

# A state-of-the-art machine learning pipeline for the analysis of spatial proteomics data

L. Gatto<sup>1,2,\*</sup>, LM. Breckels<sup>1,2</sup>, T. Naake<sup>1,2</sup>, S. Wiczorek<sup>3</sup>, T. Burger<sup>3</sup> and KS. Lilley<sup>2</sup>

<sup>1</sup>Computational Proteomics Unit and <sup>2</sup>Cambridge Centre for Proteomics, Department of Biochemistry, University of Cambridge, UK

<sup>3</sup>Université Grenoble-Alpes, CEA (IRSTV/BGE), INSERM (U1038), CNRS (FR3425), 38054 Grenoble, France

\*lg390@cam.ac.uk

<http://cpu.sysbiol.cam.ac.uk>



## Introduction

Organelle proteomics, or spatial proteomics, is the systematic study of proteins and their assignment to subcellular niches including organelles. pRoloc and pRolocGUI are R/Bioconductor packages that implement all the necessary tools for the sound and reproducible analysis and interactive exploration of spatial proteomics data from any type of experiment such as LOPIT, PCP or PCP-SILAC.

Below, we illustrate a typical pRoloc analysis pipeline:

1. Loading data into R and adding markers
2. QC: checking resolution in the data and organelle markers
3. Detection of new organelle clusters
4. Classification of unlabelled proteins
5. Results, interpretation and visualisation

### 1) Data input

We start by reading quantitative data from 10 fractions (at positions 2 to 11 in the data spreadsheet) sampled along a separation gradient from a csv file and add *Drosophila* organelle markers. This code creates an MSnSet data object (named spat) that stores the quantitative data and the metadata, and can subsequently be easily manipulated, plotted and further processed.

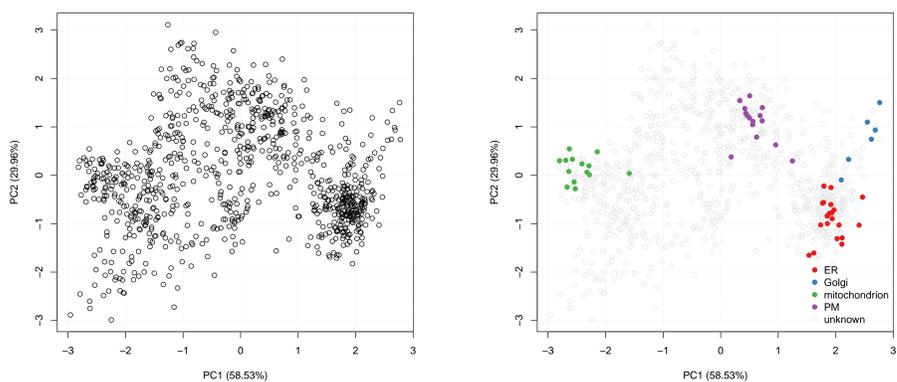
```
spat <- readMSnSet2("quant-data.csv", ecols = 2:11)
spat <- addMarkers(spat, "dmel")
```

	Fraction <sub>1</sub>	Fraction <sub>2</sub>	...	Fraction <sub>m</sub>	...	markers	...
prot <sub>1</sub>	Q <sub>1,1</sub>	Q <sub>1,2</sub>	...	Q <sub>1,m</sub>	...	unknown	...
prot <sub>2</sub>	Q <sub>2,1</sub>	Q <sub>2,2</sub>	...	Q <sub>2,m</sub>	...	organelle <sub>1</sub>	...
prot <sub>3</sub>	Q <sub>3,1</sub>	Q <sub>3,2</sub>	...	Q <sub>3,m</sub>	...	unknown	...
prot <sub>4</sub>	Q <sub>4,1</sub>	Q <sub>4,2</sub>	...	Q <sub>4,m</sub>	...	organelle <sub>2</sub>	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
prot <sub>i</sub>	Q <sub>i,1</sub>	Q <sub>i,2</sub>	...	Q <sub>i,m</sub>	...	organelle <sub>k</sub>	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
prot <sub>n</sub>	Q <sub>n,1</sub>	Q <sub>n,2</sub>	...	Q <sub>n,m</sub>	...	unknown	...
	Fraction <sub>1</sub>	Fraction <sub>2</sub>	...	Fraction <sub>m</sub>	...		
	⋮	⋮	⋮	⋮	⋮		
	⋮	⋮	⋮	⋮	⋮		

### 2) Quality control

We verify on PCA plots that there is structure in the data (i.e. we distinguish well defined clusters, left) and that the markers defined well resolved organelle clusters (right).

```
plot2D(spat)
addLegend(spat)
```



Gatto *et al.* Mass-spectrometry-based spatial proteomics data analysis using pRoloc and pRolocdata. *Bioinformatics*. 2014 May 1;30(9):1322-4.

Gatto *et al.* A foundation for reliable spatial proteomics data analysis. *Mol Cell Proteomics*. 2014 Aug;13(8):1937-52.

