



Single-Cell Data Mining

Working notes for the single cell analytics course at HHMI

These notes include scOrange workflows and visualizations we will construct during the course.

The notes were written by Martin Stražar and Blaž Zupan with a huge help from the members of the Bioinformatics Lab in Ljubljana that develop and maintain scOrange. In part, we have reused lecture notes for Orange workshops as published by the same group.

Welcome! We have designed this course for biologists interested in interactive single-cell RNA sequencing data analysis. You will see how you can accomplish single-cell data mining tasks without programming. We will use scOrange to construct reproducible, shareable and visual data mining workflows.

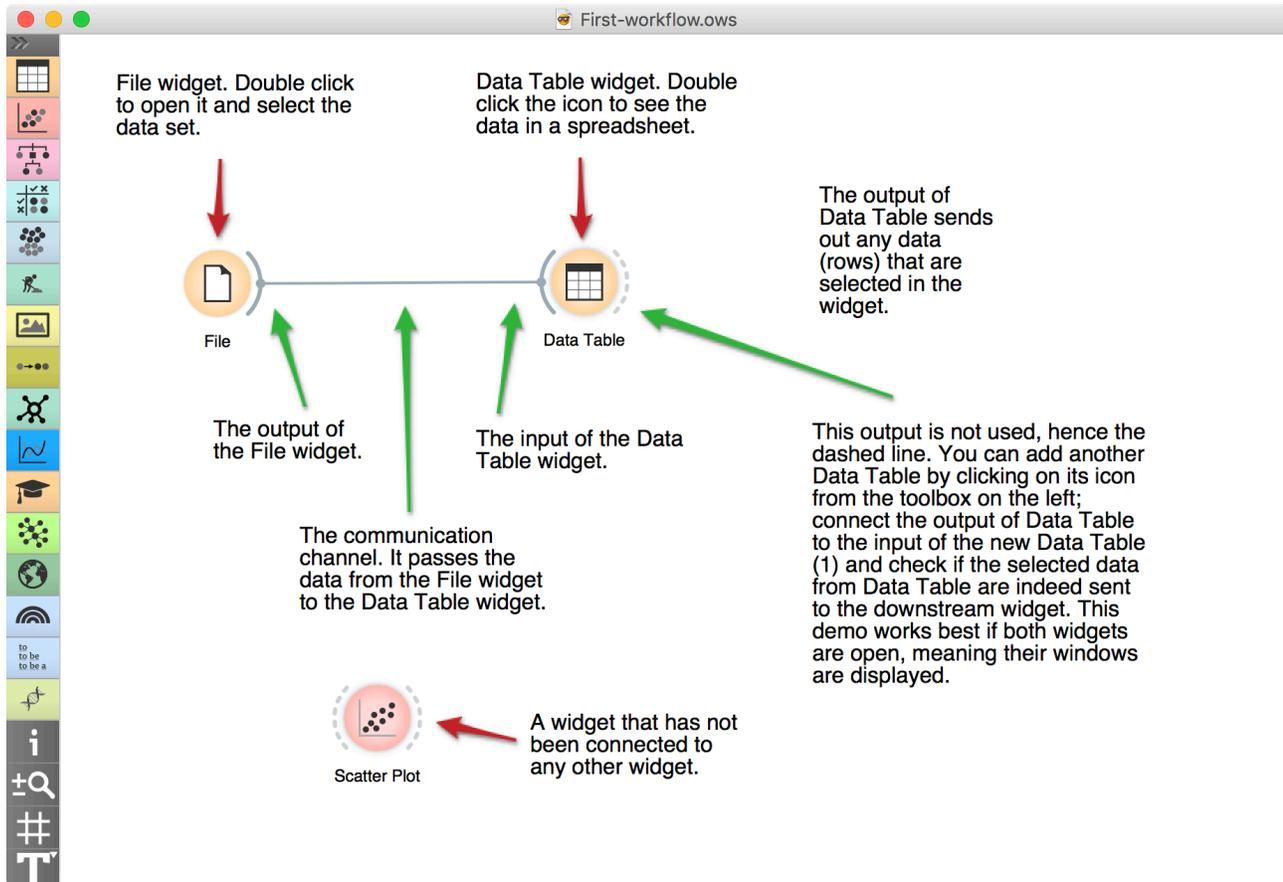
If you haven't already installed scOrange, please download the installation package from <https://singlecell.biolab.si>. The program scOrange is a derivative of Orange, an open-source data mining suite, with a pre-installed add-on for single-cell analytics. In this notes, we will refer to both of the programs as Orange, and refer to a program we will use during the workshop as scOrange when emphasizing that the functionality we describe is not included in the Orange distribution.



Attribution-NonCommercial-NoDerivs
CC BY-NC-ND

Lesson 1: Workflows in Orange

Orange workflows consist of components that read, process and visualize the data. We call them “widgets”. Widgets are placed on a drawing board (the “canvas”). Widgets communicate by sending information along a communication channel. Output from one widget is used as input to another.

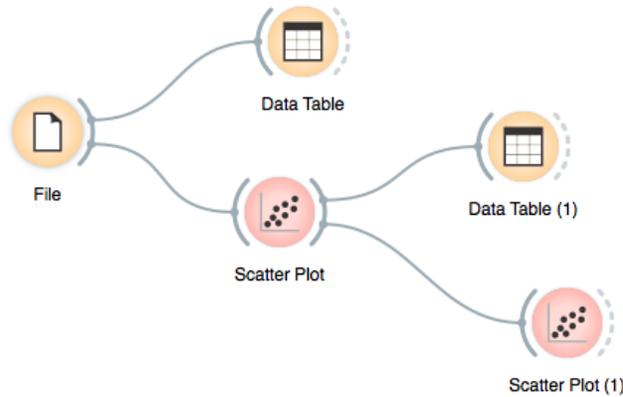


A simple workflow with two connected widgets and one widget without connections. The outputs of a widget appear on the right, while the inputs appear on the left.

We construct workflows by dragging widgets onto the canvas and connect them by drawing a line from the transmitting widget to the receiving widget. The widget’s outputs are on the right and the inputs on the left. In the workflow above, the File widget sends data to the Data Table widget.

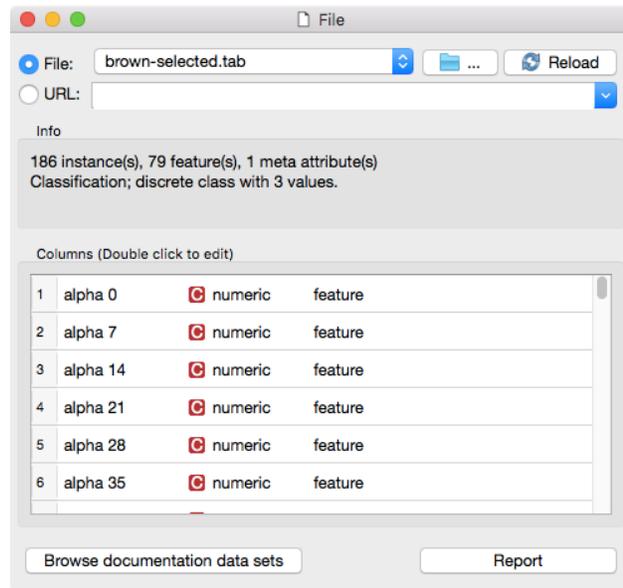
Start by constructing a workflow that consists of a *File* widget, two *Scatter Plot* widgets, and two *Data Table* widgets:

Workflow with a File widget that reads the data from disk and sends it to the Scatter Plot and the Data Table widgets. The Data Table renders the data in a spreadsheet, while the Scatter Plot visualizes it. Selected data points from the Scatterplot are sent to two other widgets: Data Table (1) and Scatter Plot (1).



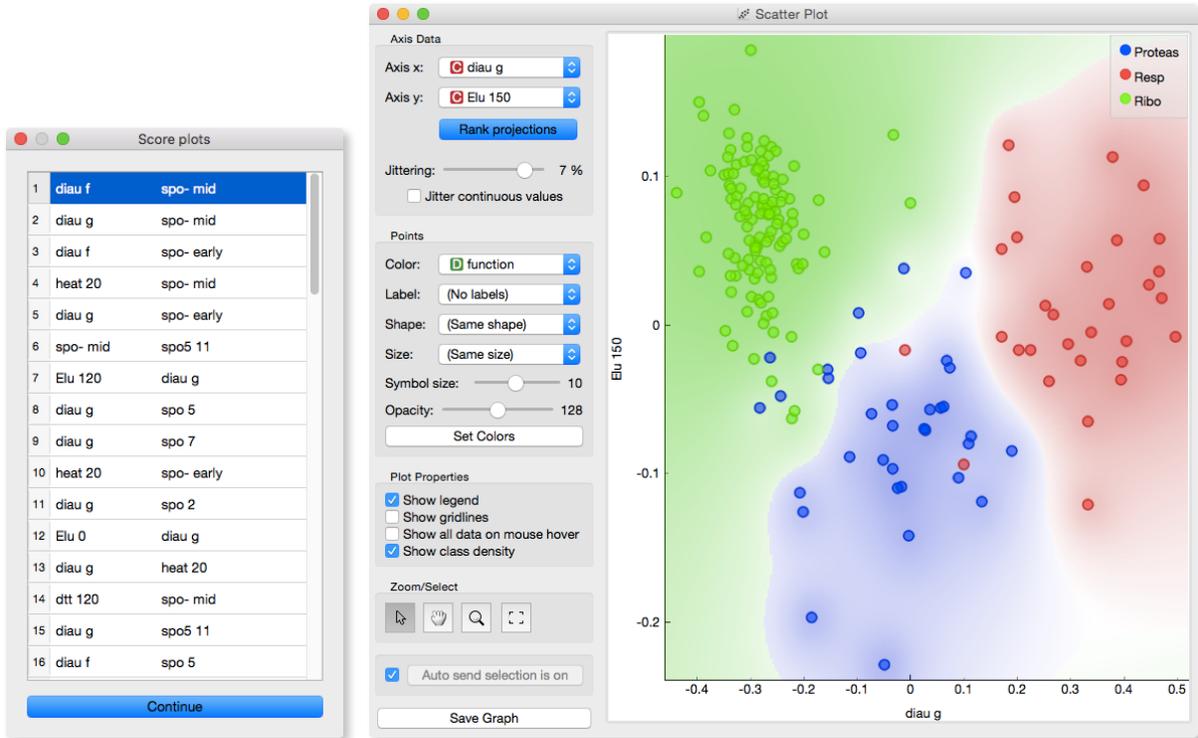
The File widget reads data from your local disk. Open the File widget by double-clicking its icon. Orange comes with several preloaded datasets. From these (“Browse documentation data sets...”), choose *brown-selected.tab*, a yeast gene expression dataset.

Orange workflows often start with a File widget. The brown-selected dataset comprises of 186 rows (genes) and 81 columns. Out of the 81 columns, 79 contain gene expressions of baker’s yeast under various conditions, one column (marked as a “meta attribute”) provides gene names, and one column contains the “class” value or gene function.



After you load the data, open the other widgets. In the Scatter Plot widget, select a few data points and watch as they appear in the widget Data Table (1). Use a combination of two Scatter Plot widgets, where the second scatter plot shows a detail from a smaller region selected in the first scatter plot.

Following is more of a side note, but it won't hurt. Namely, the scatter plot for a pair of random features does not provide much information on gene function. Does this change with a different choice of feature pairs in the visualization? Rank projections (the button on the top left of the Scatter Plot widget) can help you find a good feature pair. How do you think this works? Could the suggested pairs of features be useful to a biologist?



Lesson 2: Loading Your Own Dataset

Orange reads data from Excel, comma- and tab-separated files and urls. Try constructing a spreadsheet in Excel or in Google Sheets. If using Google Sheets, copy a shareable link and paste it into the File widget. Press Enter to load the data.

The datasets we have worked with in the previous lesson come with Orange installation. Orange can read the data from spreadsheet file formats which include tab and comma separated and Excel files. Let us prepare a toy dataset in Excel and save it on a local disk.

	A	B	C	D	E	F	G	H	I
1	Cell ID	Type	BCL2	CCR5	CD4	CD8	IL2	IL10	
2	TGATGATT	damaged	6	31	76	52	8	96	
3	TTAAGGCC	normal	24	16	20	19	78	50	
4	CCGGAATT	damaged	72	47	20	67	9	24	
5	CGCGAGTG	damaged	71	25	11	38	65	2	
6	GTAGCGTG	normal	0	90	46	40	69	21	
7	CGTAGCTG	normal	9	6	95	36	46	13	

In Orange, we can use the File widget to load this dataset.

Name	Type	Role	Values
1 Type	C categorical	feature	damaged, normal
2 BCL2	N numeric	feature	
3 CCR5	N numeric	feature	
4 CD4	N numeric	feature	
5 CD8	N numeric	feature	
6 IL2	N numeric	feature	
7 IL10	N numeric	feature	
8 Cell ID	S text	meta	

Looks good. Orange has correctly guessed that cell IDs are character strings and that this column in the dataset is special, meant to provide additional information and not to be used for any

kind of modeling. All other columns are numeric features except for the type, which is a categorical feature. This is also the feature we would not like to include in the profile of the cell and should rather consider it as a cell's class. Double-click on the “feature” in the *Role* column and change the role of the feature type to “target”. Then click the *Apply* button.

It is always good to check if all the data was read correctly. We can connect our *File* widget with the *Data Table* widget,



and double-click on the *Data Table* to see the data in the spreadsheet format.

The screenshot shows the Data Table widget interface. On the left, there is an 'Info' panel with the following details:

- 6 instances (no missing values)
- 6 features (no missing values)
- Discrete class with 2 values (no missing values)
- 1 meta attribute (no missing values)

The 'Variables' section has the following options:

- Show variable labels (if present)
- Visualize numeric values
- Color by instance classes

The 'Selection' section has the following options:

- Select full rows

Buttons include 'Restore Original Order' and 'Send Automatically' (checked).

