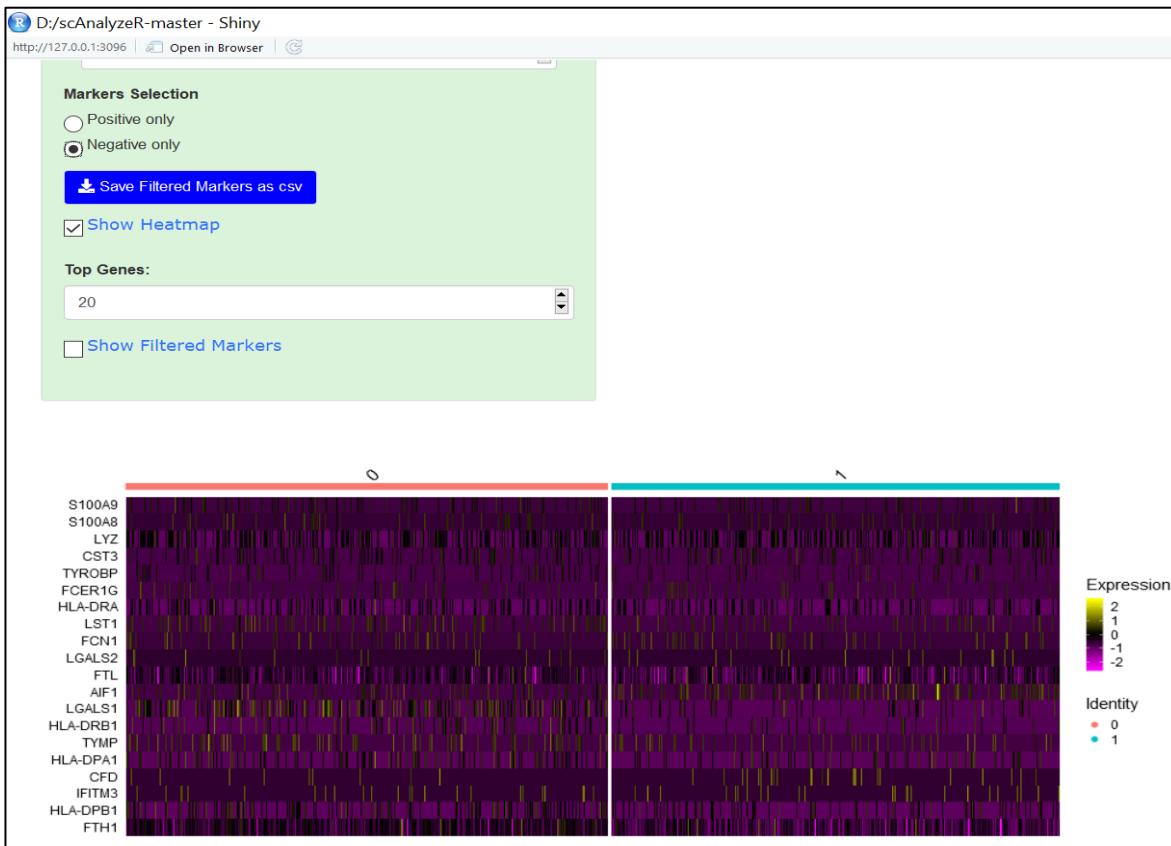


Figure 41. Differential Expression (Find markers by cluster): The heatmap illustrates expressions for top-twenty down-regulated DE genes after filtering

9. Plots

The ‘Plots’ module is divided into five essential plotting functions, i.e., Violin Plot, Feature Plot, Heatmap, and Correlation Plot. In the “Violin Plot” textbox, write one or more gene symbol(s) (one gene per line), then the violin plot will be shown automatically, which represents the gene expression across clusters (**Fig. 42-43**). The “Feature Plot” visualizes the gene expression on a tSNE plot (**Fig. 44**). The input style is the same as the “Violin Plot.”

There are two different input options for the “Heatmap,” you can either upload a text file that contains a set of gene symbol (one gene per line) or write gene symbols in the textbox that is provided on the interface, then the expression heatmap plot (by default Heatmap with clusters) will be shown for provided genes with all cells (after filtering, in pre-processing module) (**Fig. 45**). You may choose the “Heatmap without clusters” radio button that generates an expression heatmap for the same input genes without cluster indication (**Fig. 46**). To draw the dynamic PCA plot (**Fig. 47**), upload a set of genes in a text file (one gene symbol per line), then the PCA plot will be displayed. There are several advantages of the dynamic PCA plot, e.g., zooming a particular area (**Fig. 48-49**), zoom out, box and lasso selection, etc. The “Correlation Plot” submodule shows a correlation between two given genes. The ‘GAPDH’ and ‘TP53’ genes are provided, the ‘pearson’ correlation method and the linear regression line is drawn, by default setting. You can choose different gene names and by mouse clicking. If you check the “Apply Log2” checkbox, the log2 function is applied to the normalized expression data for given gene symbols, and then it draws the correlation plot between the two genes (**Fig. 50**).