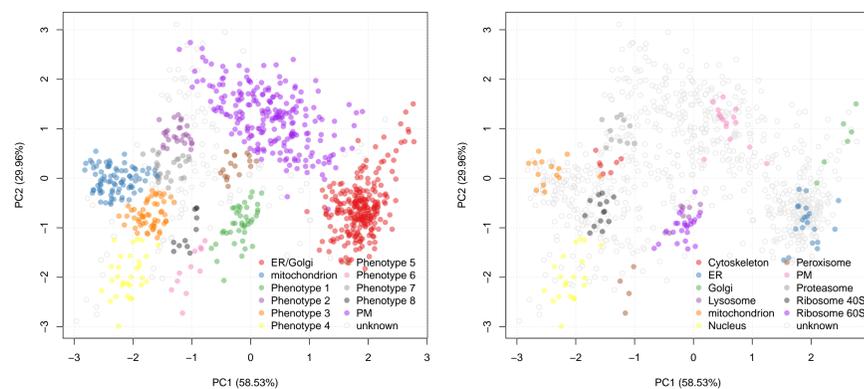
software <http://is.gd/pRoloc>  
 documentation [http://is.gd/pRoloc\\_tutorial](http://is.gd/pRoloc_tutorial)  
 GUI <http://is.gd/pRolocGUI>  
 data <http://is.gd/pRolocdata>  
 Videos <http://is.gd/R4ProteomicsVideos>



### 3) Novelty detection

Our manually curated markers do not cover the entire sub-cellular diversity. We use a semi-supervised machine learning algorithms (Breckels *et al.*, *J Proteomics*. 2013 Aug 2;88:129-40.) to identify new putative organelle clusters, called phenotypes (left), which require validation by the user (right).

```
spat <- phenoDisco(spat)
plot2D(spat, fcol = "pd")
```



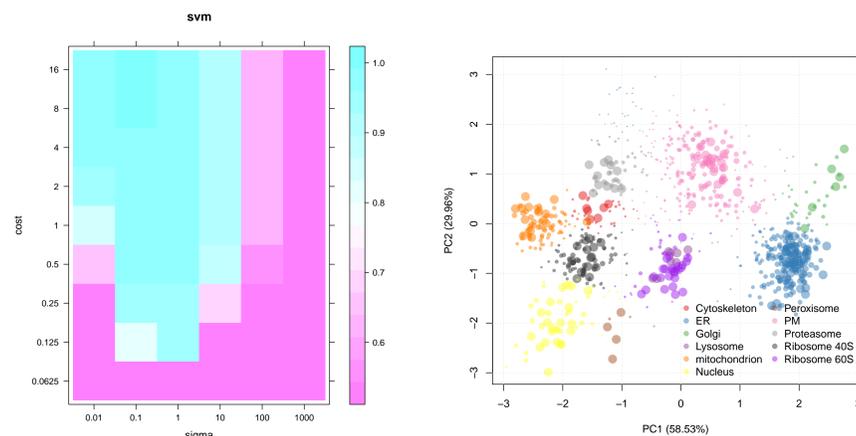
### 4) Classification

We can now classify unlabelled proteins to any of the augmented classes using a supervised machine (SVM) learning algorithms, for example, a support vector machine classifier. It is essential to tune the classification model parameters (here *sigma* and *cost*) prior to actual classification (left).

```
params <- svmOptimisation(spat, fcol = "pd.markers")
spat <- svmClassification(spat, params, fcol = "pd.markers")
```

The classification algorithm calculates classification probabilities that reflect the distance of a protein to the decision boundaries defined by the SVM model (right). Eventually, the data can be exported to a spreadsheet file.

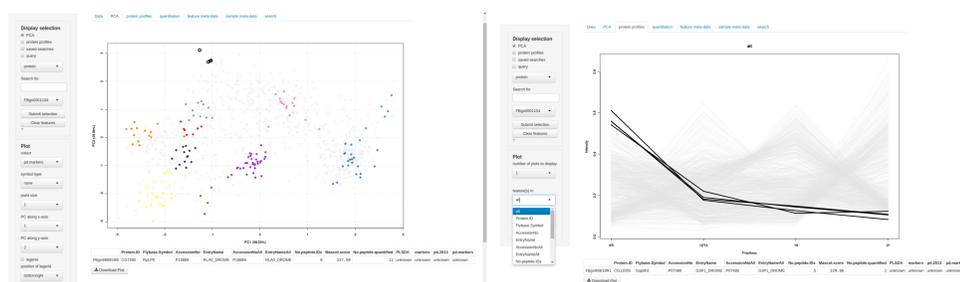
```
ptsze <- exp(fData(spat)$svm.scores) - 1
plot2D(spat, fcol = "svm", cex = ptsze)
write.exprs(spat, file = "spat-results.csv")
```



### 5) Interpretation

The graphical user interface implemented in the pRolocGUI package enables one to interactively explore the data.

```
library("pRolocGUI")
pRolocVis(spat)
```



This work was supported by the European Union 7<sup>th</sup> Framework Program PRIME-XS project and a BBSRC Tools and Resources Development Fund.