The main table displays the following data:

	Type	Cell ID	BCL2	CCR5	CD4	CD8	IL2	IL10
1	damaged	TGATGATT	6.0	31.0	76.0	52.0	8.0	96.0
2	normal	TTAAGGCC	24.0	16.0	20.0	19.0	78.0	50.0
3	damaged	CCGGAATT	72.0	47.0	20.0	67.0	9.0	24.0
4	damaged	CGCGAGTG	71.0	25.0	11.0	38.0	65.0	2.0
5	normal	GTAGCGTG	0.0	90.0	46.0	40.0	69.0	21.0
6	normal	CGTAGCTG	9.0	6.0	95.0	36.0	46.0	13.0

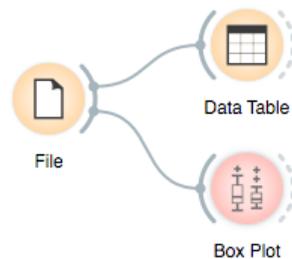
Instead of using Excel, we could also use Google Sheets, a free online spreadsheet alternative. Then, instead of finding the file on the local disk, we would enter its URL address to the File widget's URL entry box.

Nice, everything is here.

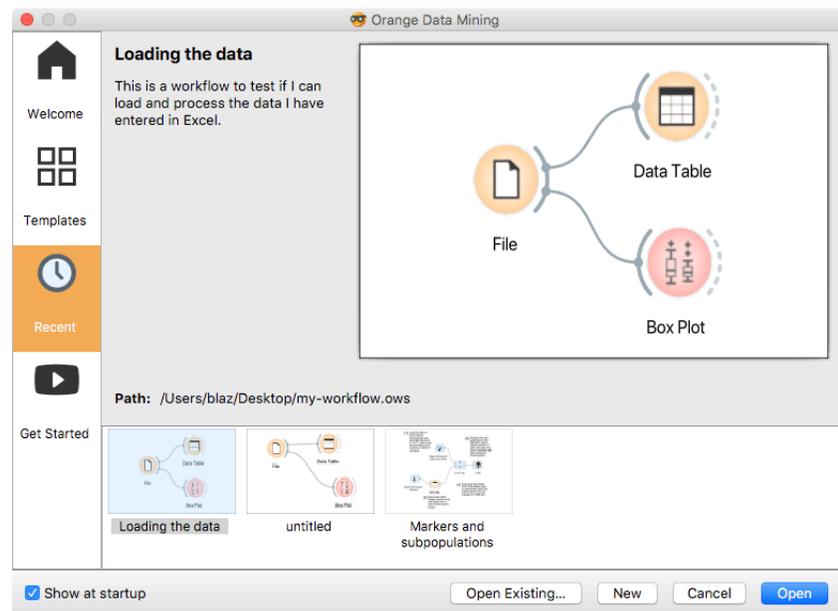
There is more to input data formatting and loading. We can define the type and kind of the data column, specify that the column is actually a web address of an image, and more. But enough for now. If you would really like to dive in for more, check out the [documentation page on Loading your Data](#), or a [video](#) on this subject.

Lesson 3: Saving and Sharing Your Work

Add a box plot to your workflow from the previous lesson. Your workflow should now look like this:



Orange can save your workflow to a file. Use Save or Save As from the File menu and save your work to your desktop. Before saving, you can describe your workflow on the information page using “i” icon for the toolbar. We have named the file as *my-workflow* and Orange would add a suffix to it. Now exit Orange and rerun it. Select *Recent* in the welcome screen and choose from one of the recently saved workflows.



You can email workflow files or share them with your colleagues. Note though that the data, unless stored on the web, has to be sent separately.

Lesson 4: Basic Exploration of Single Cell Data

Use the *Single Cell Datasets* widget and load the *Cell cycle in mESC (Fluidigm)* dataset with 182 cells and a full set of 38,293 genes.

Single-Cell Orange (scOrange) provides a quick-start interface to several recent datasets from single-cell RNA publications. In this lesson, we will look at gene expression changes in the cell cycle stages of mouse embryonic stem cells.

Single Cell Datasets

Info
18 datasets
5 datasets cached

Filter
cell cycle

Title	Size	Instances
Cell cycle in mESC (Fluidigm)	3.7 MB	182
Cell cycle in T-cells	4.3 MB	81
Cell cycle in T-cells (cell cycle genes)	220.2 KB	81
Cell cycle in mouse liver	125.9 KB	5
Cell cycle in mouse liver (cell cycle genes)	3.5 KB	5
Cell cycle in mESC (bulk RNA-seq)	287.0 KB	4
Cell cycle in mESC (bulk RNA-seq, cell cycle genes)	7.0 KB	4

Description
Cell cycle in mESC (Fluidigm) (2015), from [ArrayExpress](#)
A single-cell RNA-seq dataset comprised of 182 mouse embryonic stem cells (mESCs) with known cell-cycle phase. Cells were sorted using FACS for three different cell-cycle phases. This resulted in a filtered set of 59 cells in G1 phase. 58 cells in S phase

Send Data

You can get a first glance of the data — as before — with the *Data Table* widget.



Single cell data can use different quantification measures, like read counts or transcripts counts. Our cell cycle data includes read counts.

Data Table

Info
182 instances (no missing values)
38293 features (no missing values)
Discrete class with 3 values (no missing values)
No meta attributes

Variables
 Show variable labels (if present)
 Visualize numeric values
 Color by instance classes

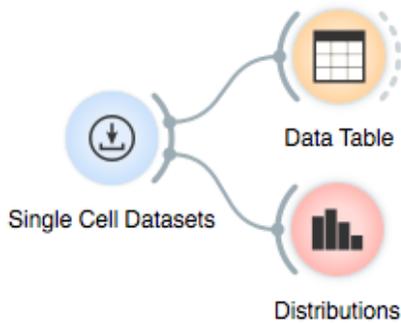
Selection
 Select full rows

	cycleStageName	Gna13	Pbsn	Cdc45	H19	Scm2
1	G1	82	0	598	0	0
2	G1	5	0	10	0	0
3	G1	87	0	36	0	391
4	G1	49	0	72	0	3
5	G1	855	0	6	0	0
6	G1	137	0	374	0	0
7	G1	733	0	782	0	32
8	G1	1	0	703	0	1
9	G1	715	0	359	0	1
10	G1	643	0	595	0	0
11	G1	1083	0	12	0	0
12	G1	5	0	1242	0	0
13	G1	1250	0	21	0	0
14	G1	38	0	804	0	0
15	G1	114	0	962	0	0
16	G1	524	0	4	0	31
17	G1	177	0	761	0	7
18	G1	6	0	115	0	12
19	G1	0	0	242	0	5
20	G1	36	0	866	0	333
21	G1	169	0	5	0	0
22	G1	42	0	1414	0	0
23	G1	0	0	50	0	0
24	G1	517	0	3	0	22
25	G1	0	0	39	0	0
26	G1	200	0	978	0	239
27	G1	1	0	886	0	23
28	G1	0	0	256	0	50

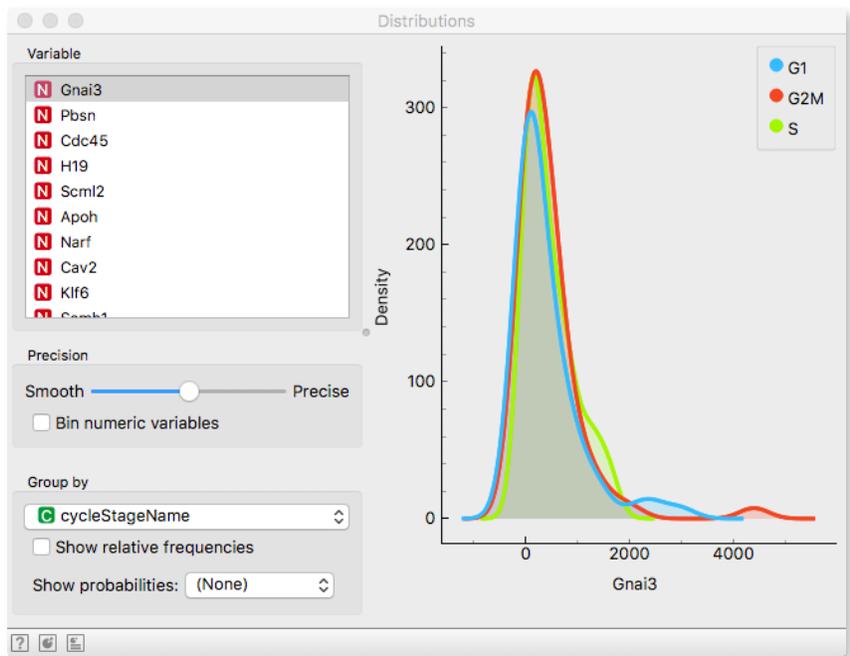
Restore Original Order
 Send Selected Rows

The program scOrange stores the data tables with cells in rows and genes in columns. For each cell, we have the count of sequence reads associated with each of the vast numbers of genes. In other words, we have a direct quantification of gene expression profiles. There are a few things to notice here — what can we say about the range of counts stored in the data table?

The counts vary wildly across cells, but also across genes — this is a typical example of technical variance that is compensated with downstream statistical methods. We can visualize the distribution of counts with the *Distribution* widget. Select different genes and see how genes generally either display very heavy-tailed distributions of counts or are not detected at all.



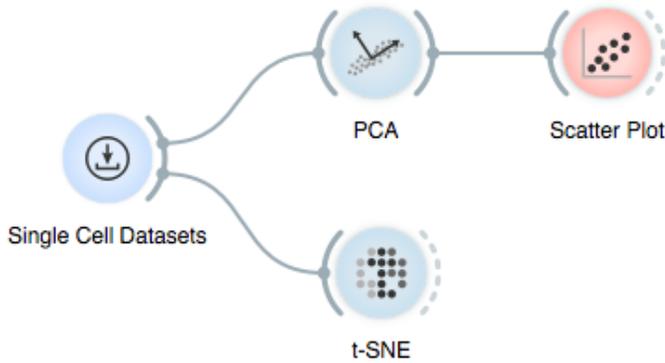
Set the *Group By* property to *cycleStageName* to see the distribution of counts for each of the three subgroups of cells.



Heavy-tailed and sparse data does not play well with conventional data analysis techniques. In the next section, we will do something about this.

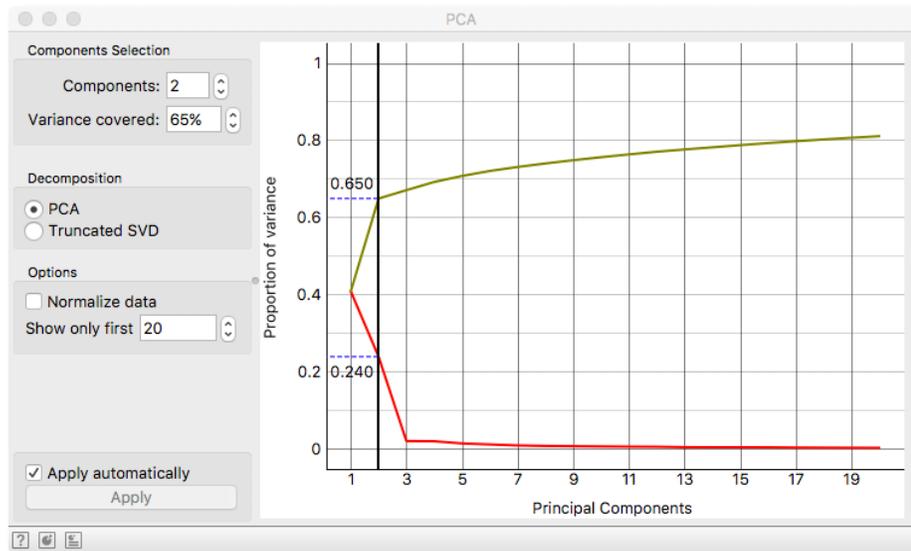
But first, a short detour.

Lesson 5: Data Projection with PCA and t-SNE



Looking at the values in a large table will probably not get us very far — except if you're a fan of a certain spreadsheet program. We can instead start with the *Principal component analysis (PCA)*, which is a way to quantify variance of different directions in multi-dimensional data. Selecting directions with the highest explained variance is one way to reduce the dataset size. If we consider only two directions — the two principal components — we can visualize the data in a scatter plot.

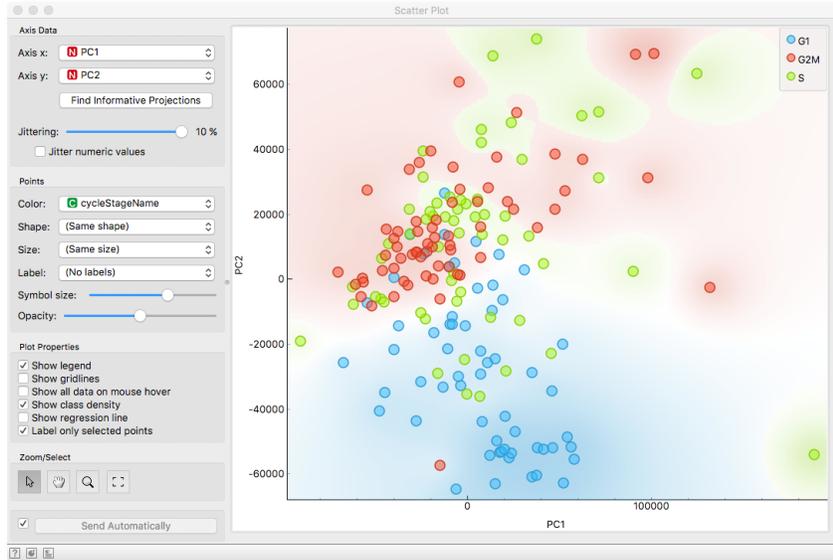
Note that 'Normalized data' checkbox is ticked off in the PCA; simply standardizing the gene expressions would amount to some loss of information on relative expression levels among genes.



The first two components explain a gigantic 65% of the variance, which is great, considering that we are dealing with 182-dimensional data!

However, the large magnitude on the axes reminds us of the very heavy-tailed distributions and raises an alarm that the projection might be dominated by a small number of high count values. We are postponing this problem to the next section.

