

Welcome to DART



DART, Decision Analysis by Ranking Techniques, is a powerful yet easy-to-use tool for analysis of datasets based on the most up to date ranking theories. It offers the capability of importing any kind of dataset, and processing it with total and/or partial ranking procedures, providing both numerical and graphical outputs (including the Hasse diagram).

This help provides also a short introduction to the Ranking Theory, explaining the mathematical meaning of the procedures used by DART software.

DART has been developed by [TALETE Srl](#), funded by the [EC Joint Research Centre](#) . It is released under the GNU [General Purpose License](#) which allows you to redistribute freely this software and gives free access to its code (for restriction about its modification and other details, please refer to the license).

Feel free to contact DART developers at info@talete.mi.it for any communication about this software.

Quick start



The first step to be done is loading data; this can be done both by importing a plain text file containing a dataset (you can access the import procedure by clicking [FILES - IMPORT](#)) or by loading a pre-existing DART file with extension .drt ([FILES - LOAD](#))

Once you have loaded a dataset, you can access to the the DATA menu, which contains two voices:

- [VIEW](#) allows you to view and analyze the loaded dataset
- [SETUP](#) open a window with four tabs, in which you have to give a setup for the loaded dataset; trough this form, you can:
 - Choose which variable are selected and which should be ignored, and fill missing values if present;
 - Choose which object are selected and which should be ignored;
 - Select a transform function and set a weight for each selected variable.

You can also access the PREPROCESSING menu, which contains several methods that can be applied to your dataset to perform a better ranking analysis:

- [SIGNIFICANT DIGITS](#) allows you to round values of your dataset
- [BINS PARTITION](#) allows you to partiotionate your dataset into a given number of bins
- [PCA](#) performs a complete Principal Component Analysis on your dataset, then allowing you to use a wanted number of principal components as new variables for your dataset
- [K-MEANS CLUSTERING](#) performs a clustering of the samples via K-Means algorithm

If you perform one or more of these preprocessing methods, you will be required to perform the data setup again, as the dataset could be changed and the previous setup suitable no more.

After you ended the setup step, you can access the [CALCULATE](#) menu, which contains three voices:

- [TOTAL RANKING](#) opens the form containing all the calculations about total ranking methods, summarized in the first grid which contains the ranking functions (desirability, utility etc.) for the selected object. In the other tabs you can access several graphs giving you different views of the results.
- [PARTIAL RANKING](#) opens the form containing all the calculations about partial ranking methods, i.e. informations about the hasse diagram such as some indices, the levels structure and the hasse matrix.
- [HASSE DIAGRAM](#) opens the hasse diagram chart; when this graph is the active windows, you can select several visualization options from the [DIAGRAM](#) menu.

Menu overview



The menu of DART software is outlined below:

- [FILES](#): Access to import and load/save functions.
- [DATA](#): Access to setup of the loaded dataset and visualization of the dataset (this menu

is active only if a dataset is loaded).

- [PREPROCESSING](#): Access to different preprocessing techniques
- [CALCULATE](#): Access to calculations for total and partial ranking method, and visualization of the Hasse diagram (this menu is active only if a dataset is loaded and a valid setup has been performed).
- [DIAGRAM](#): Visualization and manipulation options for the Hasse diagram (this menu is active only if a Hasse diagram form is active).
- WINDOW: List of active forms.
- ?: Access to DART help and about form.

Introduction



In the FILES menu you can access these voices:

[IMPORT](#): opens the import window

LOAD: loads a .drt file

SAVE: saves a .drt file

RECENT FILES: you will find the last five .drt files loaded, by clicking on a file name you will load it.

EXIT: closes DART

IMPORT menu



Through the importing procedure, it is possible to load any dataset in plain text format; datasets are intended to be a matrix of $n \times m$ numerical data, interpreted as n rows (objects/samples) for m columns (variables). The procedure is made of two main steps; after choosing the file to be imported, the first form is showed:

Import - Step 1/2

Delimiter: Tab Comma Space Semicolon Other:

Use multiple delimiter as one

Cancel

Next

This is a preview of the first 20 rows

	Col 1	Col 2	Col 3	Col 4	Col 5	Col 6	Col 7
Row 1	Molecole Beam						
Row 2	17	2	0	6	0		
Row 3	SUBSTANCE	ID	LOG(1/EC01	LOG(1/EC10	LOG(1/EC20	LOG(1/EC50	LOG(1/EC100
Row 4	Bitertanol	1	.851953239	.161745986	.058678894	.391603232	.639
Row 5	Cyproconazole	2	.729412897	.137212139	.948086655	.662436157	.450
Row 6	Diclobutrazol	3	.185857564	0.3628107	.099962105	.297037971	.592
Row 7	Difenoconazole	4	1.31173516	.769060589	.595751809	.333990475	.139
Row 8	Fenbuconazole	5	.301922225	.697904215	.505004768	.213654171	0.00
Row 9	Flusilazole	6	.306366833	.699002256	.505034048	0.21206922	.005
Row 10	Flutriafol	7	.089970332	.390649497	.544140449	.775969421	.948
Row 11	Hexaconazole	8	.424258082	.931407732	.774010836	.536282421	.359
Row 12	Myclobutanil	9	0.54792779	.035214157	.221446616	.502727558	0.71

