

# Use case: mapping sparse spatial data with TOPCAT

This use case describes a workflow related to large hyperspectral datasets.

In this example you will use data from the VIRTIS/Rosetta experiment and study the distribution of observation of the surface of comet 67P. The workflow consists in: retrieving the data files; preparing a data table restrained to spectra of interest; mapping and analyzing the footprints in TOPCAT.

## 1. Getting the data

The data are available on the PSA archive and on the PDS Small Bodies Node. The PSA has a graphic interface (<https://archives.esac.esa.int/psa/#!Home%20View>, Fig. 1) which helps select data files from their global properties. If the data are distributed in an EPN-TAP service, it can be searched and filtered using the VESPA portal (<http://vespa.obspm.fr>, Fig. 2).

The PSA interface can help identify files containing data of interest, although to a lesser extent than the VESPA portal.

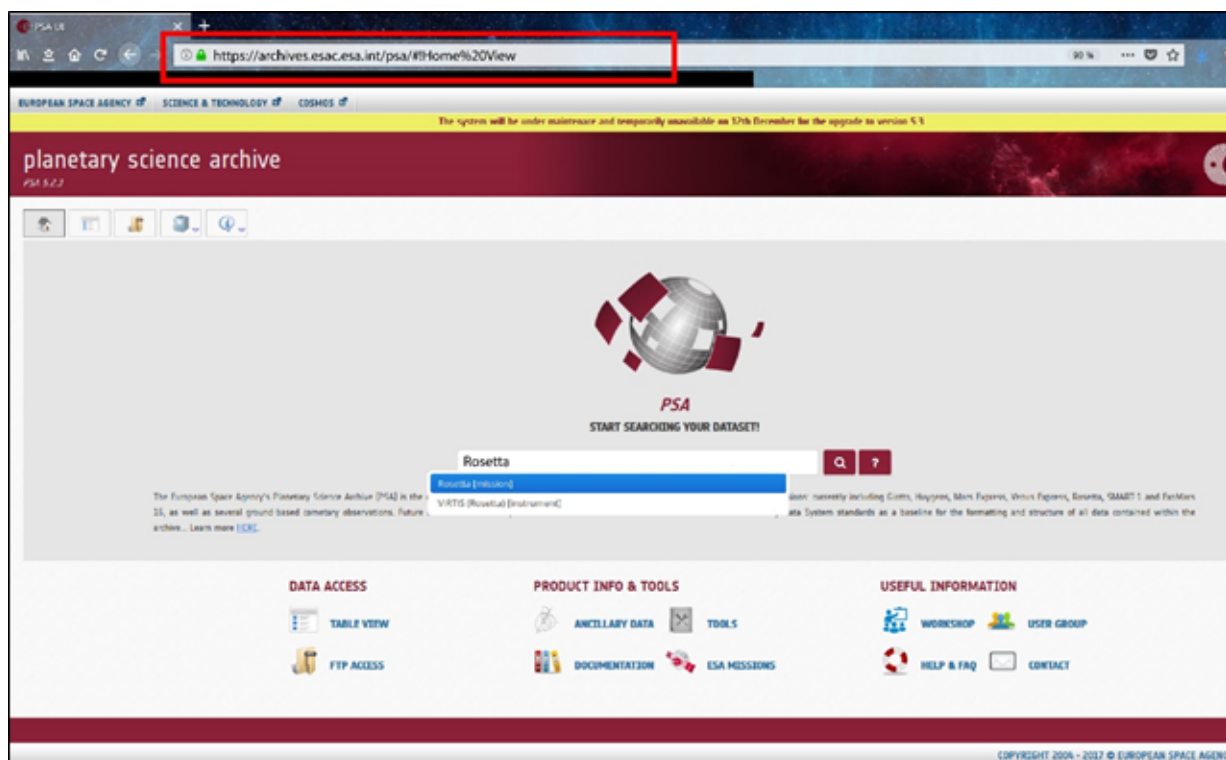


Fig. 1 PSA web interface.

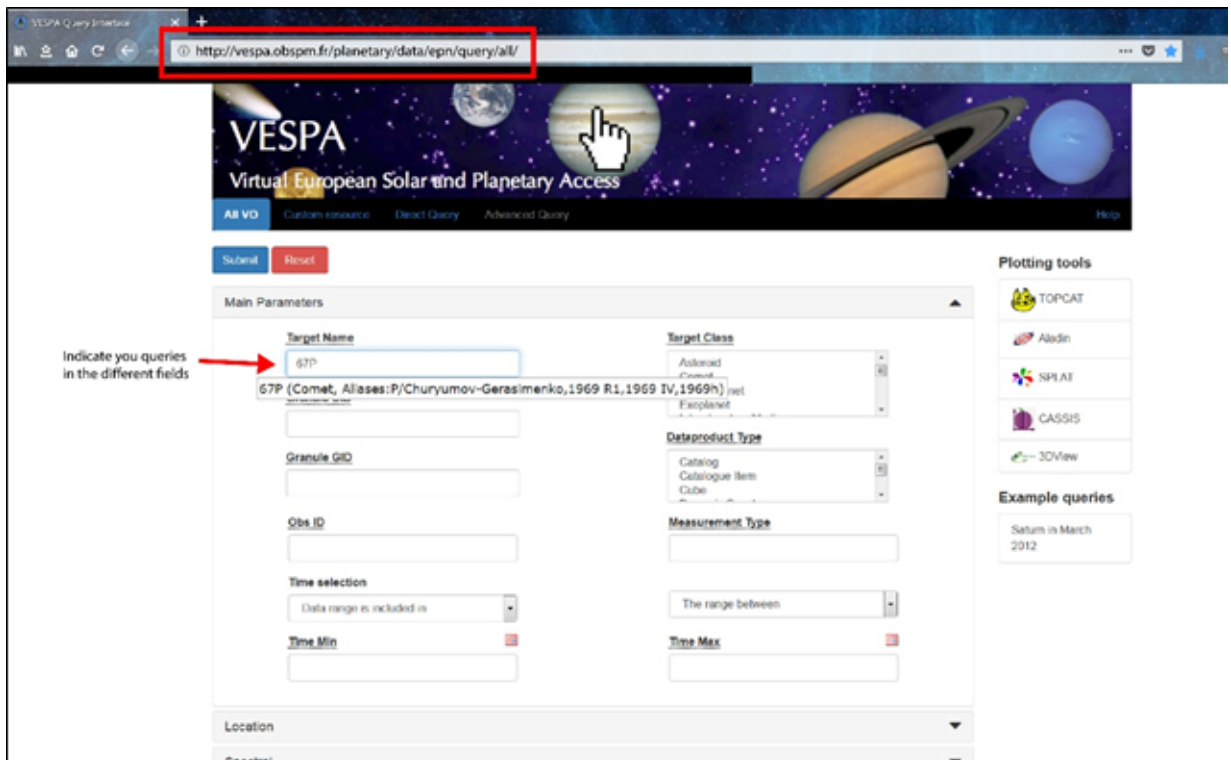


Fig. 2 VESPA web interface.