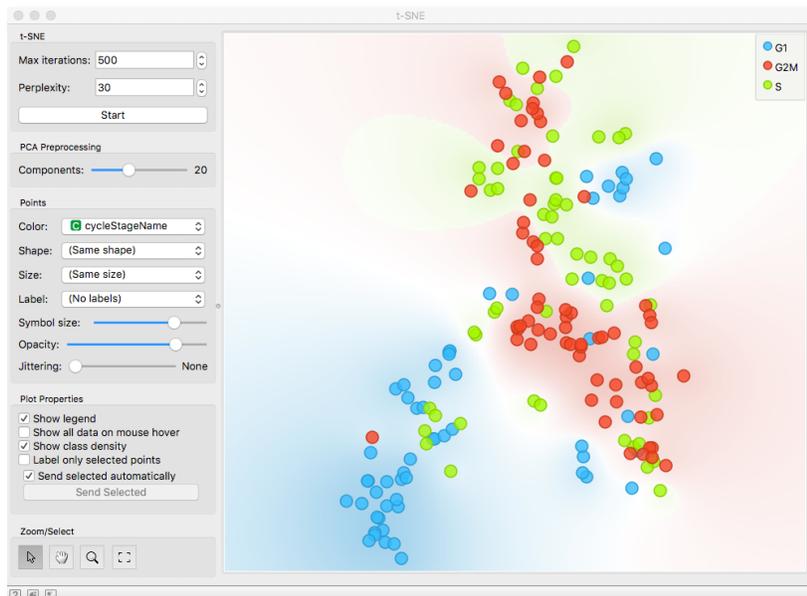
A word of caution: t-SNE essentially discards the true distances between the cells. Could the 2-D t-SNE coordinates be used in further downstream analyses, e.g. finding clusters?



The PCA + Scatter plot combo shows a mixed result: there is substantial overlap between cells in a different cell-cycle stage.

While PCA is a linear projection, the *t-distributed Stochastic Neighbour Embedding (t-SNE)* is a technique that is more suited to how humans interpret 2-D projections. With t-SNE, the main goal is to project very similar cells close together, while the distances between dissimilar cells do not significantly influence the positioning. This relaxation prevents spherical structures as those commonly seen in PCA, where the distances between points are treated equally. Hence t-SNE can yield more clustered, grape-like structures.

Try to modify the *Perplexity* parameter of the t-SNE plot. How does this value influence the visualization? What kind of values are more suitable to detect local structure and what for global structure?



Lesson 6: Data Filtering

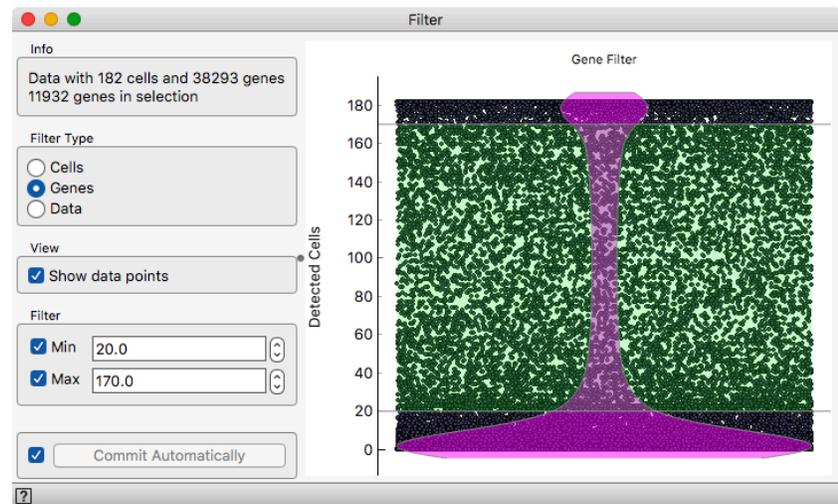
As we have previously seen, single cell data is an expression matrix with tens of thousands of genes and hundreds or thousands of cells. A natural question to ask is whether all of the data is required for further analysis.

Using the *Filter* widget, we can retain only the genes that are expressed in a certain minimal or maximal number of cells. For this example, let's set the minimal/maximal threshold to 20/170 cells (out of 182), respectively. From a total of 38,293 genes, we retain 11,932 genes. In other words, we can reduce the dataset size by 75%!



By setting a lower bound on the number of cells per gene, we remove the genes that are unlikely to significantly influence the analysis. By removing the very highly expressed genes, the *housekeeping* (ever-present) genes are removed.

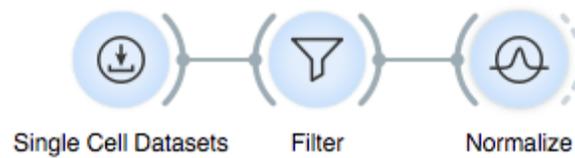
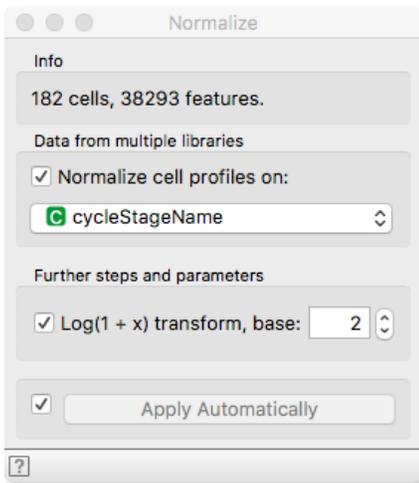
A similar filtering can be performed for cells by setting the lower/upper bound on the number of detected genes.



Lesson 7: Normalization

We can also tackle the technical variability issues. Due to a non-uniform distribution of reads per cell, some cells will *a priori* have larger counts associated to genes. A simple form of normalization is achieved when the gene count profiles are scaled to have the *same median value for all cells*. Alternatively, the *same median* constraint can be imposed on groups of cells which in this case can be defined as cells in the same cell cycle phase. We perform normalization after the data filtering.

You can also try to use *PCA* widget in the last step. How do the coordinate axes change after preprocessing?

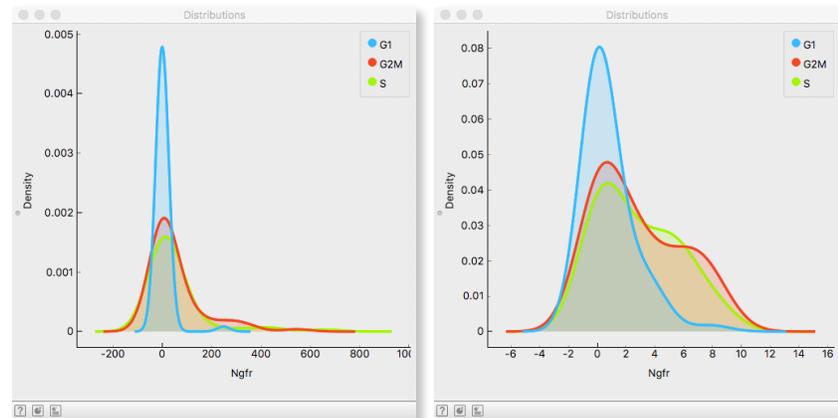


By scaling the cell expression profiles by a factor does not get rid of the heavy-tailed gene expression distributions. To achieve this, we use the following transformation:

$$y = \log(x + 1)$$

Why does this transformation make sense? What happens to genes with an expression equal to zero in the raw and transformed data? Try using *Distribution* widget to obtain the answer.

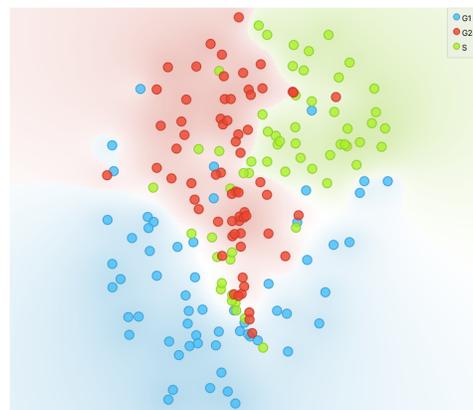
The distributions show the expression of gene *Ngfr* before and after normalization. Notice how we get rid of the long tail!



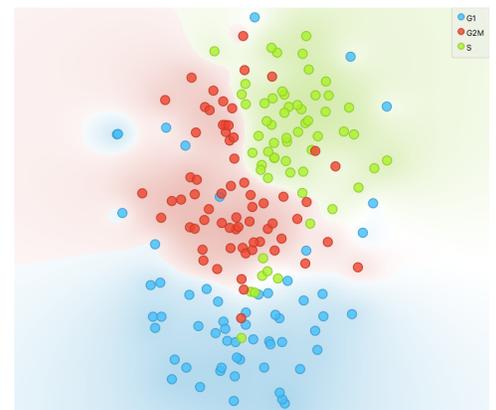
To get a feeling of what the filtering and normalization do to the data, let's create a sequence of t-SNE plots after applying each of the preprocessing steps. The preprocessing does quite a nice job, focusing on the genes that might be important for separating cells in different phases of the cell cycle!



Raw data



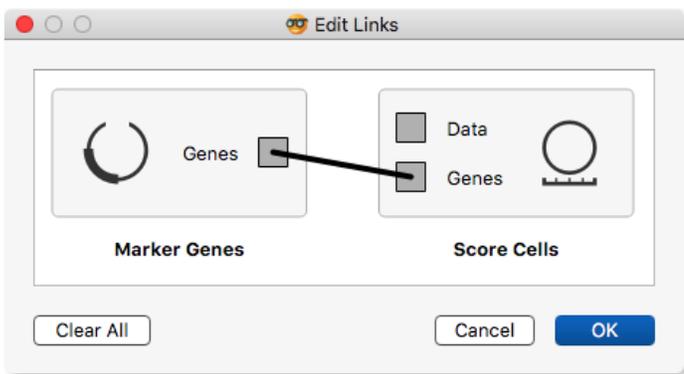
Filtered



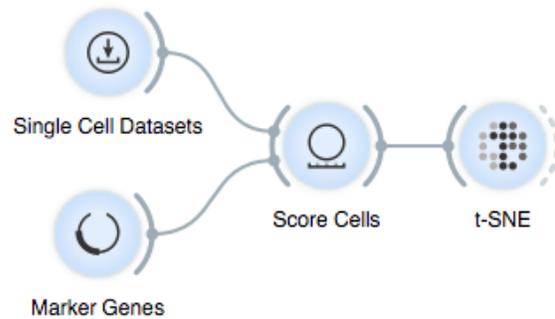
Filtered + normalized

Lesson 8: Marker Genes

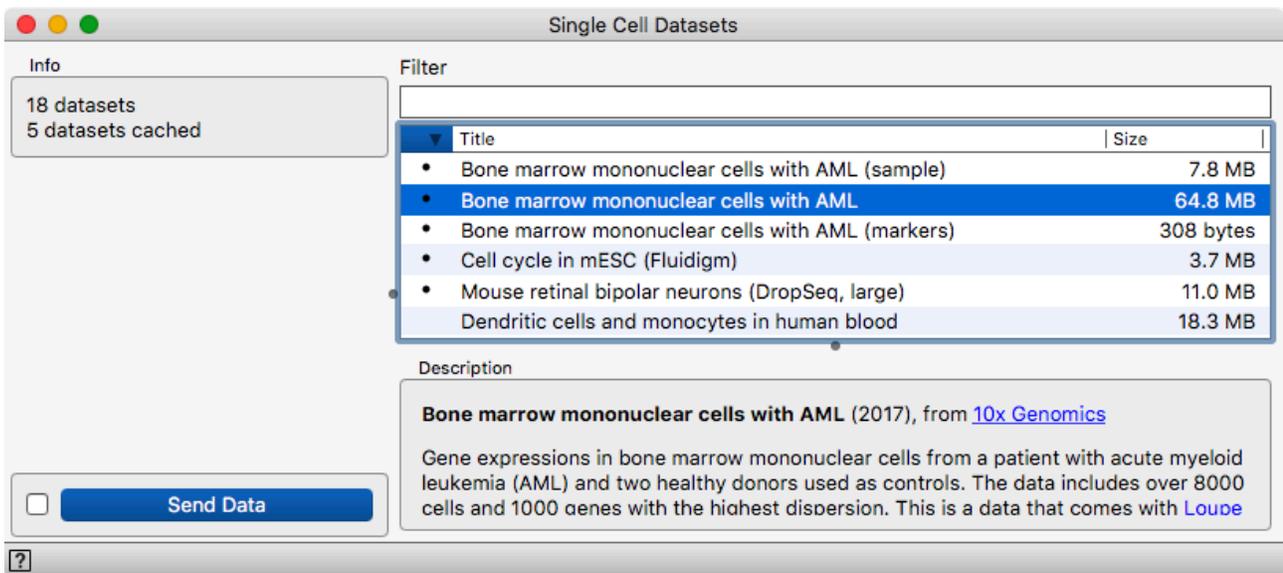
Score Cells widget has two inputs: the data and a set of marker genes. Since both inputs are data tables, how does it know which one is which? Double-click on any of the edges connecting the widgets to see the actual wiring. You can rewire the connections there.



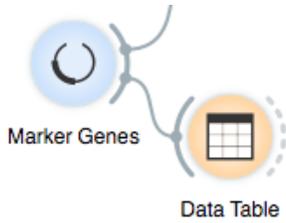
Gene markers are one of the common methods for discovering cell populations, determining cell states (e.g. cell cycle) and more. Here, we will use a sample of the data on bone marrow mononuclear cells with AML (Zheng et al., Nat Comm 2017), score the cells according to selected gene markers, and observe the scored cells in the t-SNE plot:



The sample bone marrow dataset from *Single Cell Datasets* widget includes 1000 genes and 1000 cells. There is a bigger dataset available from the same widget (8000 cells), but let us continue with the sample and define our workflow first. Then, you are free to choose larger datasets.



A widget *Marker Genes* offers a list of markers and outputs a data table with selected rows. You can test this by connecting it to *Data*



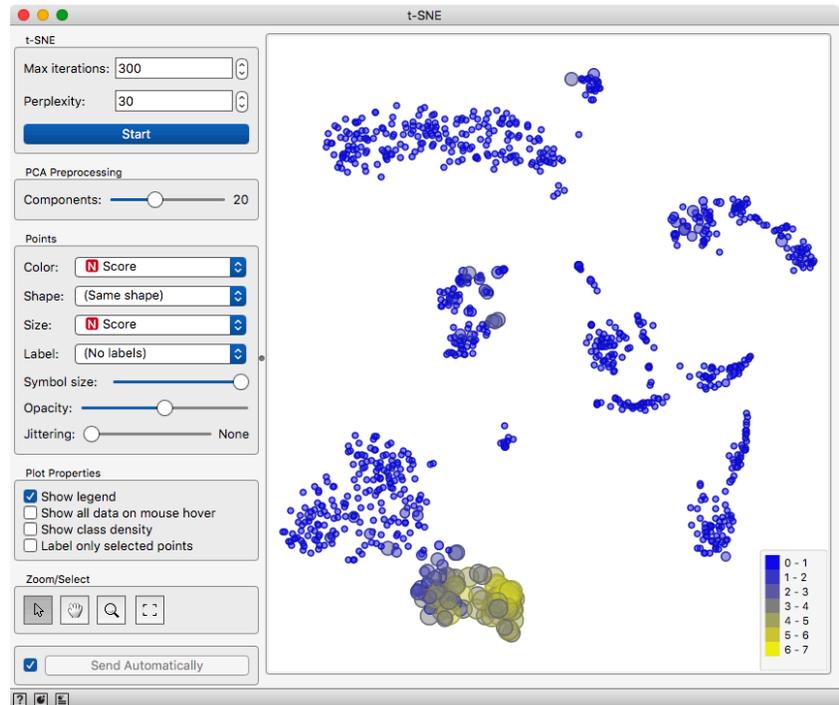
Table, selecting some rows from *Marker Genes* and inspecting the output of this widget in the *Data Table*. We can select one or several marker genes at a time:

The screenshot shows a window titled 'Marker Genes' with a dropdown menu set to 'Human' and a 'Filter...' input field. Below is a table with columns 'Name', 'ID', and 'Cell Type'. The row for 'NKG7' is highlighted in blue.