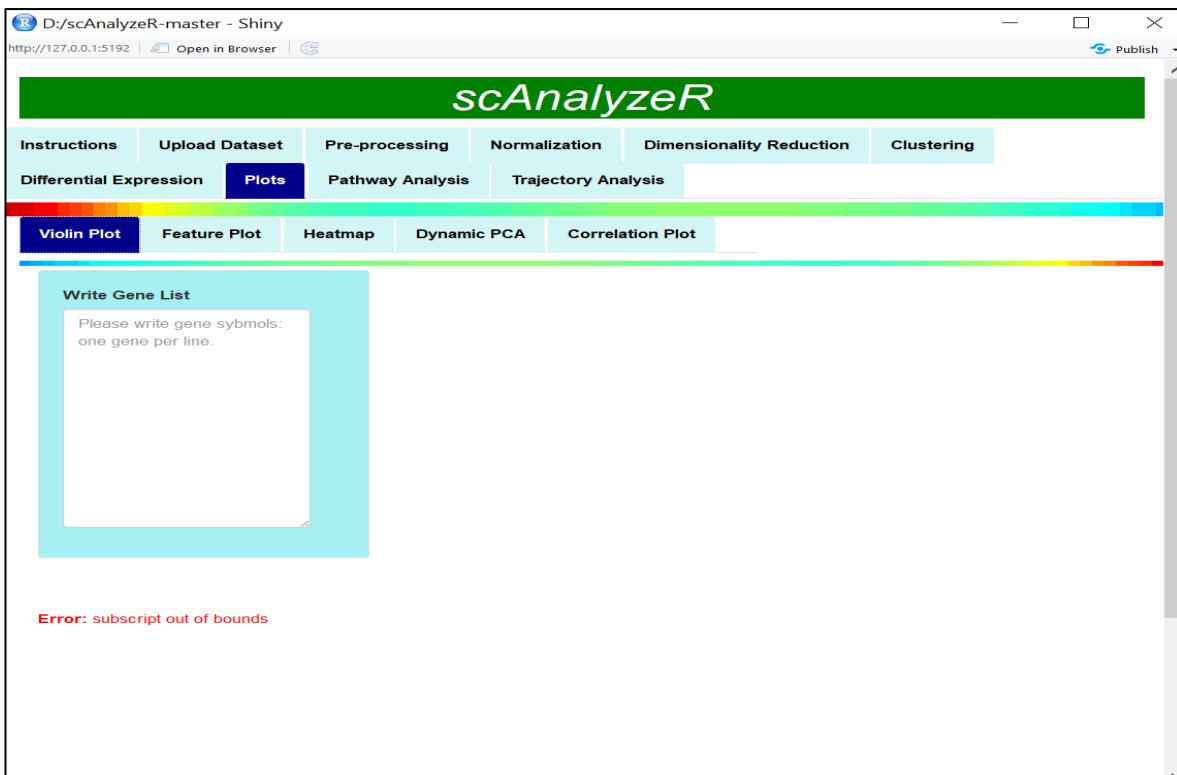


Figure 42. Violin plot: user interface

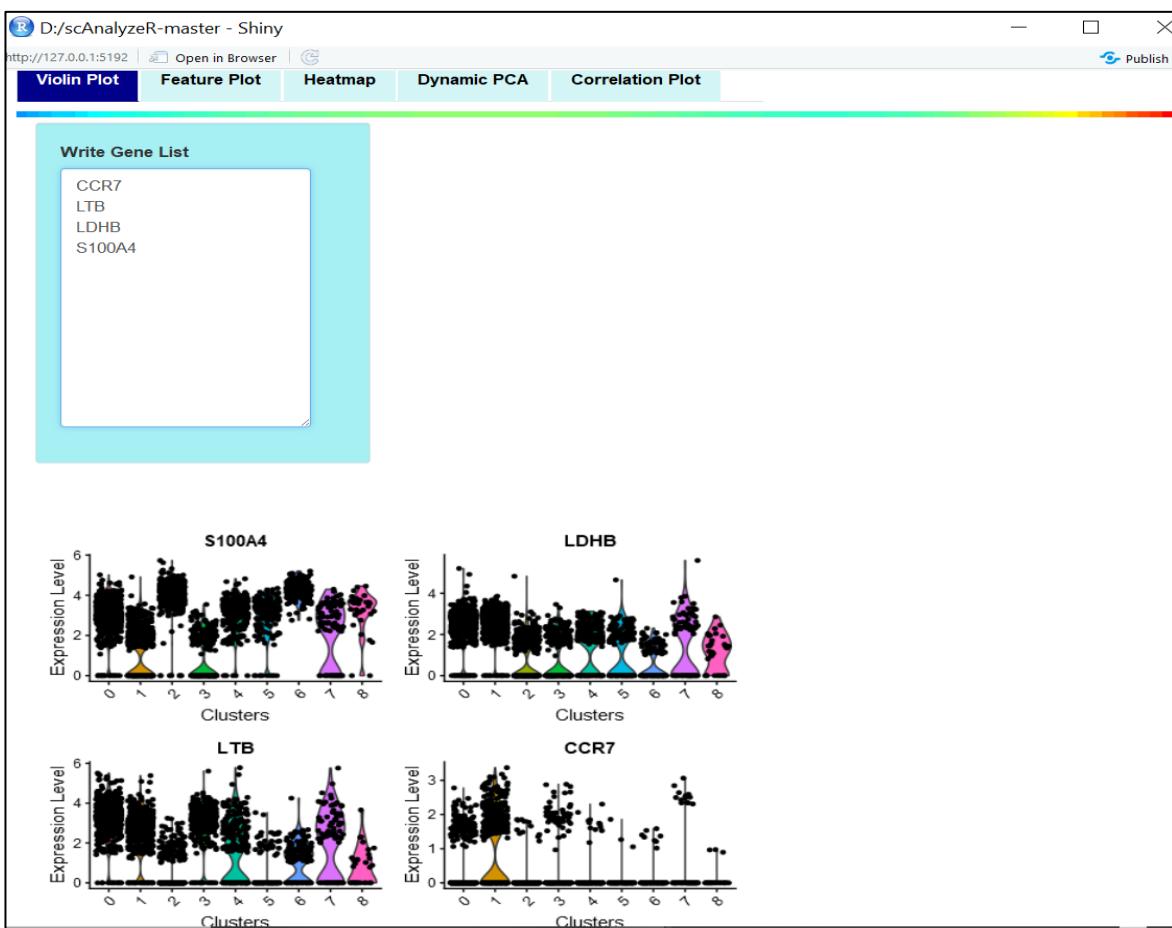


Figure 43. Violin plot for four individual genes

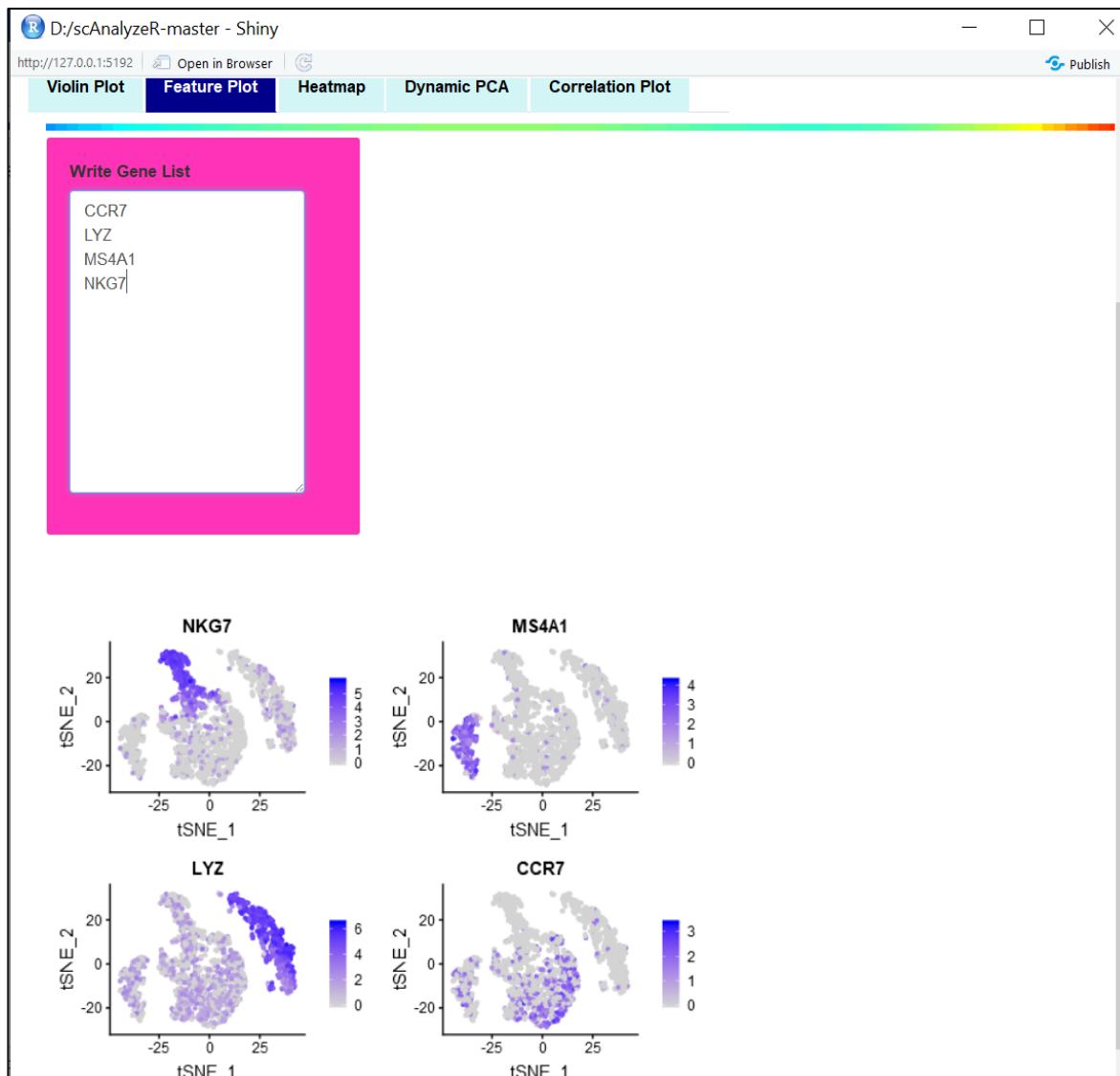