## 2. Preparing the data table

In this step, you want to prepare a table of individual spectra of interest and store it in a format suitable for TOPCAT. This file format can be VOTable or the FITS format which are both current standards, as they contain self-described data. This step is important when studying a global dataset, because it settles the tedious process of going through a number of files once for all.

There are several possibilities to prepare the FITS/VOTable for TOPCAT:

1) Data selection and writing can be done in IDL (Fig. 3) or a similar language. This is the general case where data files are retrieved from the PSA or PDS archives (files need to be open and spectra need to be filtered). IDL (or its GDL open source clone) is anyway required to read PDS files and select individual spectra. Even if the data are loaded in IDL at this point, using TOPCAT is still useful because of its superior ability to integrate data into healpix cells in the 3<sup>rd</sup> step of this tutorial.

2) If the data are distributed as an EPN-TAP service describing individual spectra (a project for this particular dataset) the VESPA portal, or a TAP query issued from TOPCAT, can alternatively be used to retrieve the selected information and write it as a table; this would mainly save the next step.

The goal is to create a table where each line corresponds to a single pixel / observation / acquisition and each column to a physical (e.g. radiance at a given wavelength) or geometric (e.g. latitude/longitude) parameter or any other information (spectral parameters, reference of the hyperspectral cube, identification number of this acquisition, period/duration of the acquisition...). Only metadata are stored in the table, because the aim is to study the footprints. Pixel / observation / acquisition here refers to a pixel in a hyperspectral VIRTIS cube. In the present case, we want to retain only acquisitions that match several criteria: FoV intercepting the surface, given illumination conditions, minimum resolution, etc.

The Fig. 3 below illustrates the formatting with IDL.

```

1: pro tuto_Vo_VIRTIS, listofcube, listofcubegeometry
2
3: for i=0,n_elements(listofcube)-1 do begin
4:   cube = virtispds(listofcube(i))
5:   geo = virtispds(listofcubegeometry(i))
6
7:   ; here you will filtered/make operations on data. For example:
8:   reflectance = mean(cube.qube(138:148,*,*),dim=1,/double) ; mean reflectance around 2.3µm
9:   incidence = geo.cube(17,*,*) ; extract the incidence
10: ; Latitude = ....
11: ; Longitude = ...
12: ; .
13: ; .
14: ; .
15: ; Thanks to the geometric cube I am able to extract information which allows to map the spectra
16
17: ; In order to write the information in a VOTable you must store these variable in an array at each iteration:
18: all_reflectance = [reflectance,all_reflectance]
19: all_incidence = [incidence,all_incidence]
20: endfor
21
22 ; then put all the array corresponding to a parameter (a futur column in your VOTable) in a structure
23 my_structure = {Reflectance: all_reflectance,$
24               Incidence: all_incidence}
25
26 ; And finally, write your VOTable thanks to write_vot.pro or to vobs_Struct2VOTable
27 vobs_Struct2VOTable,my_structure,'/here/is/my/pretty/VOTable.vot'
28 END
29

```

Fig. 3 Illustration of the data processing with IDL.

There are three steps:

1. Reading and filtering the data files
2. Packing them into a structure
3. Writing the structure as a VOTable or as a FITS table

1. The first step is data access and filtering, the most time-consuming process. PDS3 data files can be read under IDL with the `readpds` library from the PDS Small Bodies Node, or with the all-purpose `virtispds` library from LESIA (also running in the open source GDL environment). In this example, the input is a list containing the path of the hyperspectral cubes to process – the subset used here covers more than 6 months of VIRTIS/Rosetta observations in the visible range. All geometry parameters have been computed during the earlier calibration step using the SPICE library and are ready for use in the archive.

2. The information of interest is extracted from calibrated spectral cubes and geometry cubes: e.g. radiance, reflectance, incidence /emergence/phase angles, id numbers, altitude, elevation, acquisition time, cube from which it is extracted, latitude(s)/longitude(s), x/y/z coordinates... anything potentially useful to study the repartition of observations at the comet surface. Each parameter will constitute a column in the final table.

Cubes are processed in a loop; the information is extracted for the current cube, each parameter in a different array, then concatenated in a global array for all sessions. Spectra that cannot be exploited in a later phase may be filtered out during the extraction process, based on several criteria (observations angles, spacecraft altitude, signal threshold...). After the extraction process, the parameter arrays are associated in a structure, which is the required input to the table writing routines.

3. The final step consists in writing the FITS and/or the VOTable. While both formats will provide the same result, the FITS format is probably the most suitable since it is written in binary while VOTable in ASCII. Then, in the present case the table size is ~ 4GB in the VOTable format and only ~1GB in FITS format. As explain below, the loading time in TOPCAT is also faster for FITS table.

A FITS table can be write from IDL structure with:

- `mwrfits.pro` (<https://idlastro.gsfc.nasa.gov/contents.html> part of the [Astron NASA Library](#) - basic syntax is: "`MWRFITS, your_structure, filename, /CREATE`")

At least two routines are available to write VOTables from IDL:

- `write_vot.pro` ([https://github.com/ejn-vespa/IDL\\_VOTable](https://github.com/ejn-vespa/IDL_VOTable)); requires the stilts java library to be accessible (current version runs under Linux, Mac and Windows)

- `vobs_Struct2VOTable.pro` ([http://www.heliodocs.com/php/xdoc\\_list.php?dir=\\$SSW/vobs/gen/idl](http://www.heliodocs.com/php/xdoc_list.php?dir=$SSW/vobs/gen/idl) part of the SSW library)

Note that the FITS table must be written in binary (default behavior) to gain space and reading time.

Ideally, each parameter/column should be described with name, unit, and UCD for future reference and use in other tools. However, TOPCAT can ingest the data in this basic presentation.