In this form, DART tries to automatically understand which is the delimiter used in the selected files, and shows a preview of the dataset; you can manually select the desired delimiter between the most common ones (tab, comma, space, semicolon), or select the voice "other" and type the delimiter character used in the file. By checking the "use multiple delimiter as one" voice, DART will not consider multiple delimiters.

By clicking on the "next" button, the second form is showed:

Import - Step 2/2

Primary ID col:

Label row:

Missing value:

Preview of first 20 rows

Cancel

Next

	Col 1	Col 2	Col 3	Col 4	Col 5	Col 6	Col 7
Row 1	Molecole Be:						
Row 2	17	2	0	6	0		
Row 3	SUBSTANCE	ID	LOG(1/EC01	LOG(1/EC10	LOG(1/EC20	LOG(1/EC50	LOG(1/EC100
Row 4	Bitertanol	1	.851953239	.161745986	.058678894	.391603232	.639
Row 5	Cyproconazi	2	.729412897	.137212139	.948086655	.662436157	.450
Row 6	Diclobutrazo	3	.185857564	0.3628107	.099962105	.297037971	.592
Row 7	Difenoconaz	4	1.31173516	.769060589	.595751809	.333990475	.139
Row 8	Fenbuconaz	5	.301922225	.697904215	.505004768	.213654171	0.00
Row 9	Flusilazole	6	.306366833	.699002256	.505034048	0.21206922	.005
Row 10	Flutriafol	7	.089970332	.390649497	.544140449	.775969421	.948
Row 11	Hexaconazo	8	.424258082	.931407732	.774010836	.536282421	.359
Row 12	Myclobutanil	9	0.54792779	.035214157	.221446616	.502727558	0.71

Here, DART tries to automatically understand what is the ID column (i.e. the column containing the names of the objects/samples) and the label row (i.e. the row containing the names of the variables); the selected ID column and label row are highlighted in green. By

checking the "primary ID col" or "label row" voices, you are allowed to insert a manual value for the column and/or row. In the "missing value" the value used to point out missing values must be indicated; by default, the value is set to -999 (a standard value, used by most softwares).

By clicking on the "next" button, a short summary of the imported dataset is shown, and the procedure is over.

Introduction



In the DATA menu you can access these voices:

- [VIEW](#): view the current dataset.
- [SETUP](#): setup of the dataset.
- RESTORE ORIGINAL DATA: restores the dataset as it was before any preprocessing.

VIEW menu



In the view form you can examine the active dataset; for each object/sample (i.e. for each row) it is reported its id number, its name, the current status (if it is included or not) and the values for each variable. By clicking on the "var" or "obj" icons, it is shown a profile window, in which you have a graphical profiling for each variable or object.

ID	SUBSTANCE	Status	ID	LOG(1/EC01)	LOG(1/EC10)	LOG(1/EC20)	LOG(1/E
1	Bitertanol	included	1	0.8519533	0.1617460	-0.0586789	-0.39
2	Cyproconazole	included	2	1.7294130	1.1372120	0.9480867	0.66
3	Diclobutrazol	included	3	1.1858580	0.3628107	0.0999621	-0.29
4	Difenoconazole	included	4	1.3117350	0.7690606	0.5957518	0.33
5	Fenbuconazole	included	5	1.3019220	0.6979042	0.5050048	0.21
6	Flusilazole	included	6	1.3063670	0.6990023	0.5050340	0.21
7	Flutriafol	included	7	0.0899703	-0.3906495	-0.5441405	-0.77
8	Hexaconazole	included	8	1.4242580	0.9314077	0.7740108	0.53
9	Myclobutanil	included	9	0.5479278	-0.0352142	-0.2214466	-0.50
10	Padlobutrazol	included	10	-0.4980606	-0.7627438	-0.8472731	-0.97
11	Penconazole	included	11	0.7050738	0.2679431	0.1283408	-0.08
12	Prochloraz	included	12	1.9889250	1.5388980	1.3951770	1.17
13	Propiconazole	included	13	0.6501225	0.2461004	0.1170717	-0.07
14	Tebuconazole	included	14	0.9218000	0.3430213	0.1581824	-0.12
15	Tetraconazole	included	15	0.6317981	0.1142230	-0.0510700	-0.30
16	Triadimefon	included	16	0.5095825	-0.0116151	-0.1780651	-0.42
17	Triadimenol	included	17	0.3691583	-0.1283672	-0.2872572	-0.52

No. of objects: 17 No. of variables: 7

SETUP menu



The setup is the main step to be done after importing a dataset, and unless it is done you can not access the rankings methods.