Figure 44. Feature plot for four individual genes

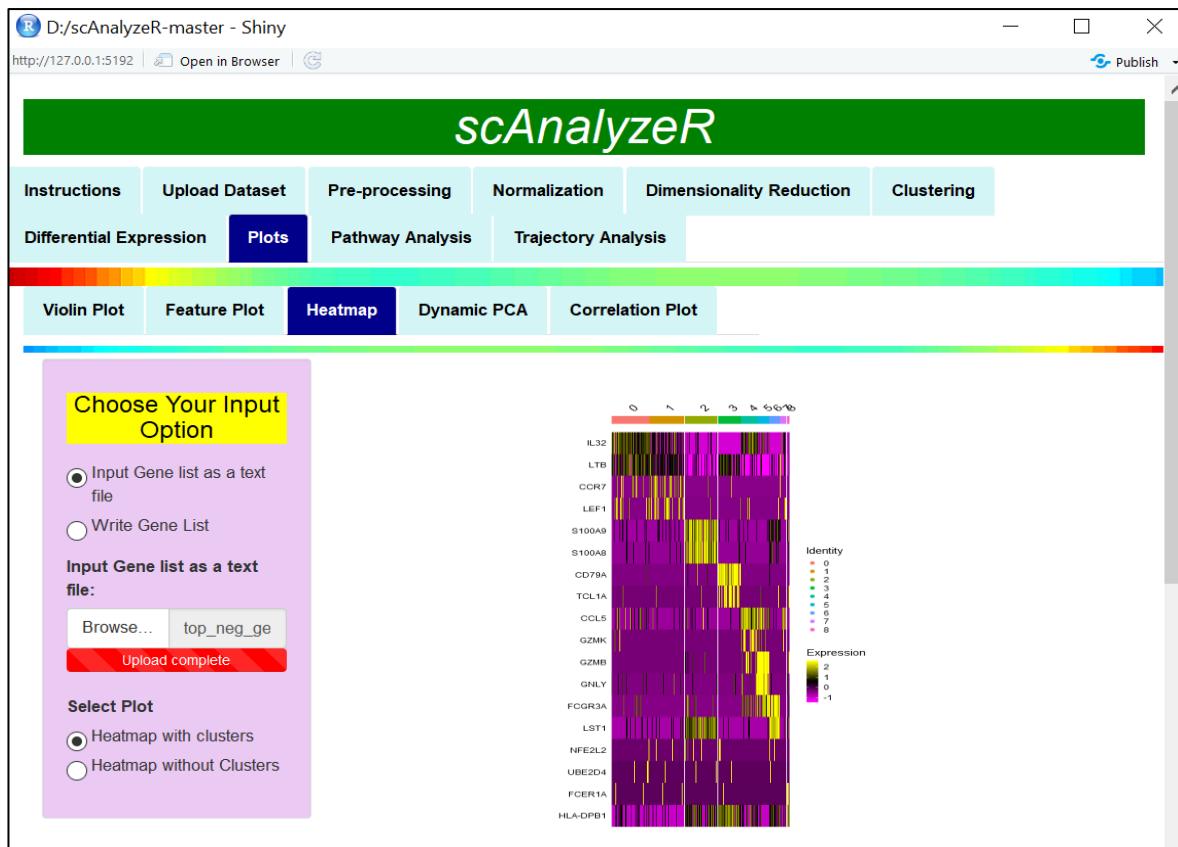


Figure 45. Heatmap plot for a set of genes(uploaded) with clusters indication

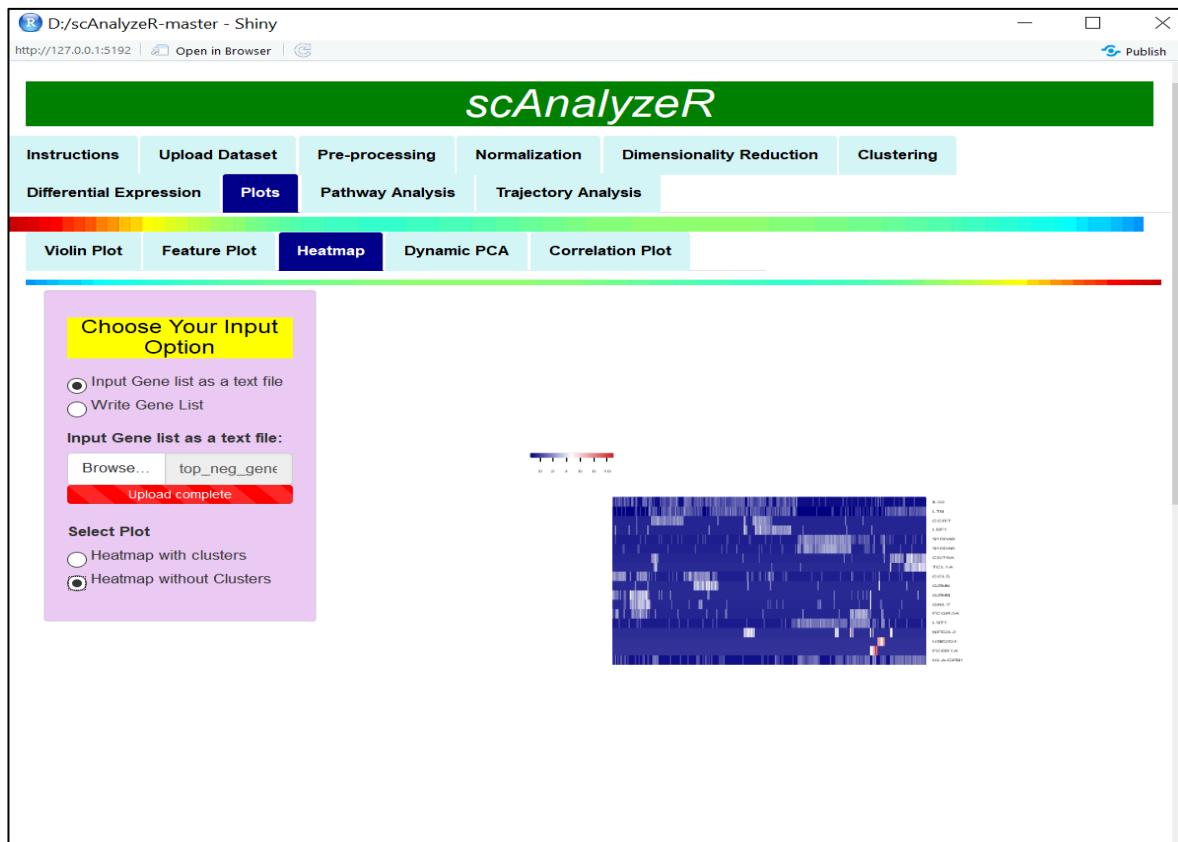