### 3. Data visualization in TOPCAT

Our table gathers all relevant information from many files, which can now be handled easily. This file can be read by the TOPCAT tool (<http://www.star.bris.ac.uk/~mbt/topcat/>) which provides very powerful plotting functions (Fig. 4).

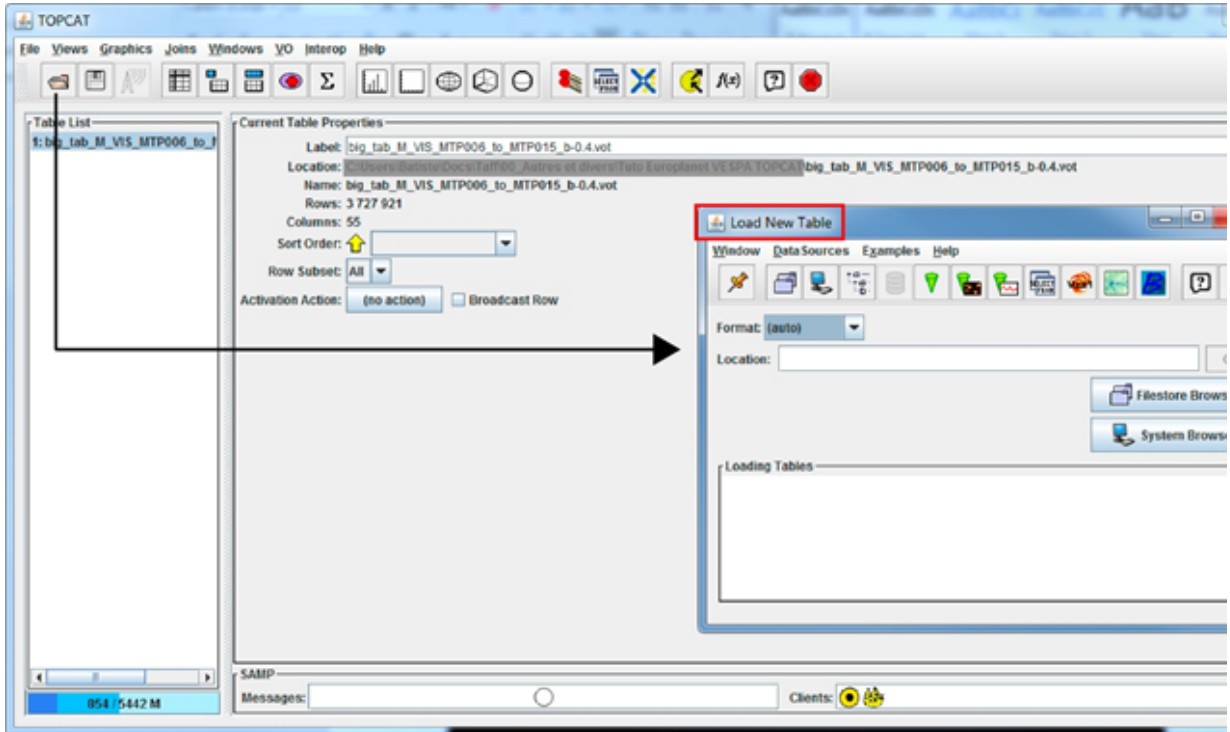


Fig. 4 Open a new table from the main interface by browsing your directory.

Our table contains 3 727 921 lines and 55 columns. TOPCAT is fast, even with huge tables. To read the VOTable on a busy Intel Core i7 @ 2.40GHz with 24Gb of RAM takes 2.15min (Java in 64 bits). However, the same table, written in FITS format, will be read instantaneously by TOPCAT. This table can be browsed from the interface (Fig. 5), which is useful to identify with precision an observation: when clicking on a dot in a plotting window, the corresponding line is highlighted in the table (and vice versa).

TOPCAT can perform further filtering of the data, so that you don't need to get back to the previous step. A standard way to achieve this is to define a new subset: click the **Display subset** icon (red and violet ellipses), then the **Define a new subset** icon (green cross) and enter an algebraic expression (you can also define a subset from the table itself using the top left icon in this window).