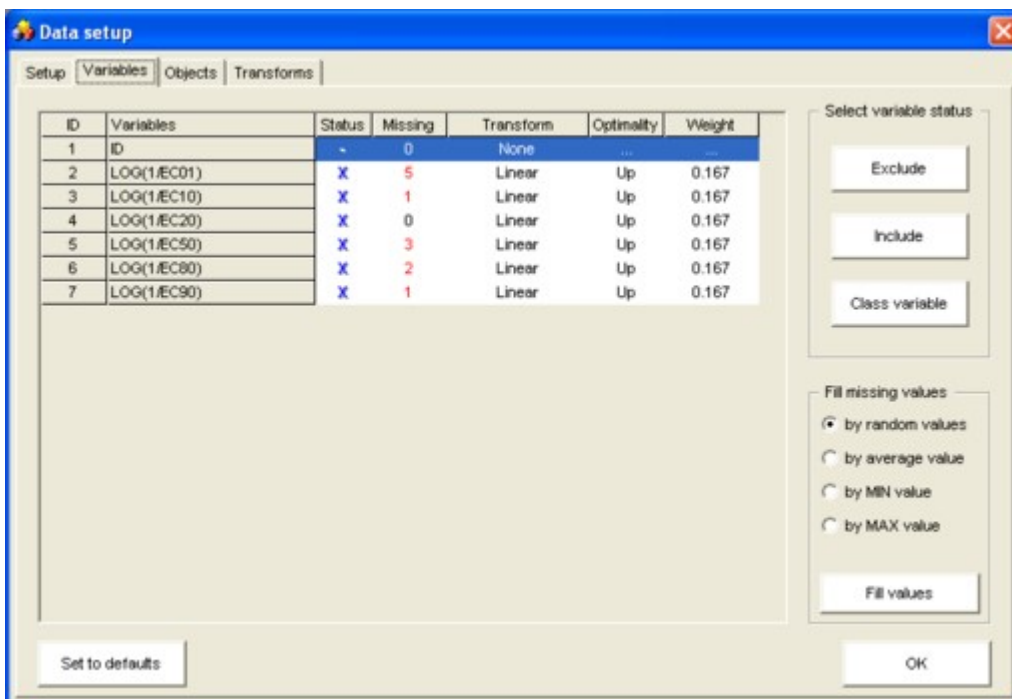
The setup window is constituted by four tabs:

- [Setup](#)
- [Variables](#)
- [Objects](#)
- [Transform](#)

[SETUP Tab](#)

The first tab reports a summary of the setup, i.e. how many variables and objects are selected and how many transform functions are set. By clicking on the "Set to defaults" button, all values are set back to defaults (all variables and objects are selected, all variables are set with linear transform function between maximum and minimum values).

[VARIABLES Tab](#)