Figure 46. Heatmap plot for a set of genes(uploaded) without clusters indication

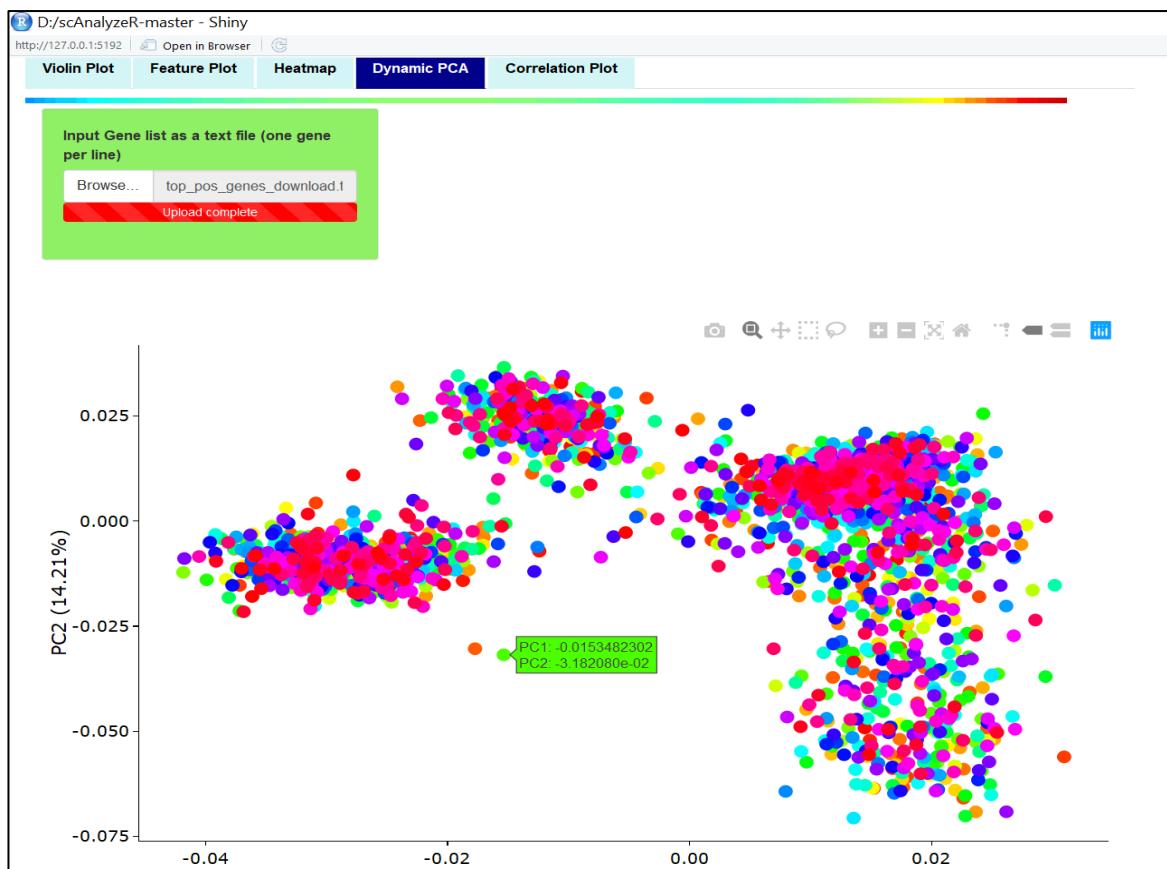


Figure 47. Dynamic PCA: PCA plot for a set of genes(uploaded) with dynamic customization options