TOPCAT(1): Table Browser

Table Browser for T\_Big\_Web\_M\_VIS\_MIP066\_to\_MIP015\_b04.rst

POS_ID	MTP	STR	CLUB	ID_CLUB	TARF	ELAPPRO_1	PHI_MID	ID_OBJ	NUMBER	FLUXDENS	AL_TITLUC	REFLECTION	LOCALLOC	PLAT_S&D	DELT_S&D	SCAM_W	SCAM_MSD	
1	1	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	245.	5	Apht	5.07	90.356	22.9912	3.8931	6.51394	4.91939	3.54801	6.17983	0.719139
2	2	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	246.	5	Apht	5.052	90.32	22.9909	3.8931	6.51397	4.92251	3.54801	6.17983	0.719139
3	3	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	247.	5	Apht	5.035	90.299	22.9793	3.8931	6.51927	4.92264	3.54801	6.17983	0.719139
4	4	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	248.	5	Apht	5.028	90.279	22.9712	3.8931	6.5217	4.92276	3.54801	6.17983	0.719139
5	5	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	249.	5	Apht	5.029	90.25	22.9632	3.8931	6.52555	4.92271	3.54801	6.17983	0.719139
6	6	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	250.	5	Apht	-5.02	90.218	22.9543	3.8931	7.89327	4.92264	3.54801	6.17983	0.719139
7	7	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	251.	5	Apht	-5.021	90.18	22.9465	3.8931	7.24337	4.92264	3.54801	6.17983	0.719139
8	8	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	252.	5	Apht	-5.022	90.158	22.9478	3.8931	7.23294	4.92264	3.54801	6.17983	0.719139
9	9	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	253.	5	Apht	-5.023	90.133	22.9404	3.8931	7.23383	4.92264	3.54801	6.17983	0.719139
10	10	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	254.	5	Apht	5.041	89.66	22.9207	3.8931	6.49327	4.92264	3.54801	6.17983	0.719139
11	11	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	255.	5	Apht	5.038	89.654	22.914	3.8931	6.49355	4.92264	3.54801	6.17983	0.719139
12	12	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	256.	5	Apht	5.039	89.638	22.9089	3.8931	6.49322	4.92264	3.54801	6.17983	0.719139
13	13	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	257.	5	Apht	5.032	89.607	22.9013	3.8931	6.49353	4.92264	3.54801	6.17983	0.719139
14	14	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	258.	5	Apht	5.042	89.589	22.8963	3.8931	6.49355	4.92264	3.54801	6.17983	0.719139
15	15	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	259.	5	Apht	5.044	89.563	22.9049	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
16	16	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	260.	5	Apht	5.044	89.57	22.9163	3.8931	6.49355	4.92264	3.54801	6.17983	0.719139
17	17	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	261.	5	Apht	5.037	89.246	22.9123	3.8931	6.49353	4.92264	3.54801	6.17983	0.719139
18	18	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	262.	5	Apht	5.037	89.229	22.9078	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
19	19	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	263.	5	Apht	5.047	89.22	22.9037	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
20	20	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	264.	5	Apht	5.042	89.223	22.904	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
21	21	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	265.	5	Apht	5.038	89.208	22.9027	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
22	22	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	266.	5	Apht	5.032	89.202	22.9013	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
23	23	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	267.	5	Apht	5.03	89.182	22.8987	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
24	24	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	268.	5	Apht	5.038	89.183	22.8943	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
25	25	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	269.	5	Apht	5.046	89.174	22.8943	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
26	26	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	270.	5	Apht	5.045	89.144	22.892	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
27	27	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	271.	5	Apht	5.043	89.143	22.892	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
28	28	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	272.	5	Apht	5.041	89.15	22.892	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
29	29	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	273.	5	Apht	5.04	89.143	22.8943	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
30	30	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	274.	5	Apht	5.04	89.134	22.8943	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
31	31	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	275.	5	Apht	5.04	89.124	22.8917	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
32	32	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	276.	5	Apht	5.041	89.116	22.8939	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
33	33	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	277.	5	Apht	5.04	89.11	22.879	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
34	34	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	278.	5	Apht	5.04	89.102	22.879	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
35	35	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	279.	5	Apht	5.04	89.096	22.8741	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
36	36	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	280.	5	Apht	5.048	89.083	22.8712	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
37	37	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	281.	5	Apht	5.037	89.083	22.8693	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
38	38	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	282.	5	Apht	5.048	89.083	22.8693	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
39	39	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	283.	5	Apht	5.038	89.082	22.8693	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
40	40	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	284.	5	Apht	5.039	89.046	22.862	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
41	41	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	285.	5	Apht	5.04	89.042	22.861	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
42	42	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	286.	5	Apht	5.041	89.038	22.86	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
43	43	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	287.	5	Apht	5.043	89.032	22.858	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
44	44	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	288.	5	Apht	5.043	89.027	22.8576	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
45	45	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	289.	5	Apht	5.043	89.028	22.8576	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
46	46	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	290.	5	Apht	5.048	89.035	22.8595	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
47	47	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	291.	5	Apht	5.042	89.035	22.8595	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
48	48	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	292.	5	Apht	5.039	89.07	22.849	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
49	49	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	293.	5	Apht	5.037	89.07	22.849	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
50	50	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	294.	5	Apht	5.041	89.07	22.849	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
51	51	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	295.	5	Apht	5.041	89.07	22.849	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
52	52	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	296.	5	Apht	5.041	89.07	22.849	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
53	53	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	297.	5	Apht	5.041	89.07	22.849	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
54	54	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	298.	5	Apht	5.041	89.07	22.849	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
55	55	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	299.	5	Apht	-5.022	89.237	22.9128	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
56	56	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	300.	5	Apht	-5.023	89.216	22.9059	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139
57	57	e	13_V1_00066F7919-CAL	0	2014-00-14723.3	0.005940	301.	5	Apht	-5.024	89.204	22.904	3.8931	6.49359	4.92264	3.54801	6.17983	0.719139

Fig. 5 Browse your table.

In the following, our goal is now to identify and extract observations in a specific location (on the Imhotep region of comet 67P) and observed under specific illumination conditions.

### 3.1 Quick

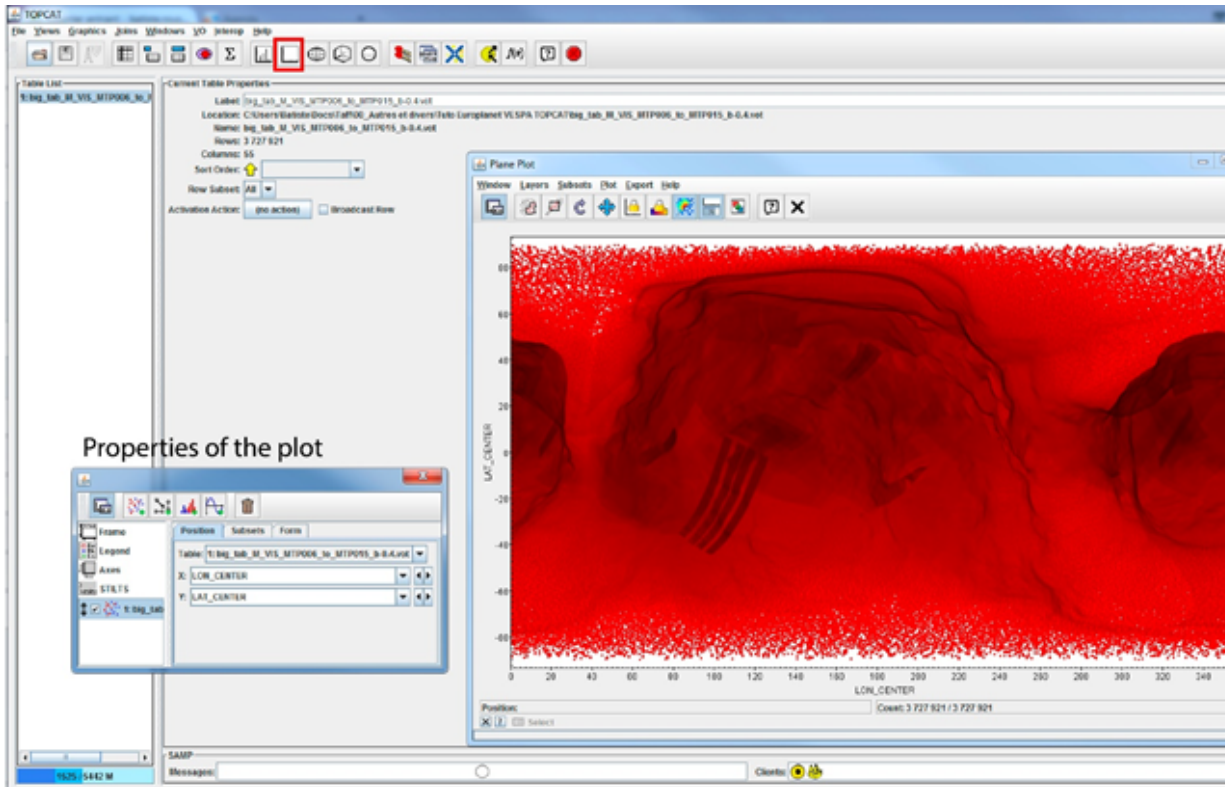


Fig. 6 Plotting longitude vs. latitude gives a...MAP ! By default, a single-color density representation is used.