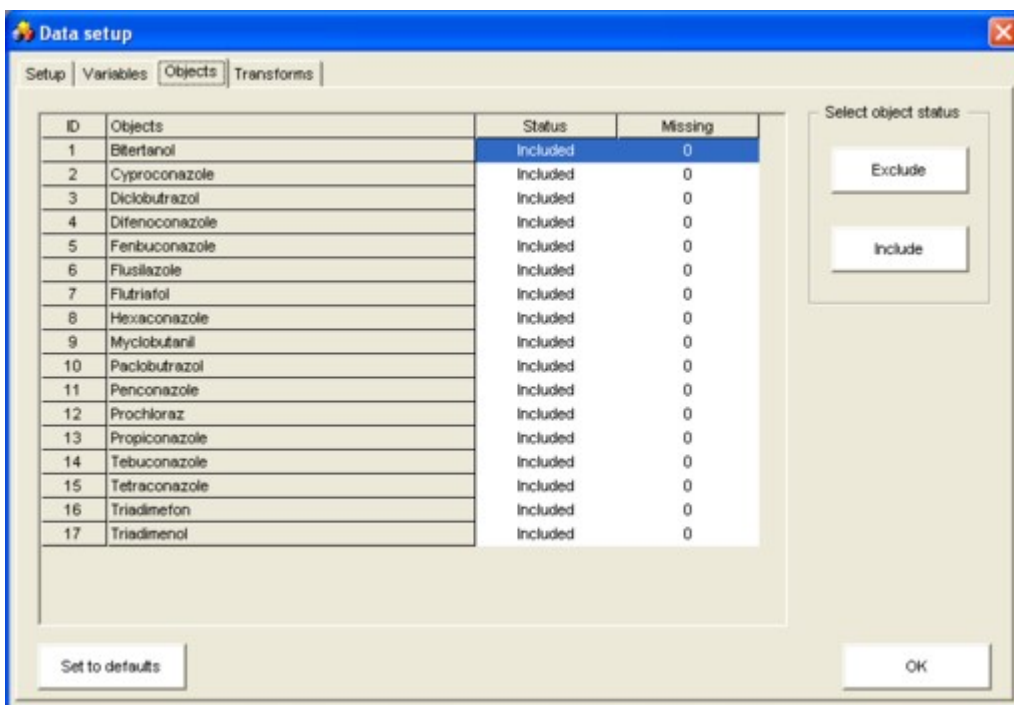


The second tab shows all variables found in the dataset and some informations for each of them (current status, number of objects with missing values for that variable, transform function selected, optimality of the transform function, weight). By clicking the "exclude" button, the variable currently highlighted in the grid is excluded (its status is shown as "-"); by clicking the "include" button, the variable currently highlighted in the grid is included (its

status is shown as "X"); by clicking the "class variable" button, the variable currently highlighted in the grid is set as the class variable (its status is shown as "C"). Note that you can set only one class variable; a class variable must contain only integer values, as its meaning is that a given object is associated with a class indicated by an arbitrary integer number.

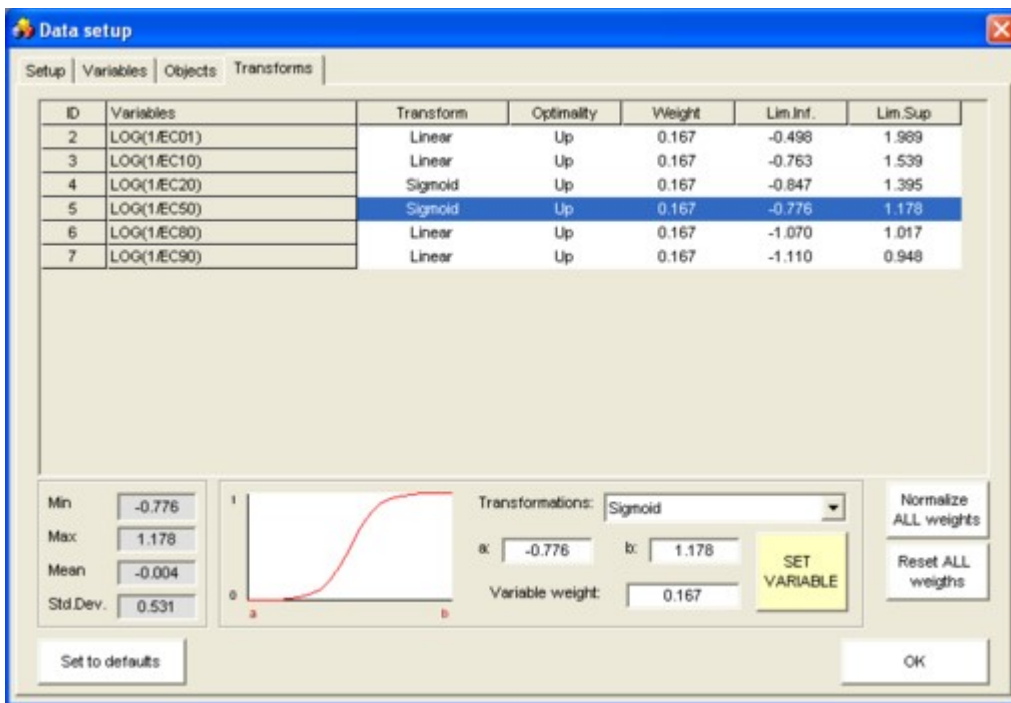
If some missing values are found, the "Fill missing values" box is shown; you can then select a method to replace missing values (by a random value or by the average/minimum/maximum value of the variable) and then proceed by clicking the "Fill values" button. Note that DART can handle dataset with missing values, so it is not necessary to fill them in order to go further; objects/samples having missing values will be indicated attaching a * character at the end of their names, to remind that all calculations performed on them can be approximative due to missing values.

[OBJECTS Tab](#)



The third tab shows all objects found in the dataset and some informations for each of them (current status, number of variables with missing values for that object). By clicking the "exclude" button, the object currently highlighted in the grid is excluded (its status is shown as "Excluded"); by clicking the "include" button, the object currently highlighted in the grid is included (its status is shown as "Included").

[TRANSFORM Tab](#)



The fourth tab shows only the previously selected variables and some informations about the transform functions for each of them (type of transform function, optimality, weight of the variable, minimum and maximum values). In the lower left panel, some statistical data are reported for the highlighted variable (minimum and maximum values, mean value, standard deviation value).