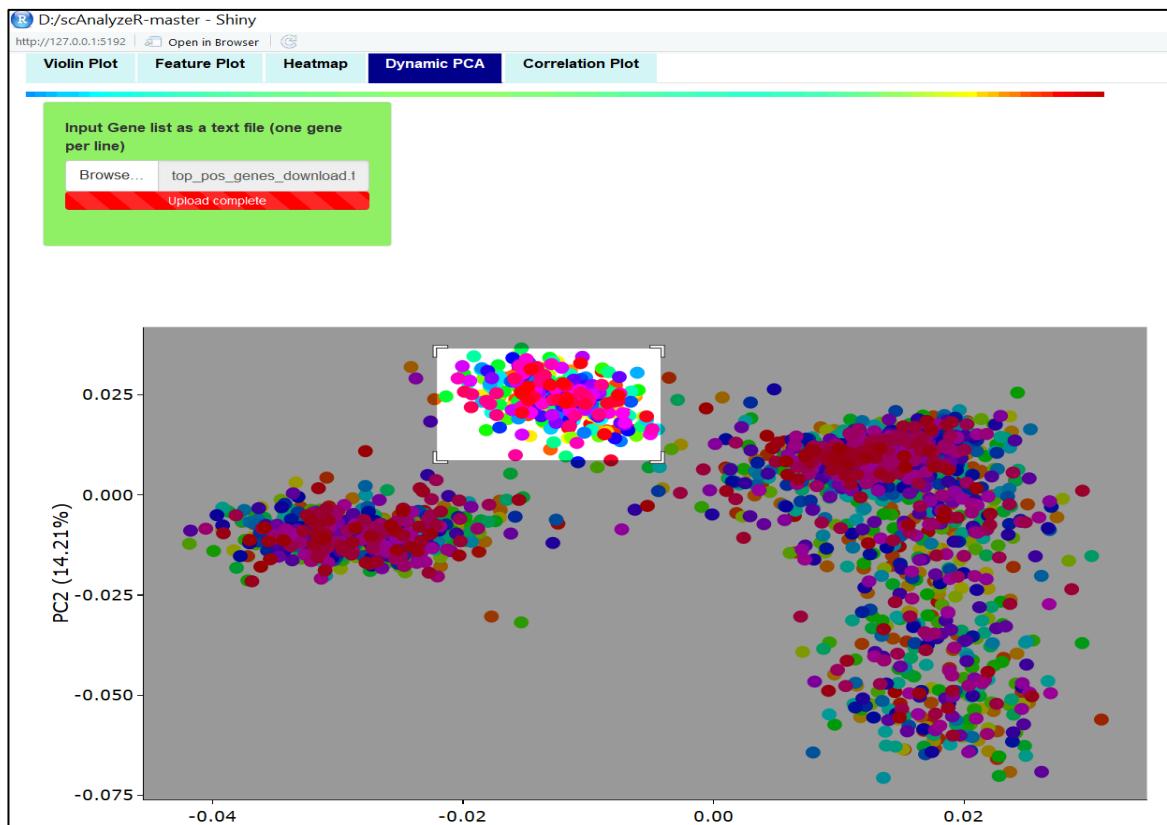


Figure 48. Dynamic PCA: Selected an area on the PCA plot for zooming cells(dots)



Figure 49. Dynamic PCA: PCA plot after zooming the selected area on Fig. 43

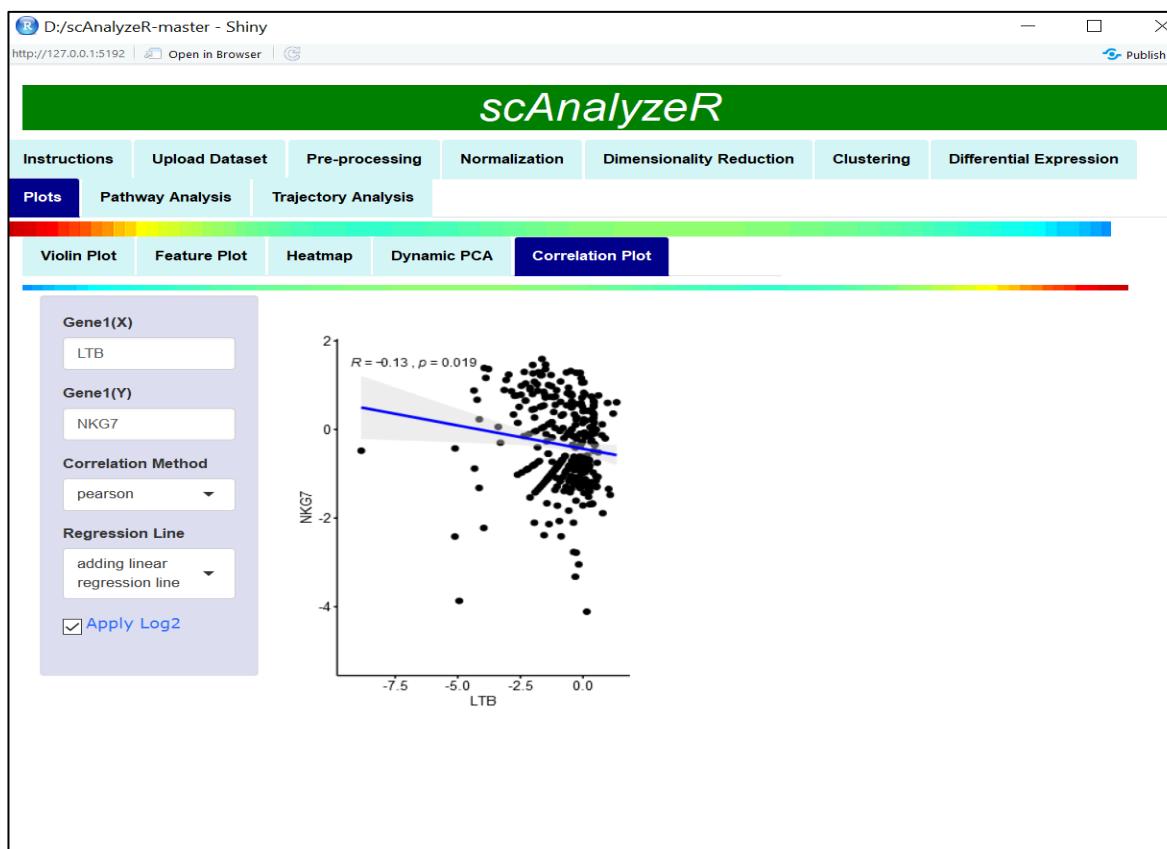


Figure 50. Correlation plot between two given genes

10. Pathway Analysis

The gene set enrichment analysis (GSEA) of previously (on ‘Differential Expression’ module) identified differentially expressed genes is performed in this analysis module. The “Cluster Specific DE Genes” is selected by default as an input gene list, and the ‘KEGG pathways’ is chosen as a source of pathway gene sets. You can either keep default selection or select different choices and click the “Compute Pathways” to proceed results. In our example, we choose default setting for computing pathways, i.e., ‘KEGG pathways’ and ‘Cluster Specific DE Genes’ (cluster 0’s DE genes) were selected (**Fig. 51**). In the “Select Pathways” tab, if you choose the ‘Positive Pathways,’ then only positive pathways are selected for all later parts of results generation, and the reverse action will be applied if the ‘Negative Pathways’ is chosen. However, the full list (both positive and negative) of pathways and filtered (filtered by padj) list of pathways can be downloaded via clicking the “Download All Pathways as a csv” and “Download Filtered Pathways as a csv” buttons, respectively.