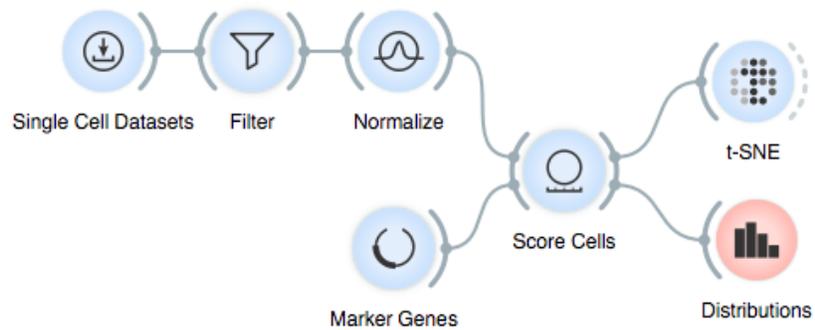
Name	ID	Cell Type
CD14	ENSG00000170458	Monocyte
CD146	?	Endothelial Cell
CD19	ENSG00000177455	B Cell
CD34	ENSG00000174059	Stem Cell/Precursor
CD3D	ENSG00000167286	T Cell
CD4	ENSG00000106610	T Cell
CD79A	ENSG00000105369	B Cell
NKG7	ENSG00000105374	NK Cell
PPBP	ENSG00000163736	Megakaryocyte
S100A9	ENSG00000163220	Monocyte

Widget *Score Cells* assigns a numerical score to each cell that is proportional to an average expression of the marker genes at the input of the widget. The score is added as a meta attribute to the cell data on the output of *Score Cells*. Check this using the *Data Table*! We can now feed this data into *t-SNE* and set the color and size of the points to the cell score:

In the t-SNE widget set the point color to cell score. Open the *Marker Genes* widget and change the selected marker. See the changes in the t-SNE widget. We have just assembled a workflow for interactive display of the marker genes!

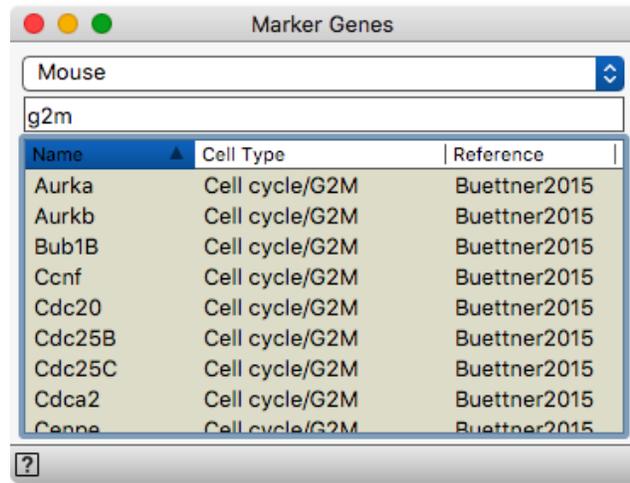


We can now also study the marker genes in the cell cycle data from the previous lesson. Make sure to change the organism in the *Marker Genes* widget to mouse. You will notice that the differences between cells are less pronounced. Use the following workflow:



We would like to compare the scores of the cells for G₂M phase markers versus the G₁/S, and observe if the distribution of the scores does indeed match the types of cells and markers identified in the original work. Open *Marker Genes* and in the mouse section of the genes select those for G₂M (use *Filter* field):

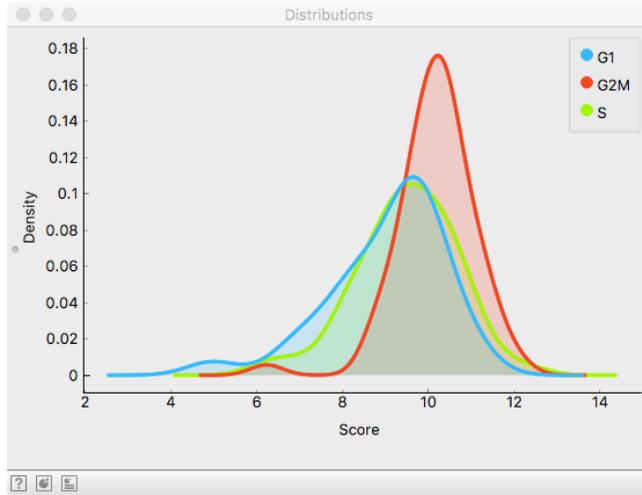
To select all marker genes for a particular cell type, specify the type in the filter entry and press Cmd-A.



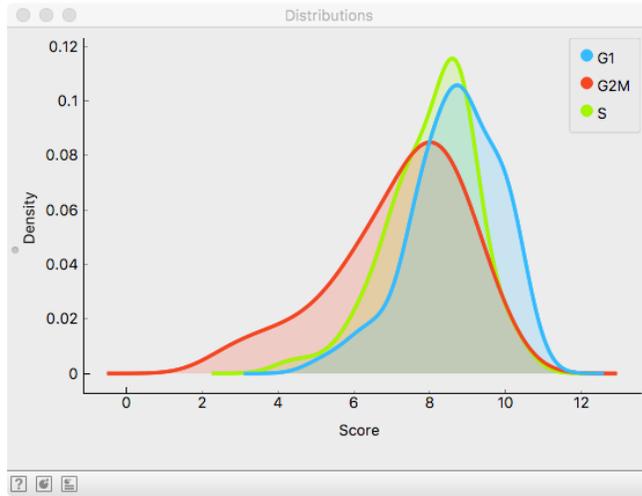
The score distributions for different types of cells indeed confirm the role of the markers: there is a difference in the mean expression of the G₂M phase markers versus the G₁/S.

We can confirm that the markers indeed differentiate between cells in the three different cycles. The differences, though, are not very substantial.

Below is the distribution we obtain for G₂M markers.



And the following is a distribution we obtain for G1/S markers:

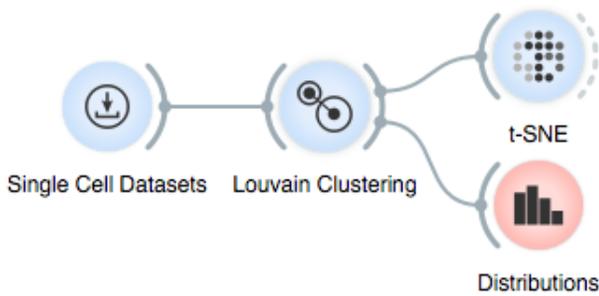


On a single gene level, the difference between cell cycle phases is very narrow but still visible. In the next section, we will introduce methods for finding combinations of genes that will hopefully give use a more clear separation of genes into different classes — cell cycle stages.

Lesson 9: Clustering and Population Discovery

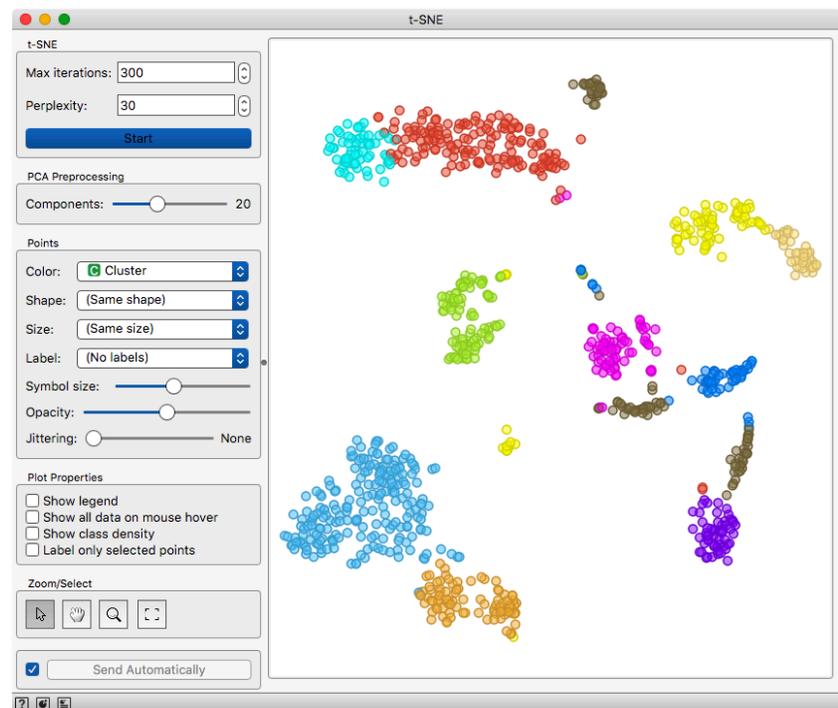
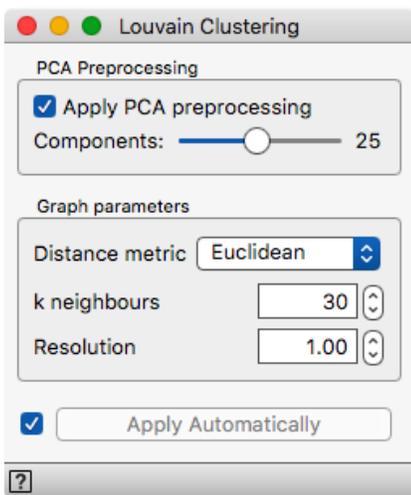
We will again use a sample from Bone marrow mononuclear cells with AML with 1000 cells that contain 1000 most variable genes.

A typical task with a huge single cell dataset is to automatically determine clusters of cells. An advanced method that does this is Louvain clustering. Given the expression matrix, Louvain clustering creates a network of cells based on pairwise distances. Then, it searches for local communities — the parts of the network that are more strongly interconnected than expected by chance (think of friendships on social networks). Each local community is a cluster and genes are assigned cluster labels accordingly.



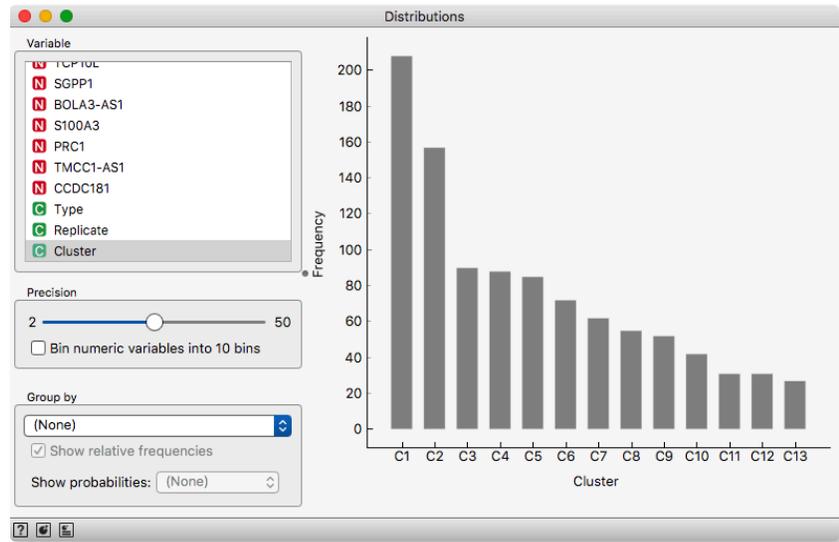
In Orange, *Louvain clustering* is in its own widget that appends a column of cluster labels to the cell data. Quite neatly, the number of clusters is determined automatically. Let us construct a workflow that displays the results of the clustering in the t-SNE plot and that examines the frequency of the cells in each of the clusters. Let us observe the t-SNE plot first.

Louvain Clustering has a number of parameters. Here, we will stay with defaults, but you can experiment, change them and see the effect in t-SNE visualization.



We have colored the points (cells) in t-SNE according to the cluster membership. Notice a nice separation of the clusters in the t-SNE plot. It looks like the cells are also well-separated in the original space of features (genes). A common mistake would be to compute the data projection first and then cluster the projected points. Obviously, then, the clusters would be separated perfectly and there would be no overlap.

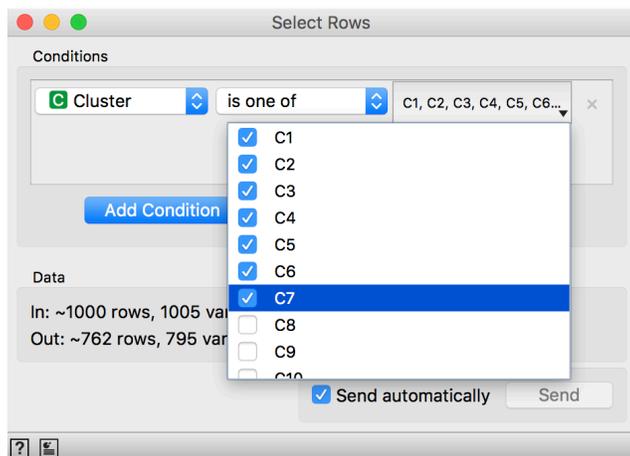
We will now use *Distributions* widget to observe the frequency of the cells within each cluster.



We can see that the number of cells in a cluster is quite unbalanced, ranging from ~220 to a few tens of cells.



We can focus on the larger clusters (and, in a way, remove outliers) by removing the clusters C8 - C13. Hint: compare the effect of this filtering using a t-SNE plot (see screenshot on the left)!