In the middle panel, it is possible to set the transform function, by selecting a transform type (to help in the choice, the behaviour of each selected transform function is shown in a small graph) and its limits "a" and "b" (i.e. the boundaries of the transform function, beyond which 0 and 1 values are always set; by default, they are set to the minimum and maximum values of the variable), and to set the variable weight. By clicking the "Set variable" button, the highlighted variable is set with the desired settings.

By clicking the "Normalize all weights", all weights will be normalized (i.e. their total sum is equal to one). Note that this procedure will be done anyway as soon as you change tab or you close the setup form (weights are always normalized). By clicking the "Reset all weights", all weights will be set as equal one to each other.

Introduction



In the PREPROCESSING menu you can access these voices, divided into two sub-menus:

On Variables:

- [SIGNIFICANT DIGITS](#): changes the number of significant decimal digits for each variable.
- [BINS PARTITION](#): performs a bins partition on variables.

- [PCA](#): performs a Principal Component Analysis.

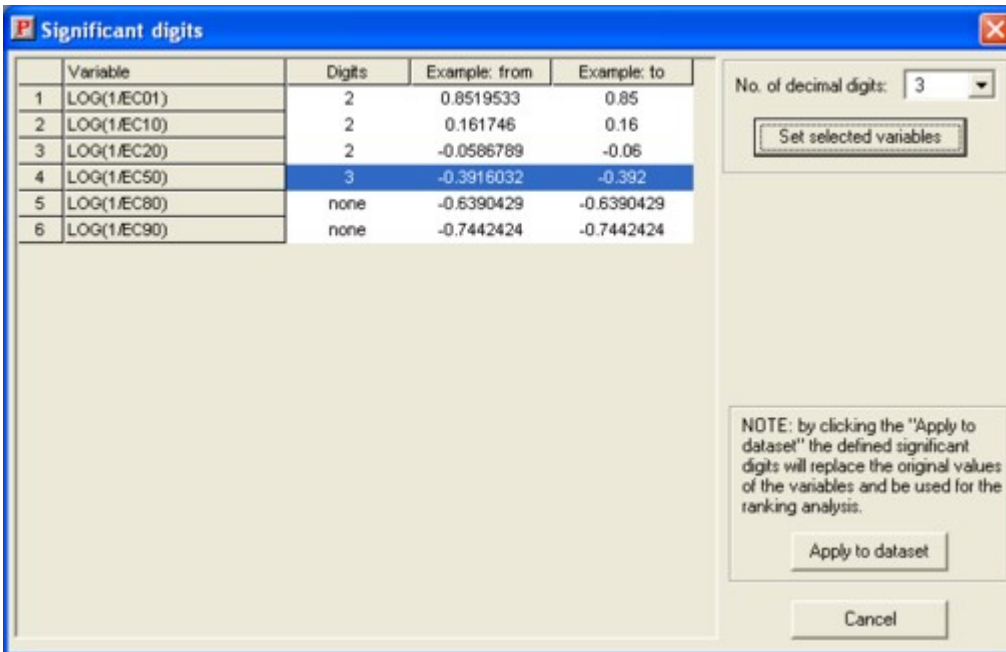
On Objects:

- [K-MEANS CLUSTERING](#): performs clustering using K-Means algorithm.

SIGNIFICANT DIGITS menu

In this form you can set the number of decimal digits for each variable. In the grid it is shown the number of digits set for each variable, followed by an example of the rounding that will be applied with the selected number of digits.

By clicking the "Apply to dataset button", the selected rounding will be applied to the current dataset.



	Variable	Digits	Example: from	Example: to
1	LOG(1/EC01)	2	0.8519533	0.85
2	LOG(1/EC10)	2	0.161746	0.16
3	LOG(1/EC20)	2	-0.0586789	-0.06
4	LOG(1/EC50)	3	-0.3916032	-0.392
5	LOG(1/EC80)	none	-0.6390429	-0.6390429
6	LOG(1/EC90)	none	-0.7442424	-0.7442424

No. of decimal digits: 3

Set selected variables

NOTE: by clicking the "Apply to dataset" the defined significant digits will replace the original values of the variables and be used for the ranking analysis.