To show the pathway plot and pathway list for filtered pathways, click the “Show Pathways(plot)” and the “Show Pathways list” checkboxes, respectively. Moreover, gene symbols of gene sets (either significant genes or all genes) and their heatmap for a particular pathway (1 by default) can be seen by clicking the “Show Gene sets” checkbox (**Fig. 52-53**).

D:/scAnalyzeR-master - Shiny
http://127.0.0.1:6670 | Open in Browser | G | Publish

scAnalyzeR

- [Instructions](#)
- [Upload Dataset](#)
- [Pre-processing](#)
- [Normalization](#)
- [Dimensionality Reduction](#)
- [Clustering](#)
- [Differential Expression](#)
- [Plots](#)
- Pathway Analysis
- [Trajectory Analysis](#)

Select DE Genes

Cluster Specific DE Genes
 Cluster(s) vs Cluster(s) DE Genes

Select Pathway Source

KEGG Pathways

Compute Pathways

[Download All Pathways as a csv](#)

Select Pathways

Positive Pathways
 Negative pathways

Cutoff: padj(<= cutoff)

0.5

[Download Filtered Pathways as a csv](#)

Show Pathways(plot)

Top Pathways:

10

Show 10 entries Search:

	pathway	pval	padj	ES	NES	nMoreExtreme	size	Sig_gen
1	Primary immunodeficiency	0.000353731871241599	0.00346685057128478	0.830462211181561	2.26965640324632	0	6	
2	Ribosome	0.00435594275046671	0.0317727588857572	0.556701030927835	2.10025202423163	6	12	
3	T cell receptor signaling pathway	0.0123393316195373	0.080530374780138	0.546198420400723	1.90089098036923	23	10	
4	Vasopressin-regulated water reabsorption	0.0261194029850746	0.161940298507463	0.990243902439024	1.32599387986836	132	1	
5	Cytokine-cytokine receptor interaction	0.0349900596421471	0.197216699801193	0.579545786454895	1.70034946185193	87	7	
6	Glycolysis / Gluconeogenesis	0.0659858601728201	0.303046172645544	0.970731707317073	1.29986592164435	335	1	
7	Cysteine and methionine metabolism	0.0659858601728201	0.303046172645544	0.970731707317073	1.29986592164435	335	1	
8	Pyruvate metabolism	0.0659858601728201	0.303046172645544	0.970731707317073	1.29986592164435	335	1	
9	Propanoate metabolism	0.0659858601728201	0.303046172645544	0.970731707317073	1.29986592164435	335	1	
10	Oocyte meiosis	0.0926118626430801	0.410138248847926	0.704433497536946	1.42240601562686	355	3	

Figure 51. Pathway Analysis: The table shows the top-ten up-regulated pathways after filtering

D:/scAnalyzeR-master - Shiny
http://127.0.0.1:6670 | Open in Browser | G | Publish

Compute Pathways

[Download All Pathways as a csv](#)

Select Pathways

Positive Pathways
 Negative pathways

Cutoff: padj(<= cutoff)

0.5

[Download Filtered Pathways as a csv](#)

Show Pathways(plot)

Top Pathways:

10

Show Pathways list

Show Gene sets

Pathway No.:

1

Select Gene Set

Significant Genes Only
 All Genes

Show Gene set Heatmap

Showing 1 to 10 of 15 entries

Previous 1 2 Next

```
[1] "IL7R"  "CD3D"  "CD3E"  "CD40LG" "LCR"   "IL2RG"
```

Figure 52. Pathway Analysis: Shows all gene list (bottom) in the top up-regulated pathway after filtering

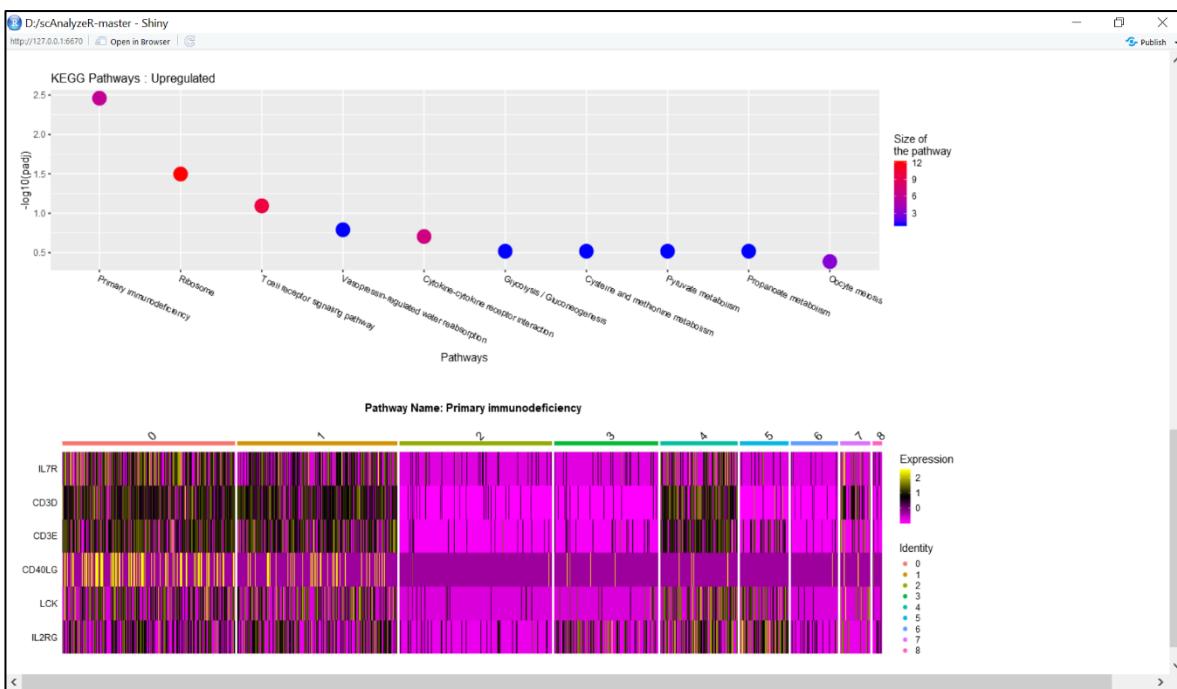


Figure 53. Pathway Analysis: Bubble plot shows the top-ten up-regulated pathways, and the heatmap illustrates expressions for the top up-regulated pathway's gene set (all genes)