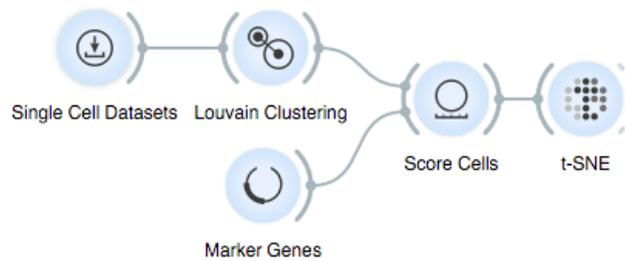


Lesson 10: Differential Gene Expression

How to tell whether clusters make biological sense? One way is to use marker genes, as before, to identify interesting clusters. And then we have to analyze what, besides the marker genes, makes the cells in the clusters special.

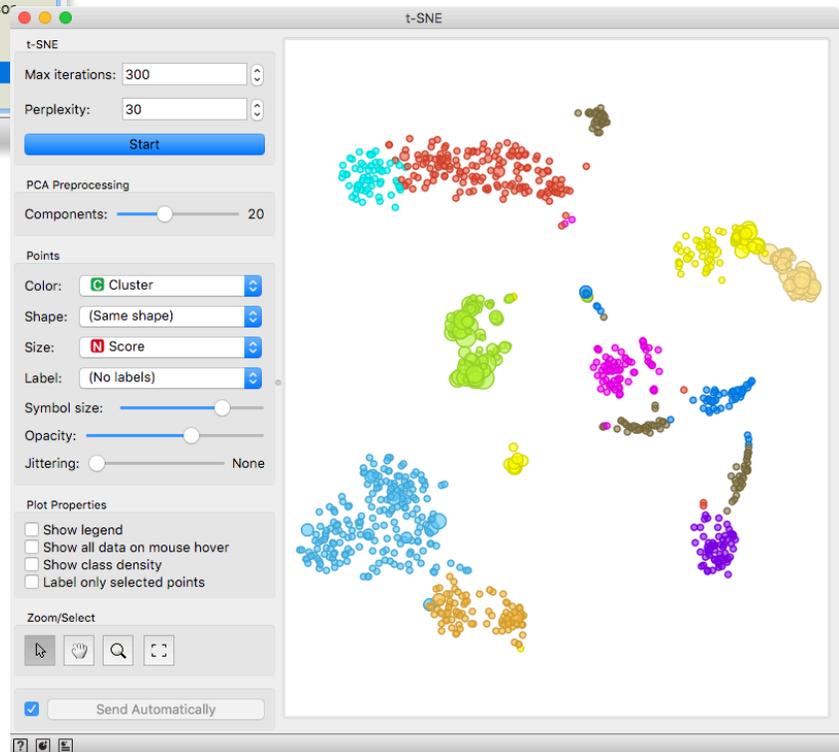
For this analysis, our starting workflow will use clustering, cell scoring and t-SNE visualization:

We could first score cells and then do the clustering, but in the current implementation the workflow on the right is faster as Louvain does not re-cluster every time we change marker genes.



Name	ID	Cell Type
CD14	ENSG00000170458	Monocyte
CD146	?	Endothelial Cell
CD19	ENSG00000177455	B Cell
CD34	ENSG00000174059	Stem Cell/Precursor
CD3D	ENSG00000167286	T Cell
CD4	ENSG00000010610	T Cell
CD79A	ENSG00000105369	B Cell
NKG7	ENSG00000105374	NK Cell

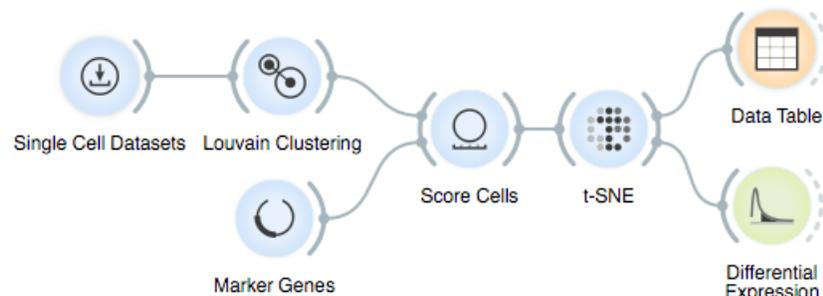
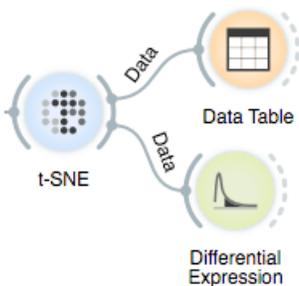
We will choose two marker genes for B-cells and set the t-SNE visualization to color the points by cluster assignments and associate the size of the points to cell's score.



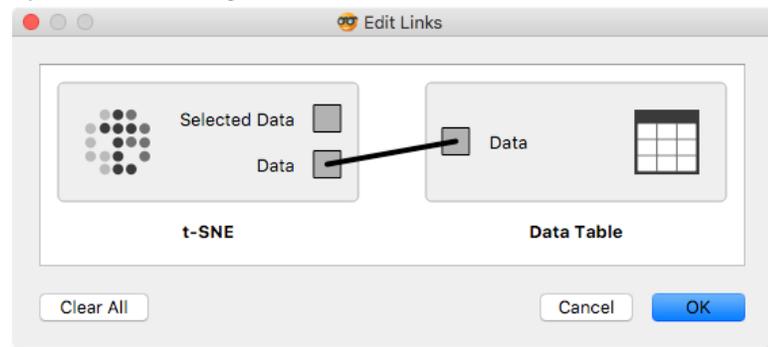
We find that B-cells are spread across several clusters (the green, yellow and light orange one). We can dig in further to find out which genes actually separate the clusters:

1. Use the mouse to select both B-cell clusters (use Shift+Cmd to append to clusters). On the output, in the Data channel, this selection will introduce the new column called *Selected* and mark the selected cells with *Selected = Yes* and all other cells with *Selected = No*.
2. Check the output data in the *Data Table* widget.
3. Connect the t-SNE output “Data” to the *Differential Expression* widget. Set the *Label* (in the *Target Labels* section) to *Selected* and choose *G1* and *G2* as target values.

The workflow visualization in Orange can also report on the types of inputs and outputs that get connected. Use *Preferences* from Orange menus and click on “Show channel names between widgets”.



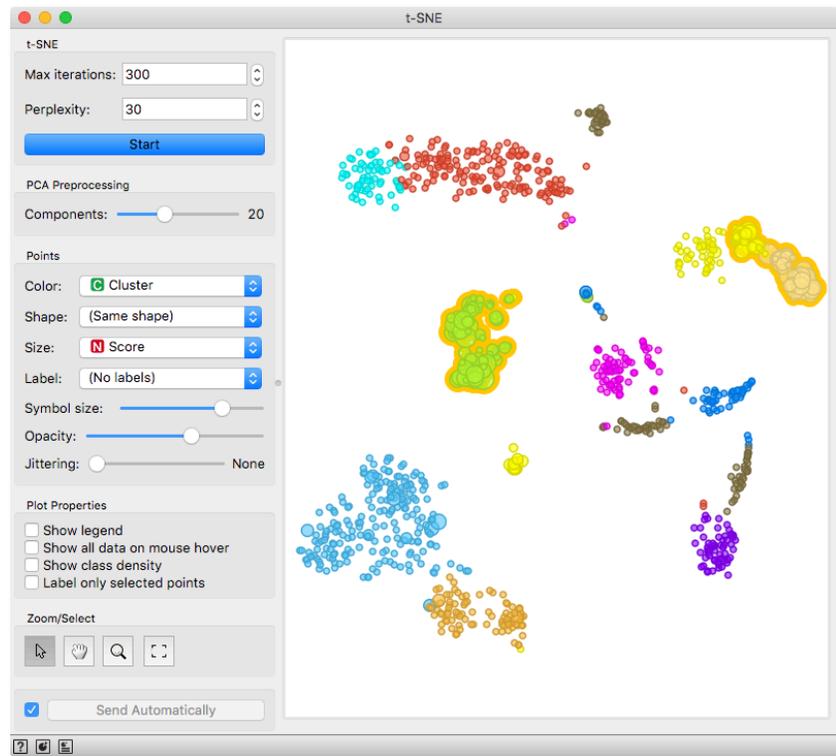
Notice that by default, Orange will connect t-SNE and Data Table through *Selected Data* channel that contains only the selected data. We need to change this, if we want to use the Data channel instead, since there are all the instances of the dataset and the data includes the Selected column. We need to rewire the connections by double-clicking on the t-SNE-Data Table channel:



We should do the same for the connection between t-SNE-Differential Expression widgets.

In the t-SNE widget, the cells in the selected two clusters are marked with orange outlines:

Orange supports several types of selection. In one type, using Shift key modifier, each selected region would be marked differently, with a different color and a different value in the column called *Selected*. Here, we would like all selected cells to be marked with the same label, hence we use Shift-Command modifier when selecting the points in t-SNE.

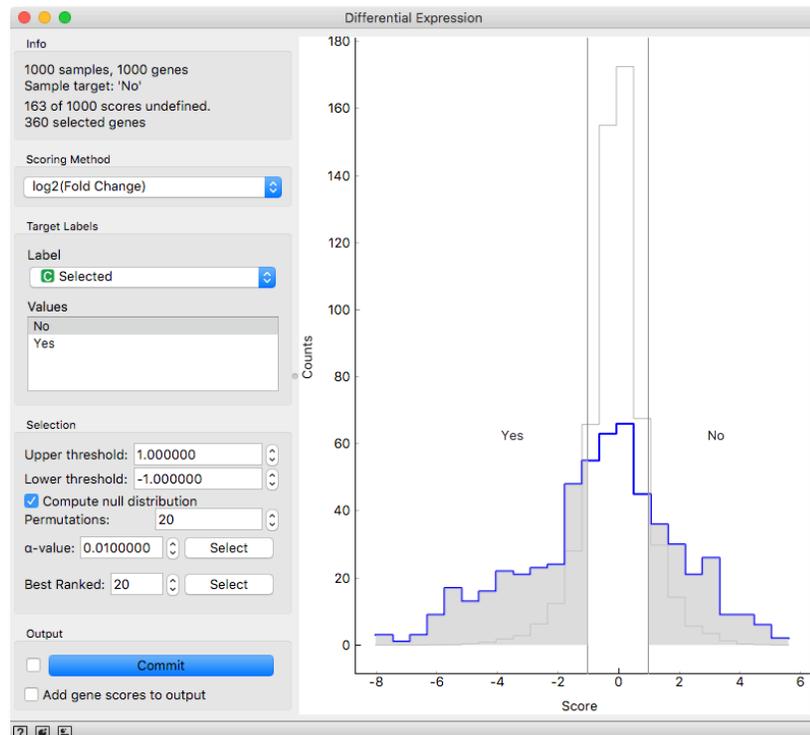


Differential expression shows the distribution of the differences of gene expression in selected and all other cells:

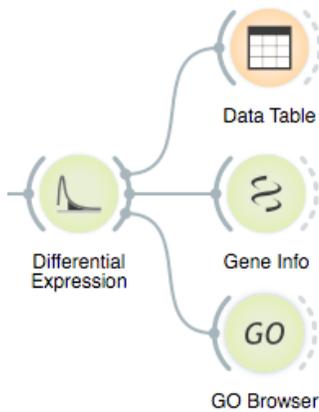
Since our genes are in the columns, *Differential Expression* would output the data with columns that include only the selected genes. Check this out in the *Data Table*.

We have set the scoring method to $\log_2(\text{Fold Change})$ and set the Label to *Selected*. *Differential Expression* also compares the

The distribution marked with grey line shows the *null distribution*, the distribution of genes scores expected by pure chance.



Differential Expression widget outputs the data with genes that are in extremes of the distribution. That is those, for which the difference in selected and non-selected cells is the largest. Genes that are most differentially expressed, lie on the left and on the right side of the two vertical splitters and their score value belongs to the shaded part of the distribution. Move the two vertical splitters such that there are only 50 selected genes.



So, where are the genes that are selected in the *Differential Expression* widget? In the output of the widget. We can observe the new dataset and analyze the set of selected genes with widgets that we connected to *Differential Expression*. Observe the data in the *Data Table*, a list of selected genes in *Gene Info* and the results of analysis of *Gene Ontology* term enrichment in *GO Browser*.

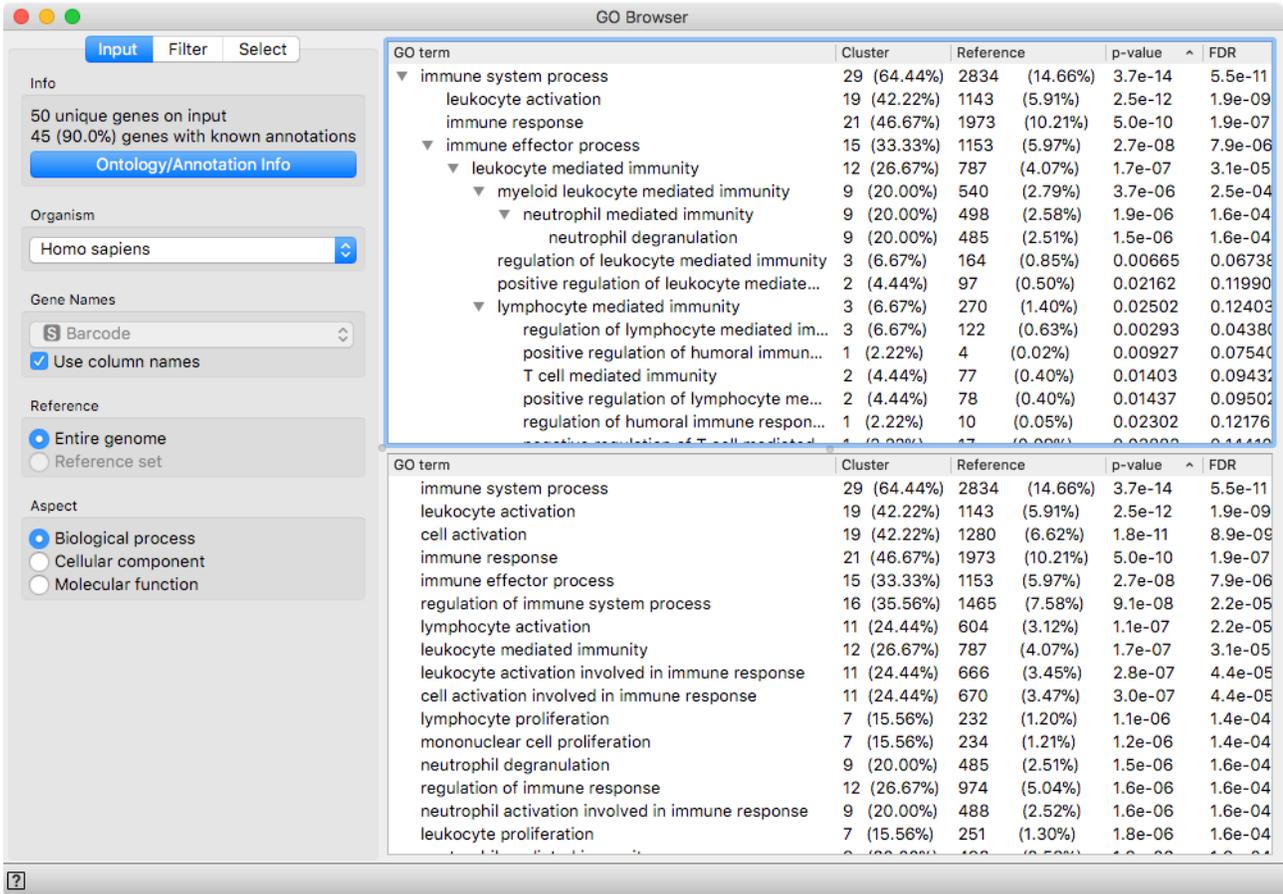
Gene Info displays the list of selected genes. Notice that not all gene names from the selection in *Differential Expression* were matched to the NCBI's database:

The screenshot shows the 'Gene Info' widget interface. On the left, there are controls for 'Info' (50 genes, 47 matched NCBI's IDs), 'Organism' (set to Homo sapiens), 'Gene names' (set to Type), and a 'Commit' button. The main area displays a table of genes with the following columns: NCBI ID, Symbol, Locus, Chromosome, Description, and Synonyms.