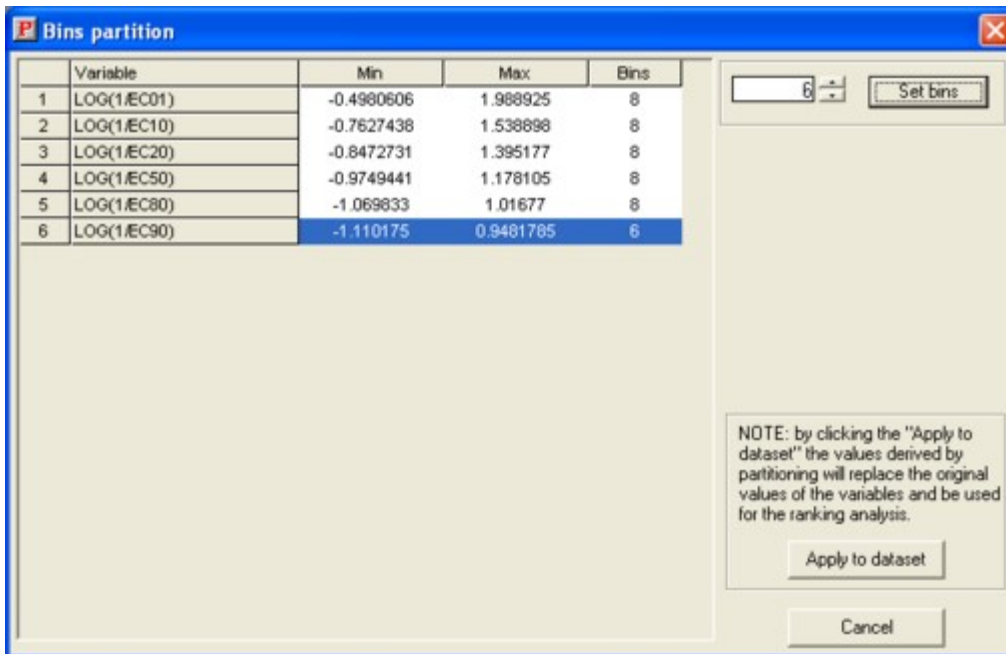
Apply to dataset

Cancel

BINS PARTITION menu

In this form you can set the number of bins to partition each variable. In the grid it is shown the number of selected bins, and the minimum and maximum values for each variable.

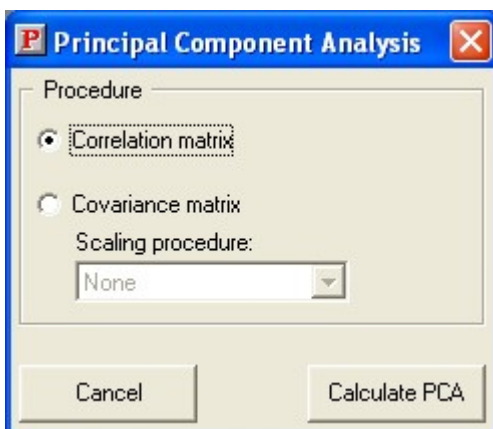
By clicking the "Apply to dataset button", the selected partitioning will be applied to the current dataset.



PCA menu



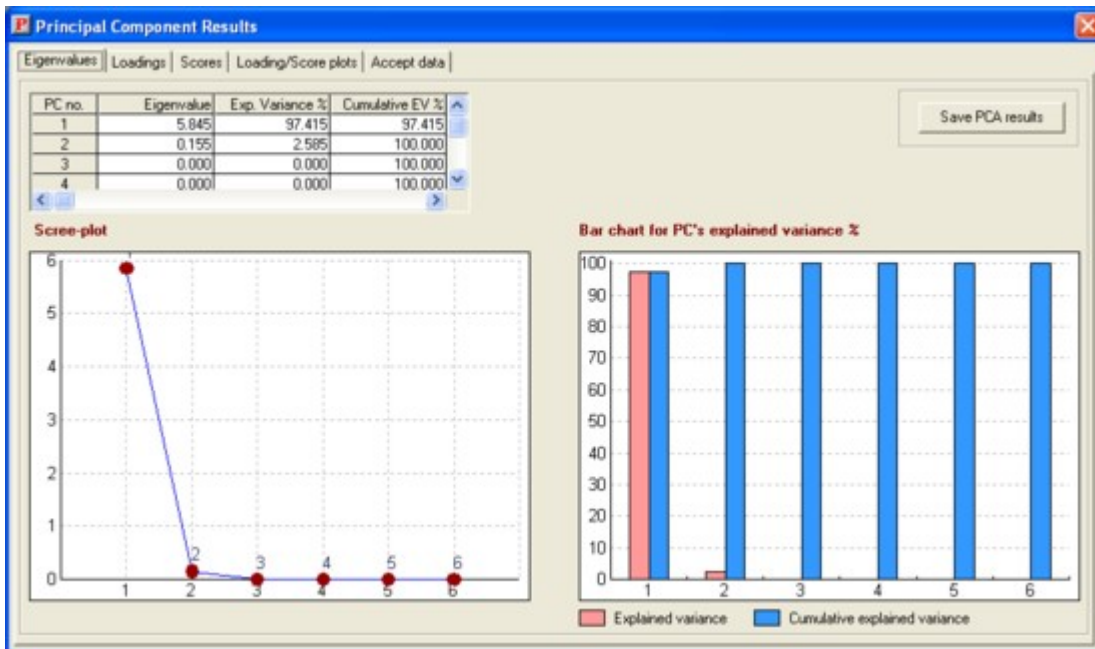
The first form shown is the setup of the Principal Component Analysis, in which you can choose whether to perform it on the correlation matrix or on the covariance matrix (in this latter case, you can choose the scaling procedure).