11. Trajectory Analysis

In the last module, cells are ordered in pseudo-time according to their expression profiles such as cell differentiation stage. Select either ‘All genes’ or ‘Highly Variable Genes Only’ (default option) from the “Gene use” tab before clicking the “Cell trajectory plot,” then cell trajectory plot will be shown with three different plot styles, ‘Clusters’ is selected by default. To see the pseudotemporal heatmap, click the “Pseudotemporal heatmap” checkbox, then the interface of the input panel will be visible. The input gene list can either be uploaded as a text file (one gene per line) or written in the input box, and can change the number of clusters or keep the default value before clicking the “OK” button. Once calculating is completed, the heatmap will be seen on the interface. In this example, the highly variable genes were used for both cell trajectory plot and pseudo-temporal heatmap (top 10 highly variable genes only) analysis (**Fig. 54-56**).

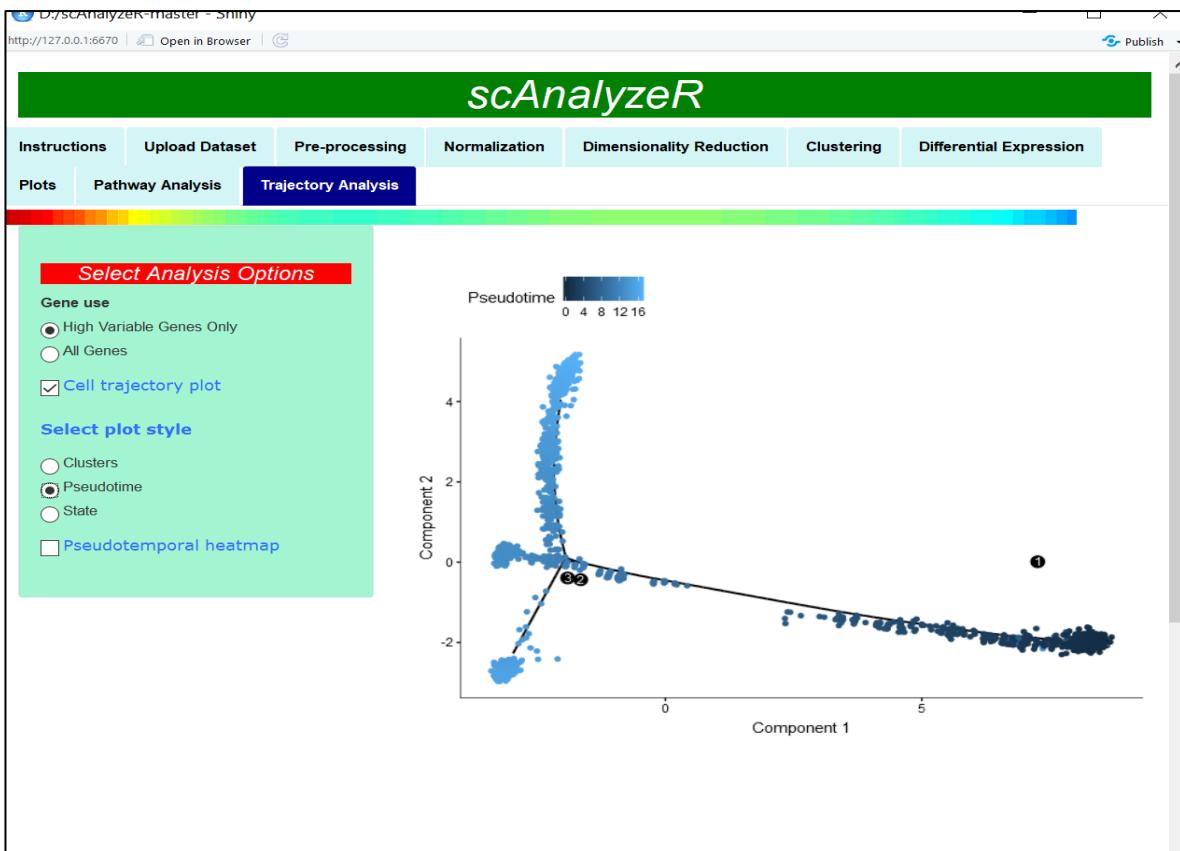


Figure 54. Trajectory Analysis: Cell trajectory plot with pseudotime

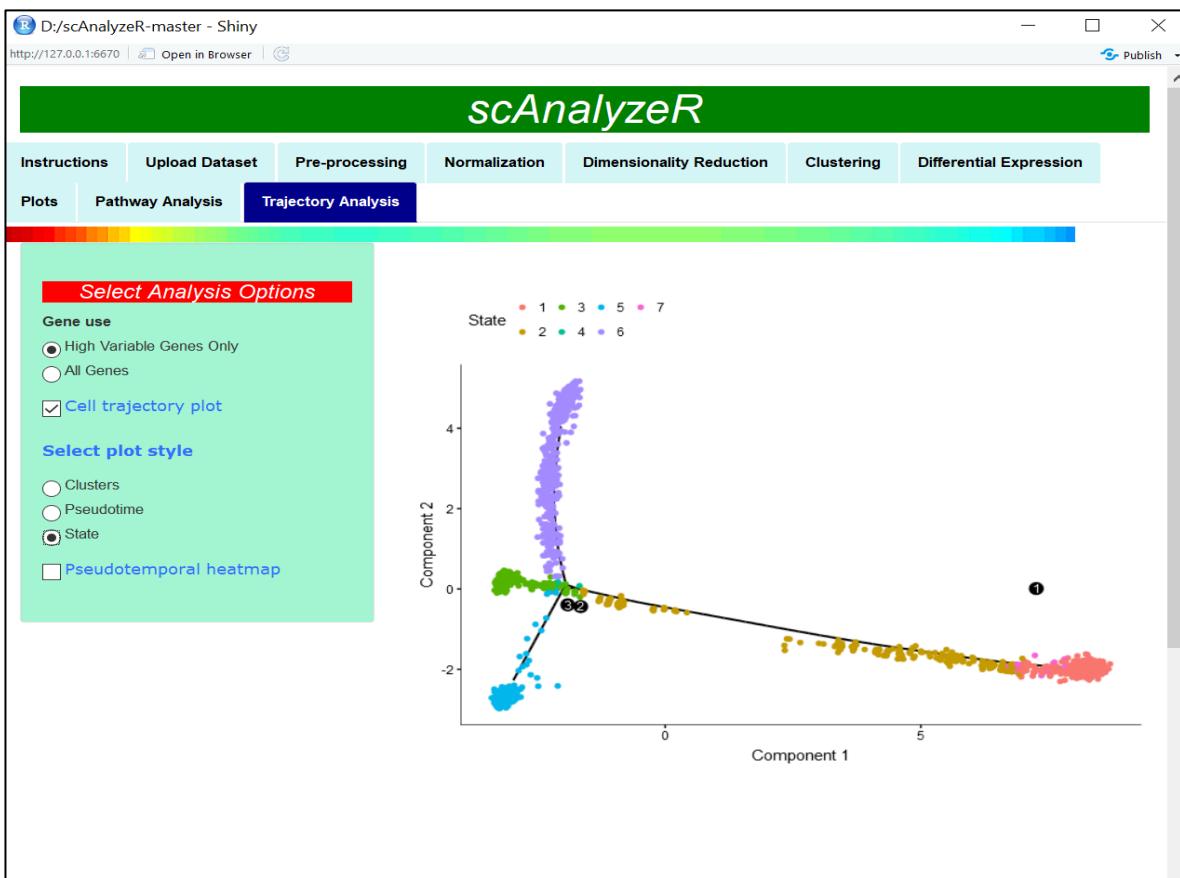


Figure 55. Trajectory Analysis: Cell trajectory plot with states