NCBI ID	Symbol	Locus	Chromosome	Description	Synonyms
9935	MAFB	-	20	MAF bZIP transcription factor...	DURS3 KRML MCTO
9911	TMCC2	-	1	transmembrane and coiled-c...	HUCEP11
971	CD72	-	9	CD72 molecule	CD72b LYB2
933	CD22	-	19	CD22 molecule	SIGLEC-2 SIGLEC2
931	MS4A1	-	11	membrane spanning 4-domai...	B1 Bp35 CD20 CVID5 LEU-1...
930	CD19	-	16	CD19 molecule	B4 CVID3
929	CD14	-	5	CD14 molecule	
911	CD1C	-	1	CD1c molecule	BDCA1 CD1 CD1A R7
83478	ARHGAP24	-	4	Rho GTPase activating protei...	FILGAP RC-GAP72 RCGAP7...
8115	TCL1A	-	14	T cell leukemia/lymphoma 1A	TCL1
79887	PLBD1	-	12	phospholipase B domain cont...	

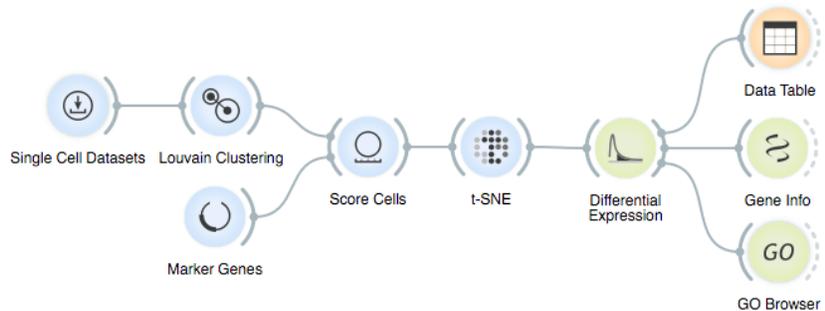
At the bottom of the table, there are two buttons: 'Select Filtered' and 'Clear Selection'.

To validate the biological meaning of differentially expressed genes, we use one final widget: *GO (Gene Ontology) Browser*. Here, we compute the enrichment of biological terms that are pertinent to our subset of selected genes. For B-cells, we expected to find genes related to the immune system. And so we did.



Make sure you have selected *Use column names* in the GO Browser to indicate where to look for gene names. Also, make sure that the right organism has been selected.

Our final workflow is as shown below. Try choosing other markers, select other clusters of cells, and observe the changes in the content of the widgets in our workflow.

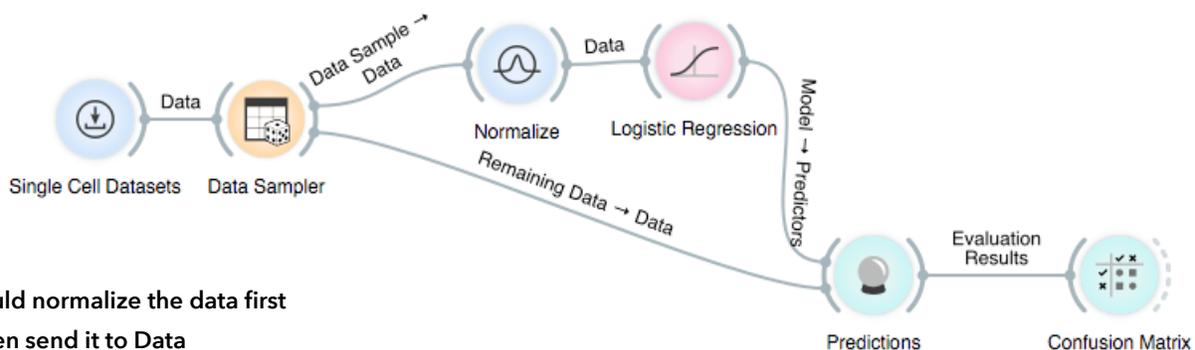


Lesson II: Predictive Modeling

Once cell types or other cell classification are determined, automatic classification models can be used to predict the type of new, unlabelled cells. In this lesson, we will build a simple linear classification model to predict cell cycle stage and learn about the techniques to score the accuracy of our models.

Load the familiar Cell cycle in mESC (Fluidigm) dataset with 182 cells and a reduced set of pre-selected 564 genes. You can experiment with a dataset with a full collection of genes later.

To predict, we need two dataset: one on which to build a model and a separate dataset on which to make predictions. We can simulate this process by splitting our cell cycle dataset to two: the *training* and *test* set — the former will be used for model inference and the latter for evaluation. Evaluation will compare the predicted cell cycle stage to the “true” one in the test set.



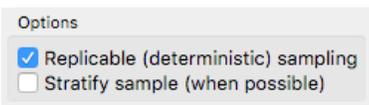
We could normalize the data first and then send it to Data Sampler. But that would be conceptually wrong. Why?

We will use *Data Sampler* widget to split the data such that 70% of the data will go into the training set and the remaining 30% of the data to the test set. *Data Sampler* uses two output channels, and we have to make sure that the right data channel is used. Here is our workflow with annotated communications channels:

We can force *Data Sampler* to always return the same sampling under the same sampling parameters by choosing *Replicable (deterministic) sampling*.

The modeling branch of the workflow first normalizes the data, and sends the normalized data to the *Logistic Regression*, a widget that builds the predictive model. Our data includes a special column called *CycleStageName* and it is this variable that will be modeled from variables that profile the cell, that is, from gene expressions.

The prediction branch of our workflow takes the remaining data and uses the developed model to predict the cell cycle stage. The predictions are made in the *Predictions* widget. Logistic regression predicts probabilities for all the classes in the input data and we



To compare predictions with the true class, compare Logistic Regression column to the column CycleStageName. The widget *Predictions* also outputs a data table with prediction results that store probabilities and predictions and these may be analyzed in other Orange's widgets.

can observe these in the *Predictions*:

	Logistic Regression	cycleStageName	GnaI3	Cdc45	Ccr
1	0.36 : 0.00 : 0.64 → S	S	0	581	0
2	0.00 : 0.01 : 0.99 → S	S	466	110	0
3	0.00 : 0.00 : 1.00 → S	S	492	78	0
4	0.99 : 0.00 : 0.01 → G1	G1	200	978	1
5	1.00 : 0.00 : 0.00 → G1	G1	49	72	1
6	0.42 : 0.46 : 0.12 → G2M	G2M	4	427	0
7	0.89 : 0.00 : 0.11 → G1	G1	1	886	1
8	0.00 : 0.00 : 1.00 → S	S	40	212	0
9	0.90 : 0.02 : 0.07 → G1	G2M	814	194	0
10	0.00 : 1.00 : 0.00 → G2M	G2M	12	16	0
11	0.00 : 0.68 : 0.32 → G2M	G2M	47	29	0
12	0.00 : 0.94 : 0.06 → G2M	G2M	1283	533	1
13	0.00 : 1.00 : 0.00 → G2M	G2M	29	1086	0
14	0.00 : 0.06 : 0.94 → S	S	303	323	0
15	0.02 : 0.98 : 0.00 → G2M	G2M	4411	0	0
16	0.01 : 0.00 : 0.99 → S	S	179	515	1
17	0.99 : 0.01 : 0.00 → G1	G1	5	1242	0
18	0.92 : 0.02 : 0.05 → G1	G1	855	6	1
19	0.01 : 0.00 : 0.99 → S	S	844	2160	0
20	0.00 : 0.11 : 0.89 → S	S	126	650	0
21	0.00 : 1.00 : 0.00 → G2M	G2M	322	1490	1

Notice that the predictions (Logistic Regression column) are mostly correct. There is one error in row 9 where G1 was predicted instead of the true class G2M and where the probability assigned to the true class by logistic regression was only 0.02.

It would be hard to find all correct and false predictions from the table provided in *Predictions* widget. A better way to summarize the results of prediction is with *Confusion Matrix*:

		Predicted			Σ
		G1	G2M	S	
Actual	G1	17	1	2	20
	G2M	1	16	0	17
	S	0	3	14	17
Σ		18	20	16	54

You may try to change the dataset now to the one that includes the full complement of genes. What do you notice?

Out of 54 cells, logistic regression mis-predicted the class for seven cells.

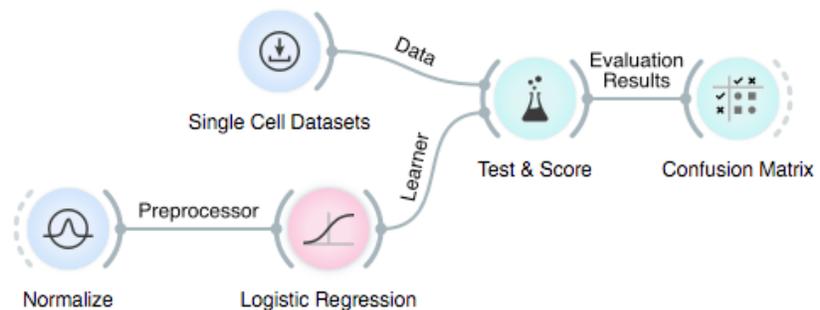
Lesson 12: Cross-Validation

The results of the classification depend on the training set and the accuracy of the test set, of course, depends on the composition of the test data. For speed, change the dataset to the cycle cell data with a subset of genes. Then, open the *Data Sampler* and the *Confusion Matrix* widget and press Sample button in the Data Sampler. See how the classification results change with every fresh instance of the data sampling?

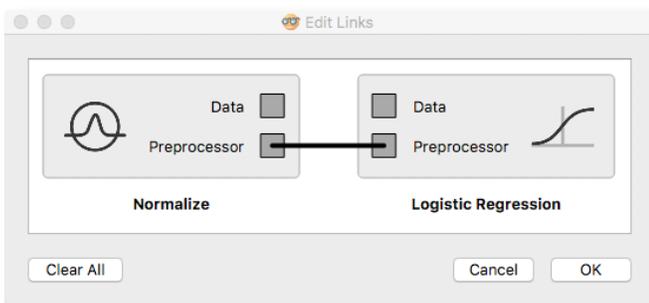
Classification results depend on sampling. When reporting how good is our data and modeling approach, we should sample many times and report on average accuracy. There is a *Test & Score* widget that does this for us and a particular sampling technique called cross-validation that became a standard within machine learning. In cross-validation, the data is split to folds and in each model-and-testing iteration, each fold is used for a test dataset exactly once. The cross-validation workflow is:

Logistic Regression is only one of the many learners in Orange. Another very popular one is *Random Forest*. Try to include it in our workflow!

Make sure that Normalize sends a preprocessor to Logistic Regression. Orange will ask what signals these two widgets should exchange upon connecting them. You can always change the type of connection by double-clicking on it.



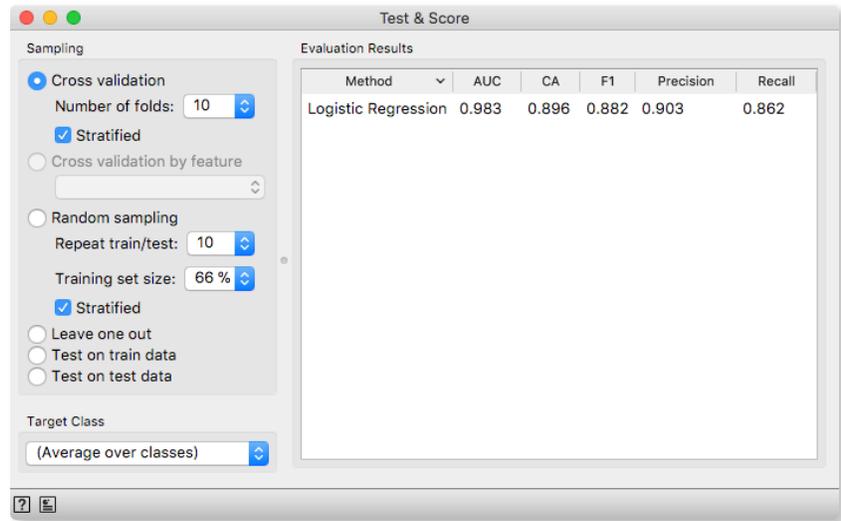
It is the *Test & Score* widget that does data sampling, model construction, and testing. Hence, *Test & Score* needs a method that is used for modeling and this is provided by *Logistic Regression* as a *Learner*. Notice also that we should execute normalization only on the training data, as considering the test data at that stage would be like seeing the future and hence cheating. Normalization is hence a part of the learning and a component that is executed prior to logistic regression. In Orange, this is solved by providing normalization as a preprocessor that logistic regression — or any other modeling technique — includes in the model development procedure.