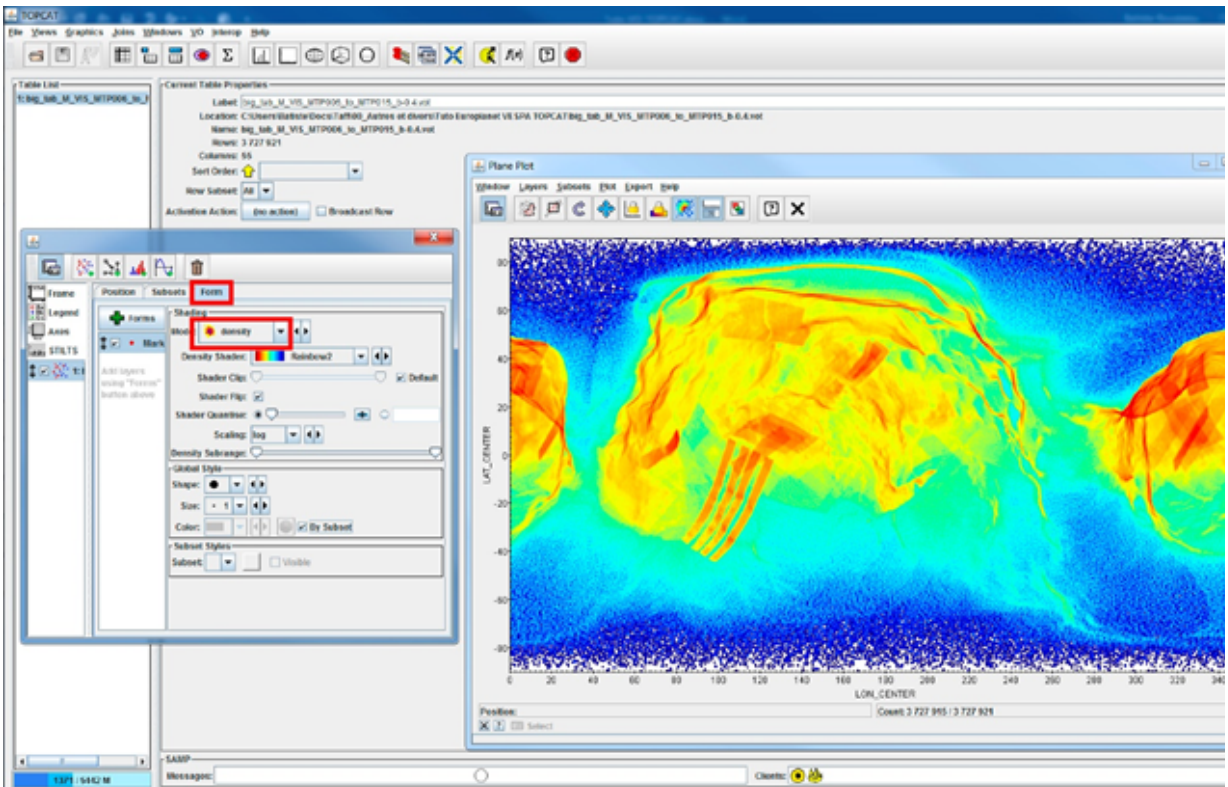


Fig. 7 Density map. This option is available in the form tab as well as other customization options.

Spatial information can be mapped with other projections using dedicated tools. Fig. 8 shows three map projections using the **sky plotting** option based on latitude and longitude. The projection type is available in the **Axes** settings (left menu of Sky plot window) + **Projection** tab, while the coordinate grid can be set in the **Grid** tab. In the example, the radius (from the 3D shape model) at each FoV center is plotted. Tools such as **3D plotting window using Cartesian coordinates** and **3D plotting window using spherical polar coordinates** (see top of the Fig. 8 for the access) allow you to use x/y/z coordinates (latitude, longitude and radius, respectively) to display a 3D representation.







Fig. 8 Different map projections of the "Sky plotting window". The radius is represented here as the auxiliary data.

### 3.2 Healpix maps

A powerful capacity of TOPCAT is to use the **Healpix** scheme to obtain a map with a homogeneous resolution / regular sampling. Healpix is a hierarchical tessellation of the sphere providing a grid of cells of the same size, available at various scales. Data points from our table are integrated in the Healpix cells, and different parameters can be displayed: sum of values, median or mean, min or max value, point density... TOPCAT is able to perform this entire task thanks to the **"Sky density"** form in the **Sky plotting window**.

As for the maps presented in Fig. 8, lon/lat coordinates are filled in with the columns "LON\_CENTER" and "LAT\_CENTER" (which correspond here to the name of our variables) to build the map. The Sky density is available from the **Sky plotting window** (step 1 in Fig. 9) in the **Form** tab. Uncheck or remove the Mark form, select **Sky density** in the **"Forms"** menu (step 2 in Fig. 9). A weight parameter can be chosen or left empty (step 3 in Fig. 9). In step 4, you can vary the size of the Healpix cells (the angular size is displayed at the bottom of this window, step 6). When a variable is chosen as weight (e.g. reflectance) the **"Combine"** option (step 5 in Fig. 9) allows you to display the sum, median, or other quantities computed into the Healpix cells. Finally, in step 6, it is possible to import/export the Healpix map into a table.