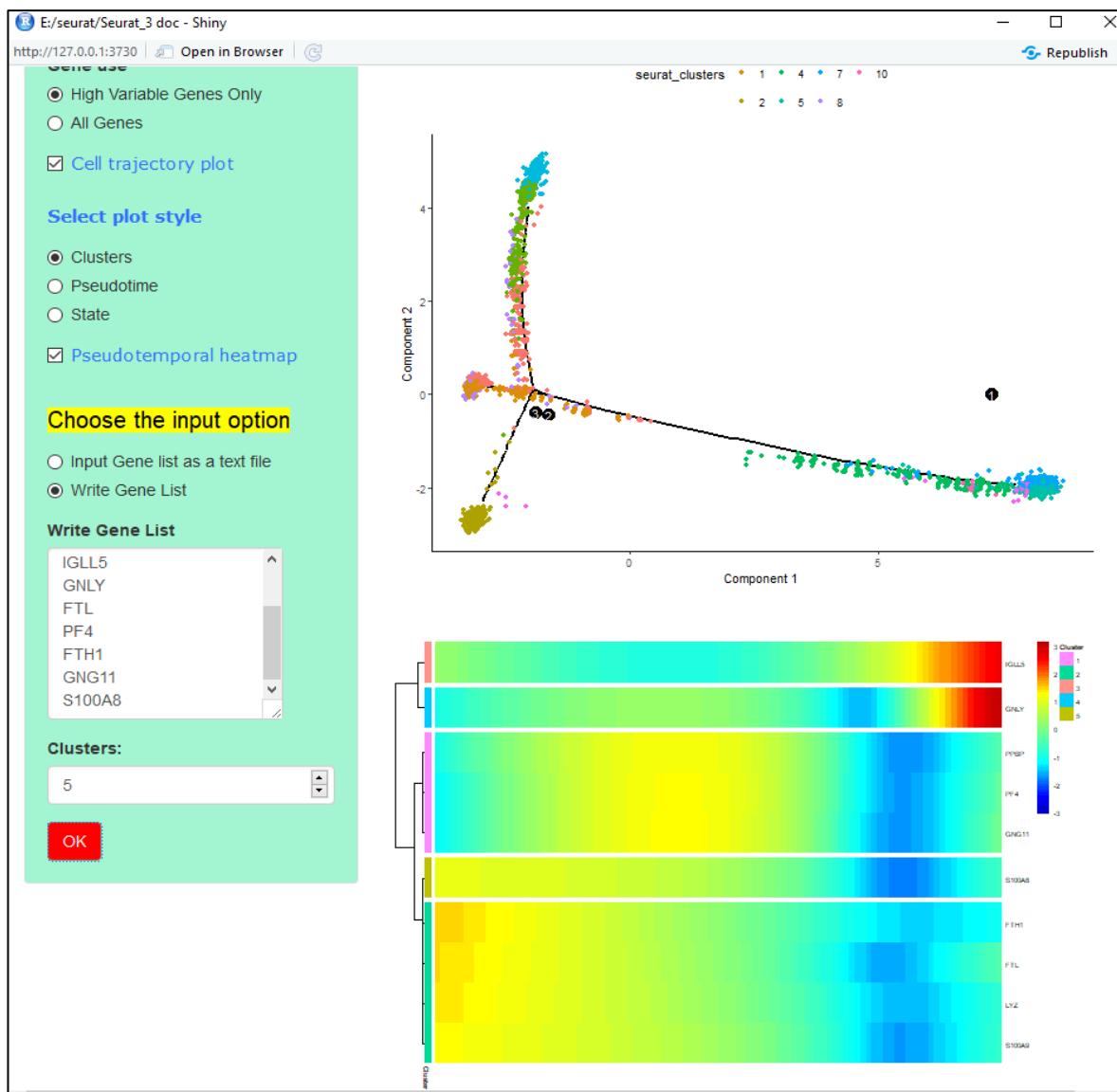


Figure 56. Trajectory Analysis: Cell trajectory plot with clusters (top), and pseudotemporal heatmap plot (bottom)

12. References

1. Nguyen, Q.H., et al., *Profiling human breast epithelial cells using single cell RNA sequencing identifies cell diversity*. Nat Commun, 2018. **9**(1): p. 2028.
2. Aizarani, N., et al., *A human liver cell atlas reveals heterogeneity and epithelial progenitors*. Nature, 2019. **572**(7768): p. 199-204.
3. Zheng, H., et al., *Single-cell analysis reveals cancer stem cell heterogeneity in hepatocellular carcinoma*. Hepatology, 2018. **68**(1): p. 127-140.
4. Su, X., et al., *Single-cell RNA-Seq analysis reveals dynamic trajectories during mouse liver development*. BMC Genomics, 2017. **18**(1): p. 946.
5. Rizvi, A.H., et al., *Single-cell topological RNA-seq analysis reveals insights into cellular differentiation and development*. Nat Biotechnol, 2017. **35**(6): p. 551-560.
6. Vuong, N.H., et al., *Single-cell RNA-sequencing reveals transcriptional dynamics of estrogen-induced dysplasia in the ovarian surface epithelium*. PLoS Genet, 2018. **14**(11): p. e1007788.