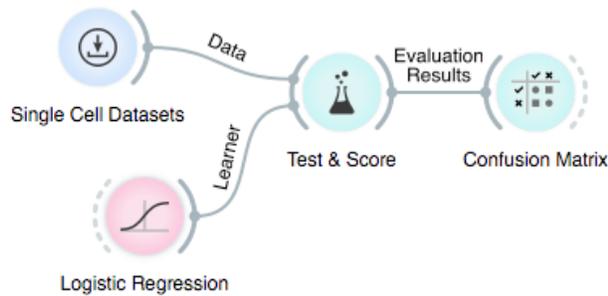
And the results? Open the *Test & Score* widget. It reports on several classification accuracy statistics. The column CA stands for classification accuracy: our predictions were right about 90% of the time. Another often used score is the area under the ROC curve. AUC considers class probabilities, is a discriminative measure and is more reliable than classification accuracy as the score does not depend on class distribution in the dataset. Our models, on average, scored very well:

Test & Score evaluates the predictive performance by counting how many of the test labels our classifier got correct. In a multi-class setting, this is done on a one-versus-all basis.

Each of the learning algorithms has parameters. Double-click on *Logistic Regression* or *Random Forest* (if you have used this learner) to change them.



For a final test, remove the normalization. How do accuracy results depend on this component of data preprocessing?

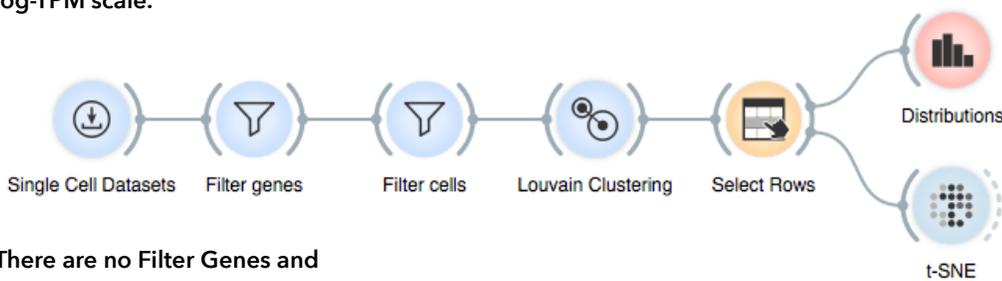


Lesson 13: End-to-End Analysis: Reproduction of Cell-Type Classification

Load the data on mouse bipolar neuron cells (from Shekar *et al.* 2016). Take the larger of the two datasets with 12,000 cells. These data come pre-normalized to log-TPM scale.

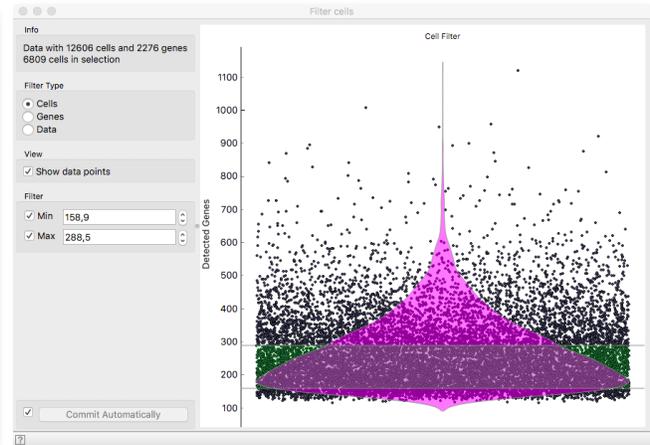
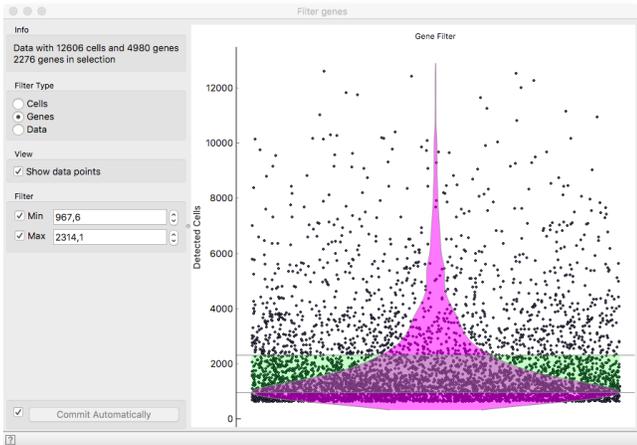
In this task, we will first reproduce the clustering as performed in the paper on mouse retinal bipolar neurons (Shekar *et al.*, 2016). Later, we will apply the cluster prediction on a new, independent dataset of mouse neuronal cells.

We start by filtering cells and genes to retain cells with enough genes expressed and also the genes that are found in a sufficient fraction of cells.

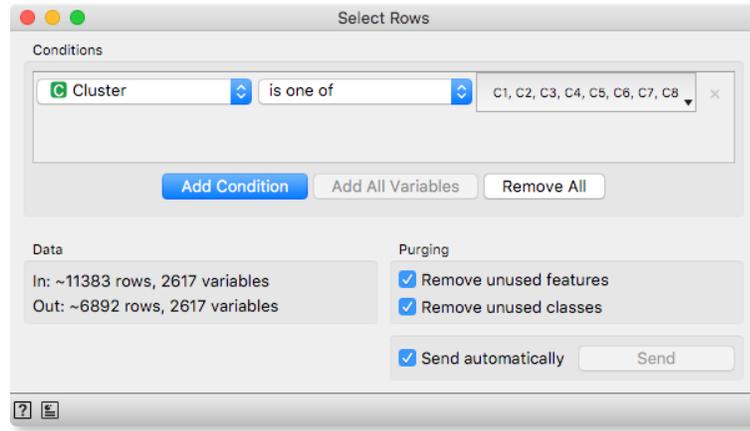


There are no Filter Genes and Filter Cells widgets, we just renamed the Filter widget to denote which kind of filtering we have performed in each of the steps.

This again allows us to retain the data harboring the most information.

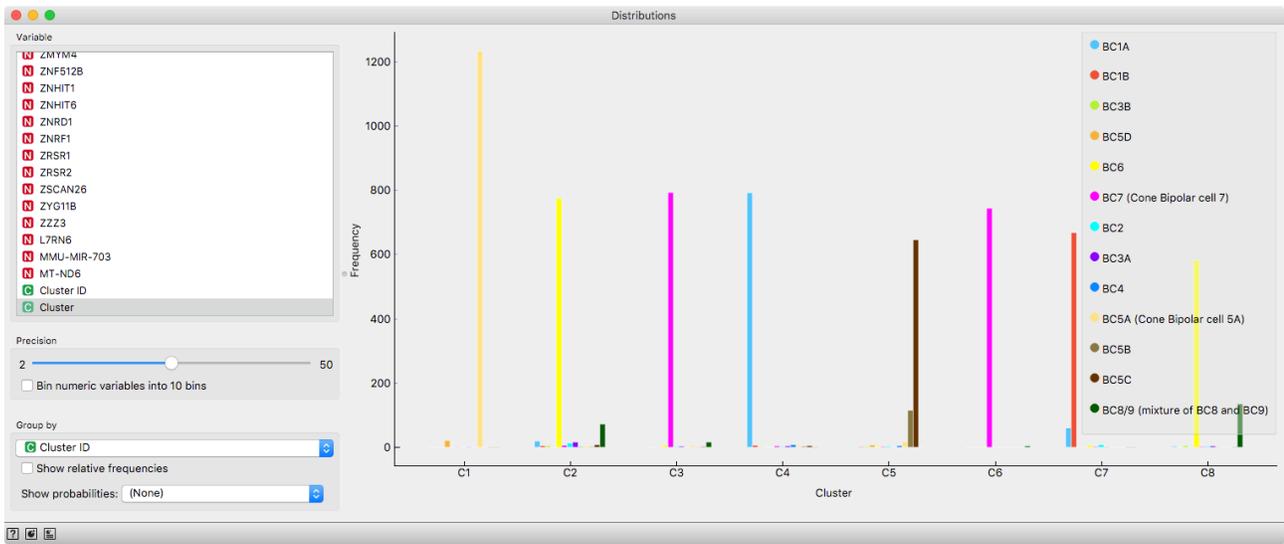


Next, we use the *Lowvain Clustering* to assign a cluster label to each cell. The clusters containing the highest number of cells can be isolated with the *Select Rows* widget by choosing the first, say, eight clusters:



Another mean of validation of clustering results is by showing a t-SNE plot and using the point color (clusters) and shape (original cell class) information.

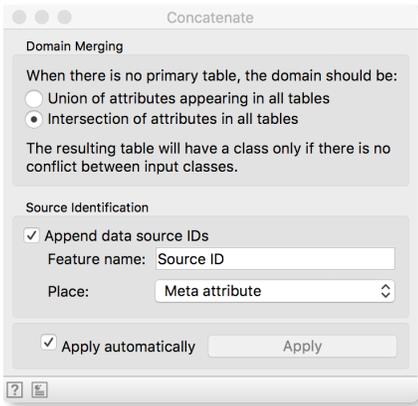
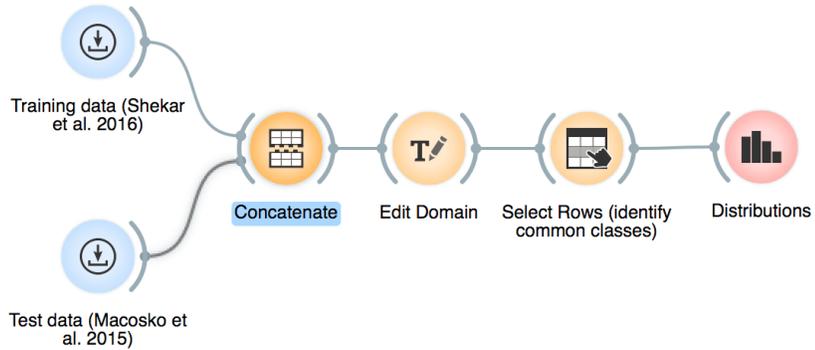
The clustering results can be validated by comparing them to the original labeling provided by the authors. This can be seen in the Distributions widget, as there is a true cluster (*Cluster ID*) always sticking out in each of the assigned clusters (*Cluster*).



Lesson 14: End-to-End Analysis: Combining two datasets

To wrap up, we will build a model to transfer the knowledge obtained on one dataset to another, independent dataset. Our first task is to map the two datasets to use the same representation: find the intersection between the available genes.

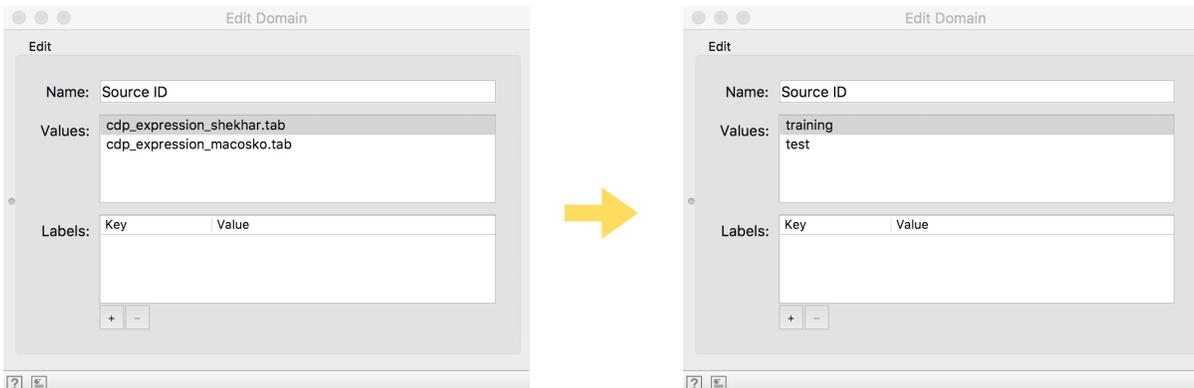
The *Concatenate* widget joins two tables by placing one on top of another. We are interested in the intersection of the available genes.



Load two datasets on mouse bipolar neuron cells (BCs) from Shekar et al. (2016) and Macosko et al. (2015). First, we use *Concatenate* widget to find the intersection of genes in both tables. This step is crucial to later provide the classifier only with the genes that are on disposal in both the training and the test data.

You can use an intermediate *Data Table* to see how many genes were retained. The *Edit Domain* widget is used to give the data sources meaningful names — let’s name the Shekar and Macosko datasets as training and test, respectively.

Using *Edit Domain*, we assign human-readable data source names to cells.

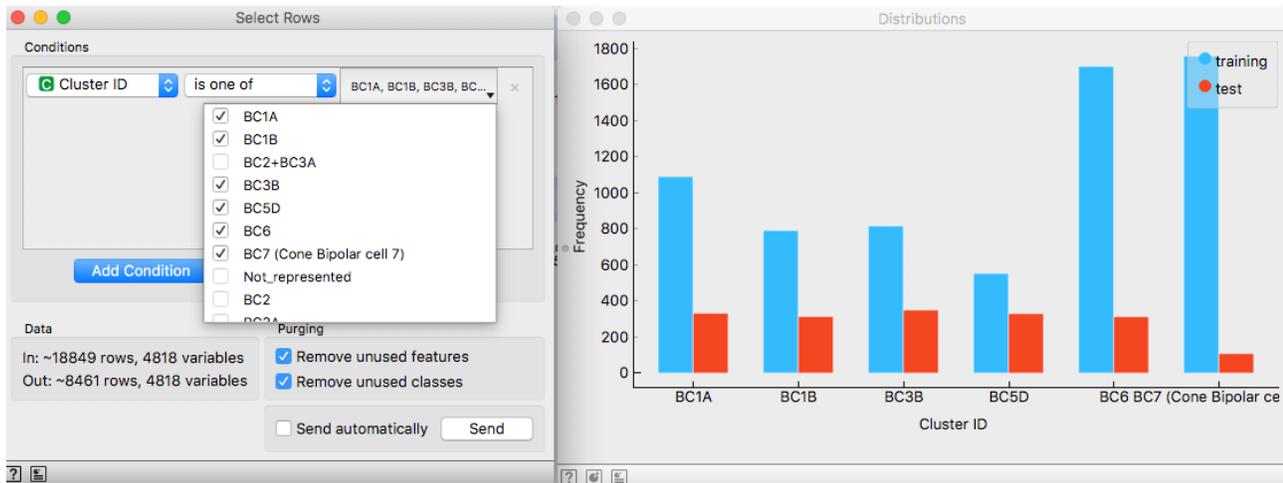


The labelling of the cell types is of course not always available in practice. However, in this idealized scenario, we assume that the same cell types are present in both training and test data.