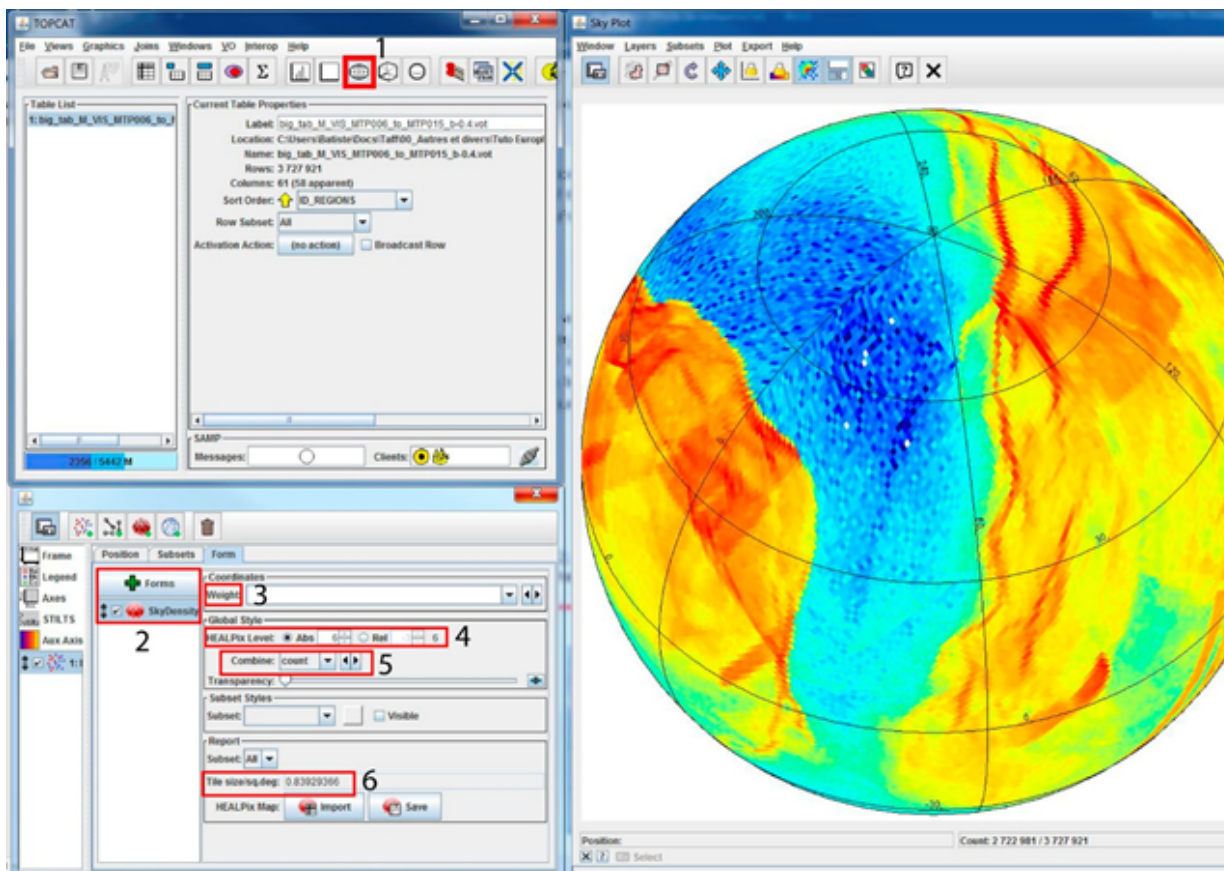


Fig. 9 An example of a density map using the Sky density function.

TOPCAT provides an alternative method to handle healpix maps (through the **"Add a new Healpix layer control to the stack"** tools). However, this method requires that the data are already averaged in the healpix cells, typically after a specific database query.

### 3.3 Filtering data

You now want to study separately the observations of the Imhotep region, which is approximately centered on longitude/latitude 0°. There are several ways to do this:

1. It is possible to create a subset graphically by enlarging the plot and defining a new subset based on visible points, thanks to the option **"Define new row subset containing only visible points"**. However, this is only relevant for rectangular regions, and not complex ones that are predefined from other parameters.
2. Arbitrary regions can be defined in freehand selection / lasso mode, thanks to the option **"Draw a region on the plot to define a new row subset"** (click, select, and click again). To make this operation easier, you can plot the identification number of each region (a parameter included in this table, see maps in Fig. 8) as auxiliary data. Regions will be colored independently. However, freehand selection is adapted to isolated groups of points, which is not the case here.

3. The easier way in this case is to use the region ID from the table to select only lines corresponding to Imhotep. Click Subset in the main TOPCAT window then click the green + icon (new subset), enter a subset name and expression `ID_REGIONS==13` (which stands for Imhotep).

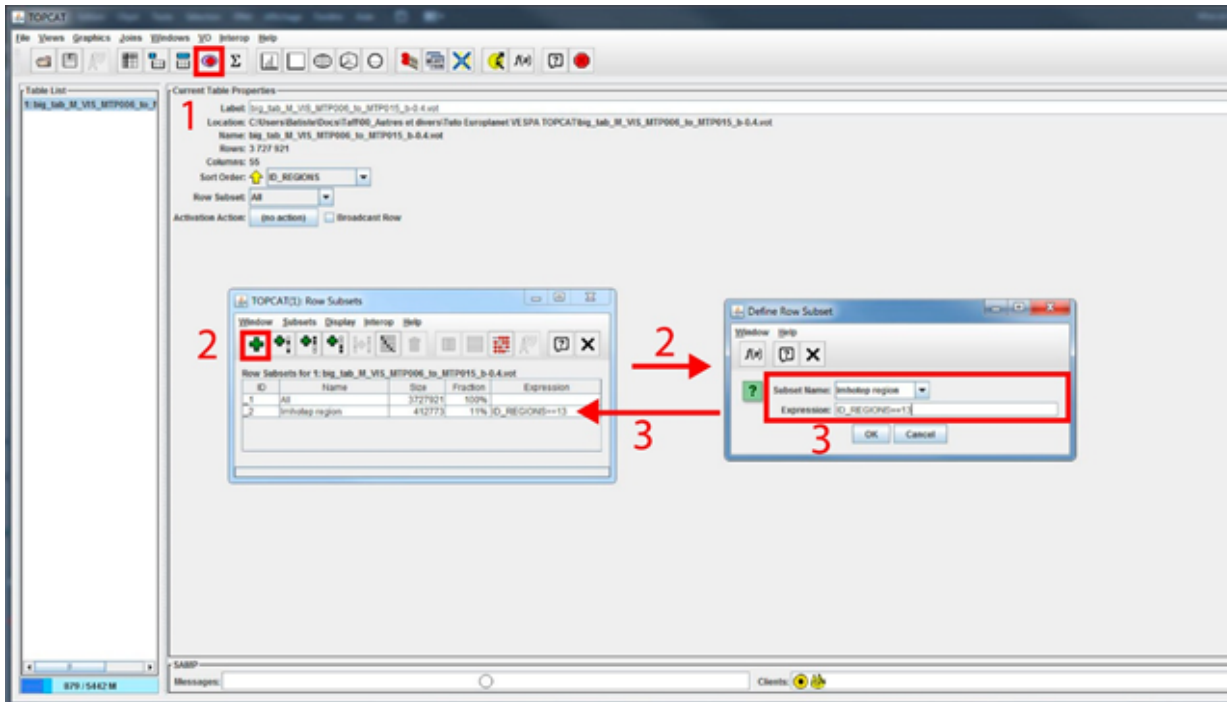


Fig. 10 How to define a subset, quickly and easily.