By clicking the "Calculate PCA" button, the analysis is performed and the main PCA form is shown. This form is constituted by five tabs:

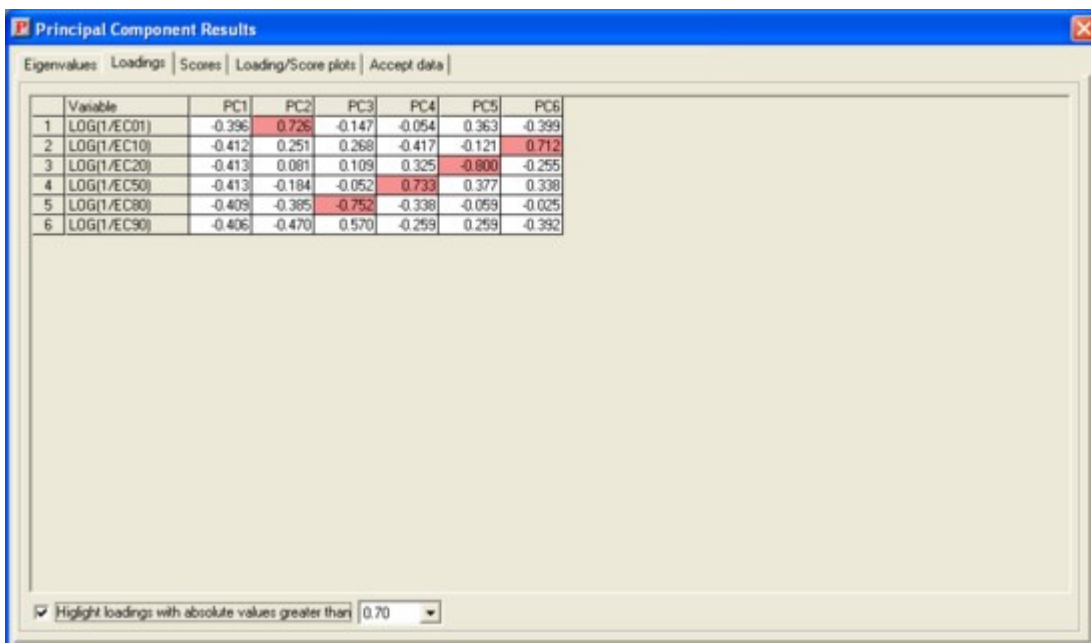
- [Eigenvalues](#)
 - [Loadings](#)
 - [Scores](#)
 - [Loading/score plot](#)
 - [Accept data](#)
-

[EIGENVALUES Tab](#)



The first tab reports all the information about the eigenvalues and the explained variance of all components. In the upper grid are reported the values of eigenvalues, their explained variance and their cumulative explained variance. In the lower part it is reported the scree plot (with the eigenvalues value of each component) and the bar chart with the explained variance and cumulative explained variance of each component. By clicking the "Save PCA results" it is possible to export the results of the analysis in plain text format.

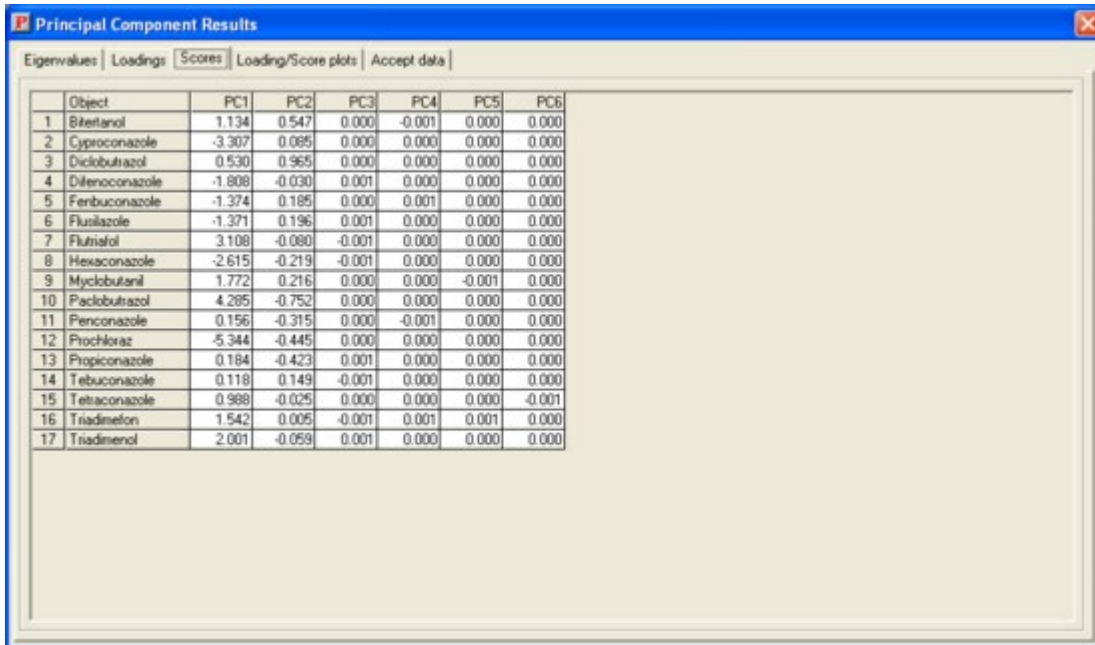
[LOADINGS Tab](#)