Use *Select Rows* and *Distributions* (opened simultaneously) to identify cell types present in both training and test data.

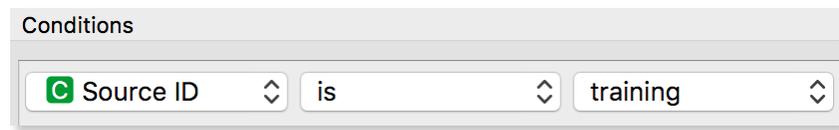
Some differences exist in the available cell types in the two datasets. To quantitatively test the models, we require the set of cell types to be equal in training and test data.

The filtering of cell types is done with *Select Rows* — on classes, to enable a later quantitative evaluation. We only retain cells with a class value (Cluster ID) that is present in both datasets; the *Distributions* widget is of great help here.



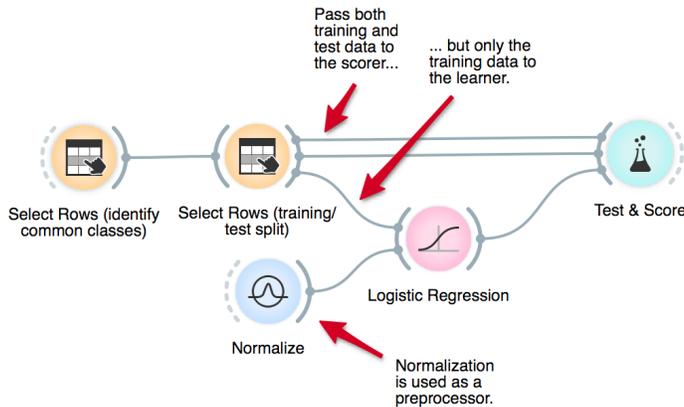
Lesson 15: End-to-End Analysis: Predicting the Clusters on an Independent Dataset

The lesson continues with the data we have prepared in the previous lesson and assemble the most involved workflow so far. Both of the following parts will continue from where the previous lesson left off — the final *Select Rows* widget (used to select the intersecting class values).



Start by splitting the data back into *training* and *test* portions, by using *Select Rows* and the attribute *Source ID*. A *logistic regression* model will be trained on training data to discriminate between the cell types.

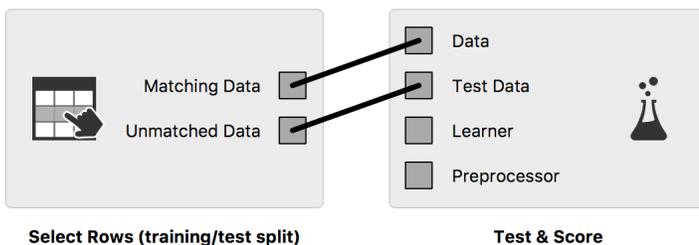
Quantitative model evaluation branch. It continues from the *Select Rows* step used in the previous lesson.



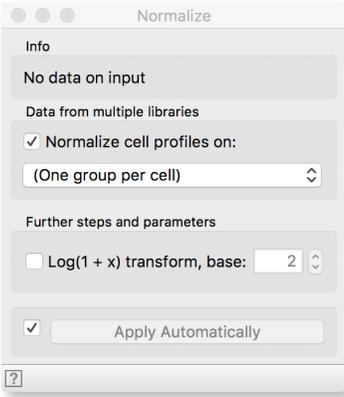
Compute the prediction accuracy with *Test&Score*, by providing it with the data and the model. Here, we use the full power of signals in Orange. Connect the Matching Data from *Select Rows* to Data in *Test&Score*. Similarly, connect Unmatched Data to Test Data. Also, use *Normalize* (note: data is already log-scaled), this time as a preprocessor. This will ensure that normalization is consistently performed before model training and prediction.

Connections (signals) from *Select Rows (training/test split)* to *Test & Score*.

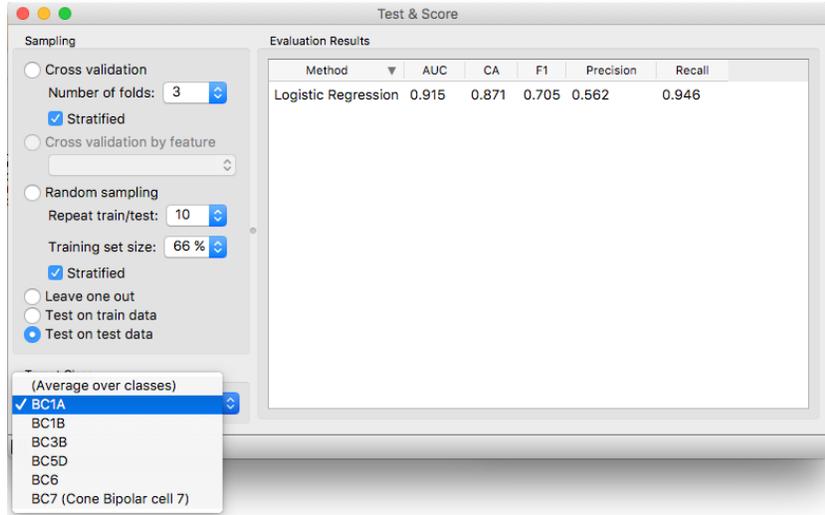
Set "Target class" to one of the classes as seen in the figure below. Some clusters are more separable (e.g. BC₅D) than others (e.g. BC₁B). Could you speculate why?



Settings in the *Normalize* widget, used as a Preprocessor.

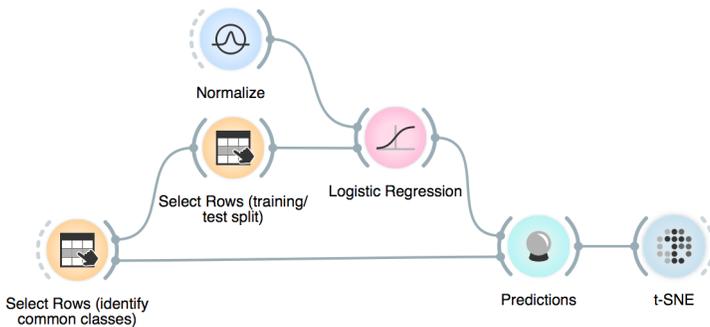


Spoiler alert: see the final t-SNE projections and the complete workflow on the next page.



Finally, let's evaluate the predictions visually. Pass all the data (training + test) through the *Predictions* widget to let the classifier assign a label to each cell. Then, do a projection on of all the data on a common t-SNE plot to visually compare the correspondence

between both datasets. Use shape, color and size of points to display: 1) the source of each cell, 2) the class predicted by Logistic Regression and 3) the probability of one of the classes, respectively. Also, you can use the *Data Sampler* before t-SNE to reduce computational time.

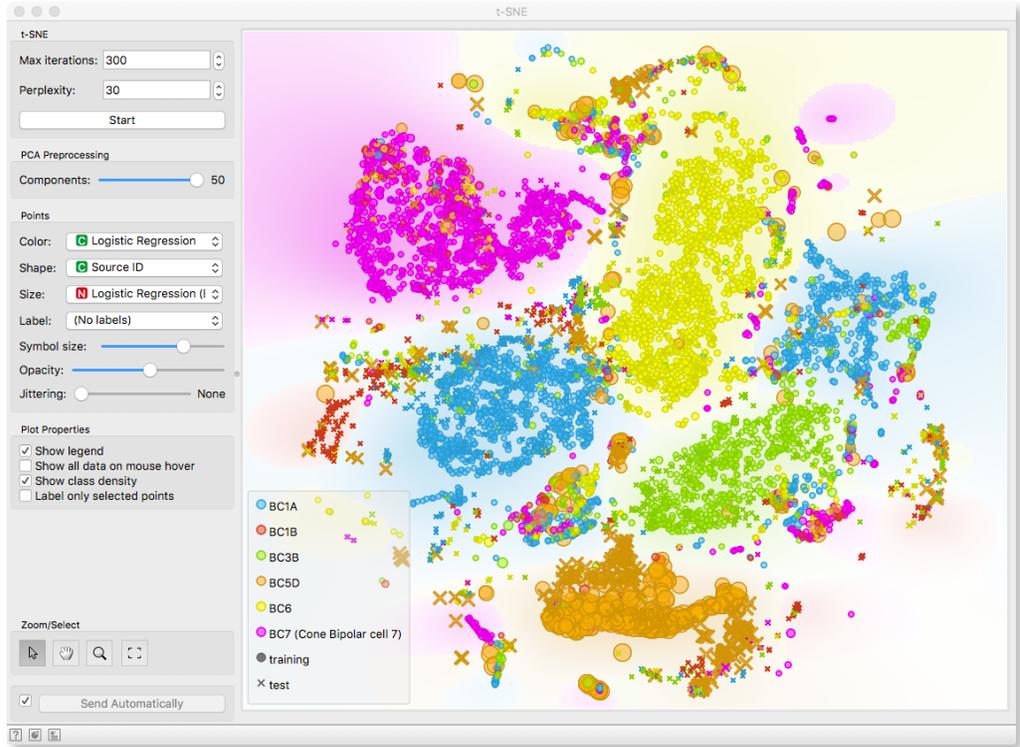


Note: the predictions for the training data will be close to 100% correct, so they can reliably be considered as ground truth. Is this surprising?

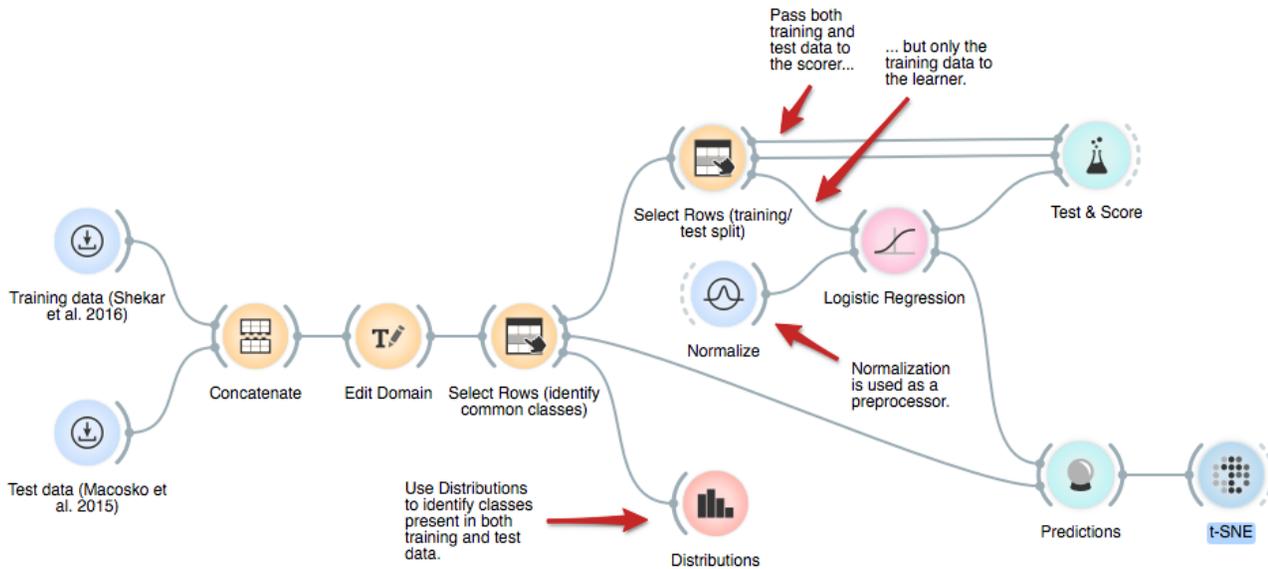
The visual evaluation branch. It continues from the *Select Rows* step placed in the previous lesson.

Does the projection agree with the quantitative evaluation in *Test&Score*? How does the prediction performance change with changing the number of available clusters?

Visualizing predictions in a t-SNE projection. Colors and shapes represent model predictions and data source, respectively.



Overview of the full workflow.

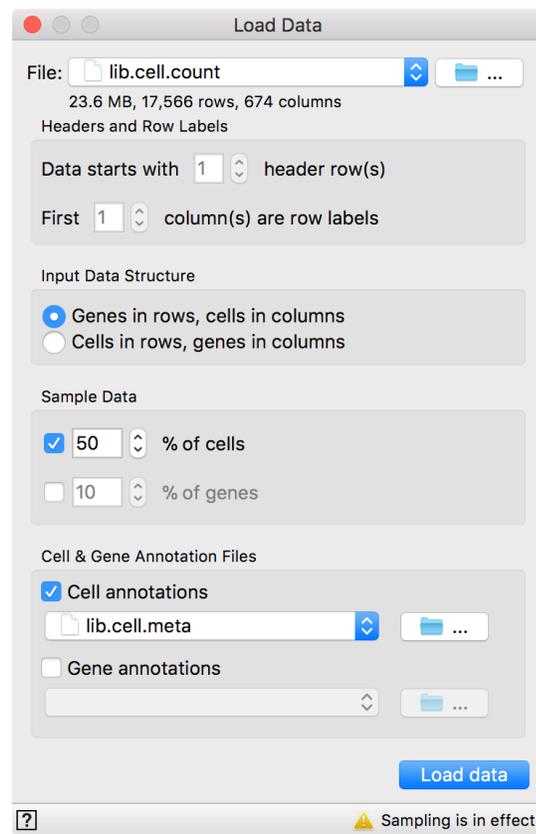
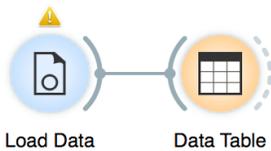


Lesson 16: Loading Single Cell from Tab-Delimited Files

In the lessons above we have worked on already prepared single-cell datasets from the recent literature. Repositories of single-cell data are shaping up and they often contain large datasets available in some format. Most often, tab-delimited files are used.

Data may often be split to data matrix and gene or cell annotation files. scOrange comes with a Load Data widget that can read such files and that supports data sampling. In particular, sampling of the rows may speed up the data loading.

A typical single-cell data exploration workflow starts with data loader and review of the data table. Notice the warning sign above Load Data widget. This is to warn the user that only a sample of the data has been loaded.



For the End

Our single-cell Orange Data Mining course ends here. We covered quite a lot of topics and hope we have taught you some crucial algorithms that should be on the stack of every single-cell data scientists. The goal was not to turn you into one but to get you familiar with some basic techniques, tools, and concepts.