Once the subset is defined, you can highlight it in the table browser or display only this part of the data using the "Subset" item in the global TOPCAT menu, when the table window is selected (Fig. 11). It is also possible to create a new table containing this dataset by selecting the subset and exporting it in any format (Fig. 12). We can save this new table or the entire session to load it later.

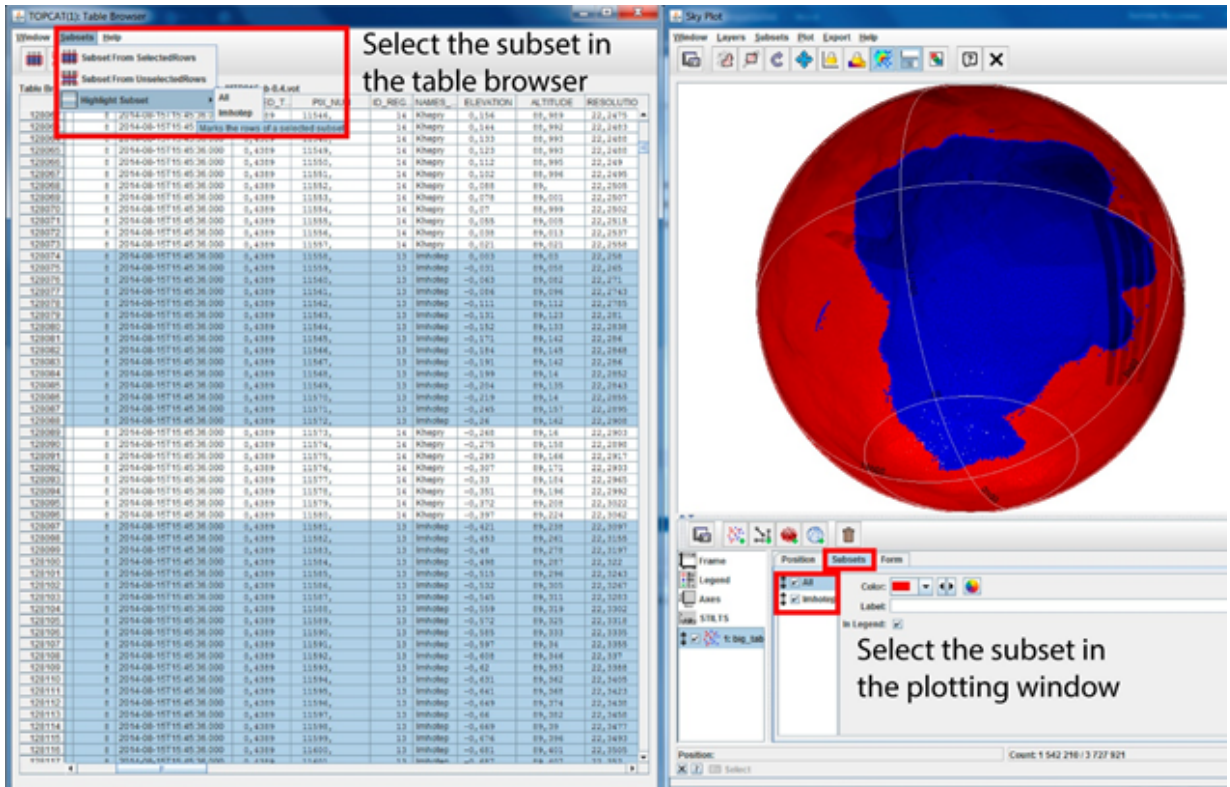


Fig. 11 After the creation of a new subset, it is possible to highlight it in the table browser (left) or plot it separately (right).

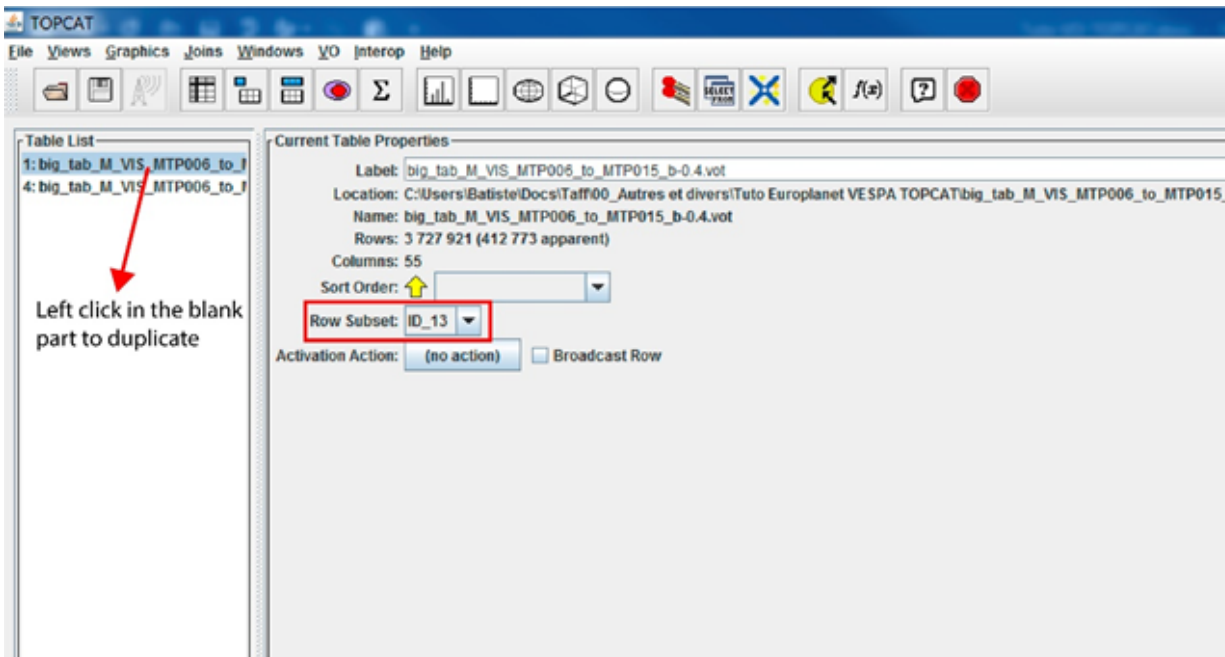


Fig. 12 Duplicate table or a subset by click and drag in the blank part.

You now want to get an idea of the illumination conditions in the selected subset. It is possible to plot **histograms** as shown in Fig. 13. Many options are available to explore the data: superposition of histogram, fitted curve... Here, the dataset ranges from 3.5° to 179° in incidence angle and from 0° to 129° in emergence angle. Another solution to obtain this information is to use the tool **Display statistics for each column** from the main window. From this window an appreciable number of quantities can be displayed (see Fig. 14).

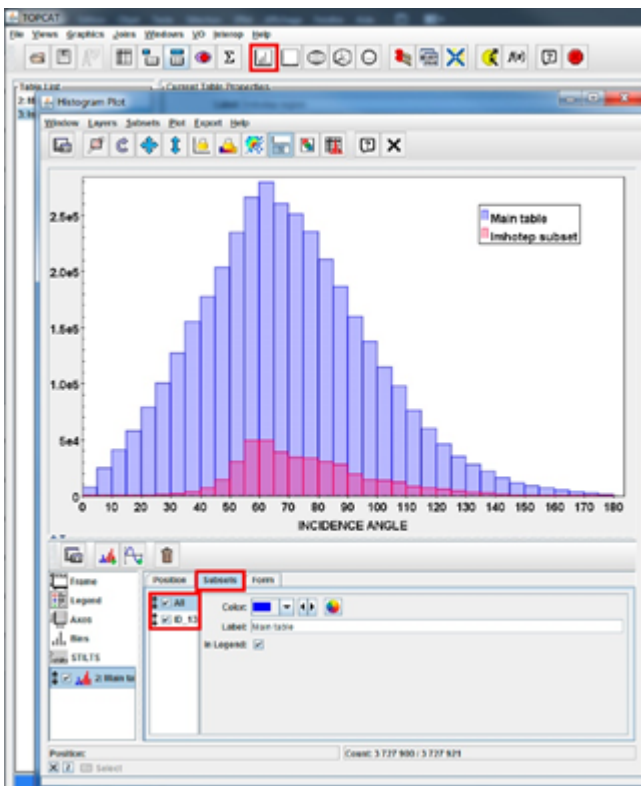


Fig. 13 Many options are available with the histogram and other plotting modes.



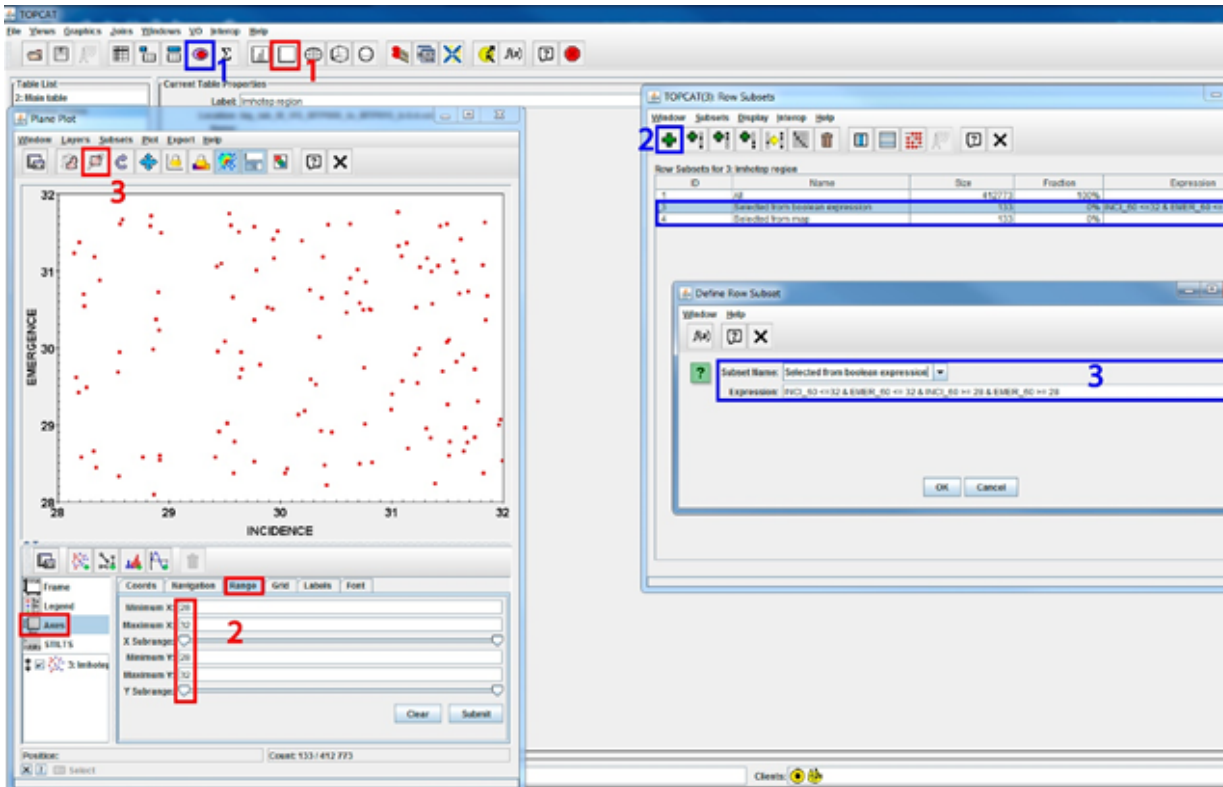


Fig. 15 Two ways of defining a similar subset.

When a specific dataset is defined, it is possible to save it or to broadcast the data to other VO tools. This method uses the SAMP protocol to communicate between the different applications. For example, data can be sent to the MATISE tool (<https://tools.asdc.asi.it/MatisseNoPermission.jsp>) which can project data on the 3D shape model of 67P; to Aladin, e. g., to access other mapping options; to CASSIS for spectral plots, etc.