The second tab reports the values of the loadings of the variables in a grid. It is possible to highlight loadings with absolute values greater than a desired value, by using the checkbox

in the lower part of the panel.

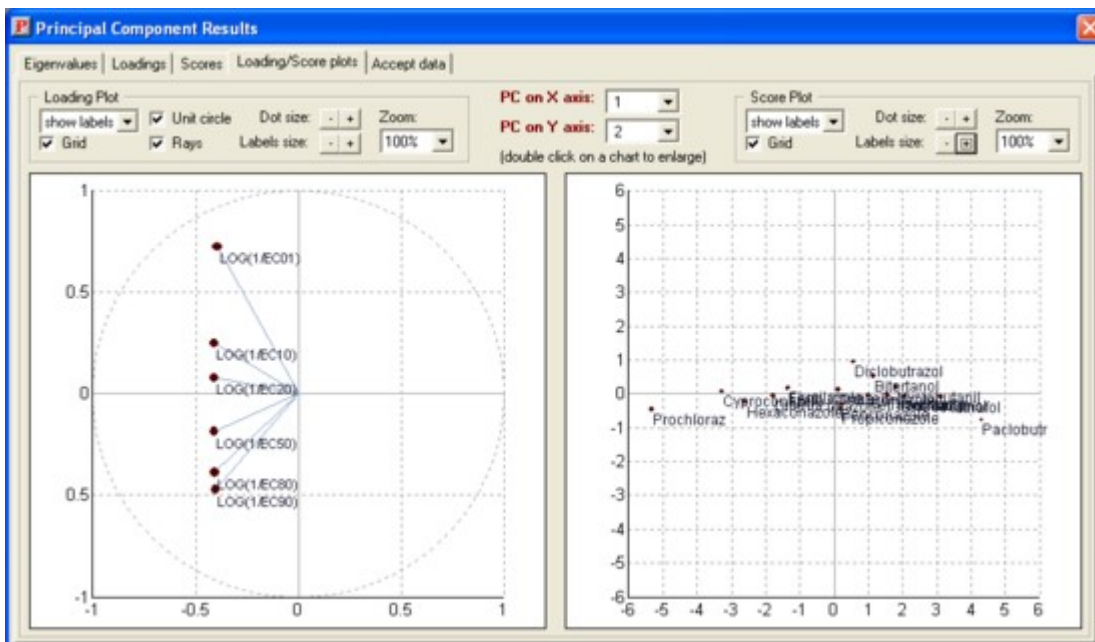
SCORES Tab



Object	PC1	PC2	PC3	PC4	PC5	PC6
1 Bifentanol	1.134	0.547	0.000	-0.001	0.000	0.000
2 Cyproconazole	-3.307	0.085	0.000	0.000	0.000	0.000
3 Dicloubutrazol	0.530	0.965	0.000	0.000	0.000	0.000
4 Difenoconazole	-1.808	-0.030	0.001	0.000	0.000	0.000
5 Fenbuconazole	-1.374	0.185	0.000	0.001	0.000	0.000
6 Flutriazol	-1.371	0.196	0.001	0.000	0.000	0.000
7 Flutriazol	3.108	-0.080	-0.001	0.000	0.000	0.000
8 Hexaconazole	-2.615	-0.219	-0.001	0.000	0.000	0.000
9 Myclobutanil	1.772	0.216	0.000	0.000	-0.001	0.000
10 Paclobutrazol	4.285	-0.752	0.000	0.000	0.000	0.000
11 Penconazole	0.156	-0.315	0.000	-0.001	0.000	0.000
12 Prochloraz	-5.344	-0.445	0.000	0.000	0.000	0.000
13 Propiconazole	0.184	-0.423	0.001	0.000	0.000	0.000
14 Tebuconazole	0.118	0.149	-0.001	0.000	0.000	0.000
15 Tetraconazole	0.988	-0.025	0.000	0.000	0.000	-0.001
16 Triadimefon	1.542	0.005	-0.001	0.001	0.001	0.000
17 Triadimenol	2.001	-0.059	0.001	0.000	0.000	0.000

The third tab reports all values of the scores of the objects in a grid.

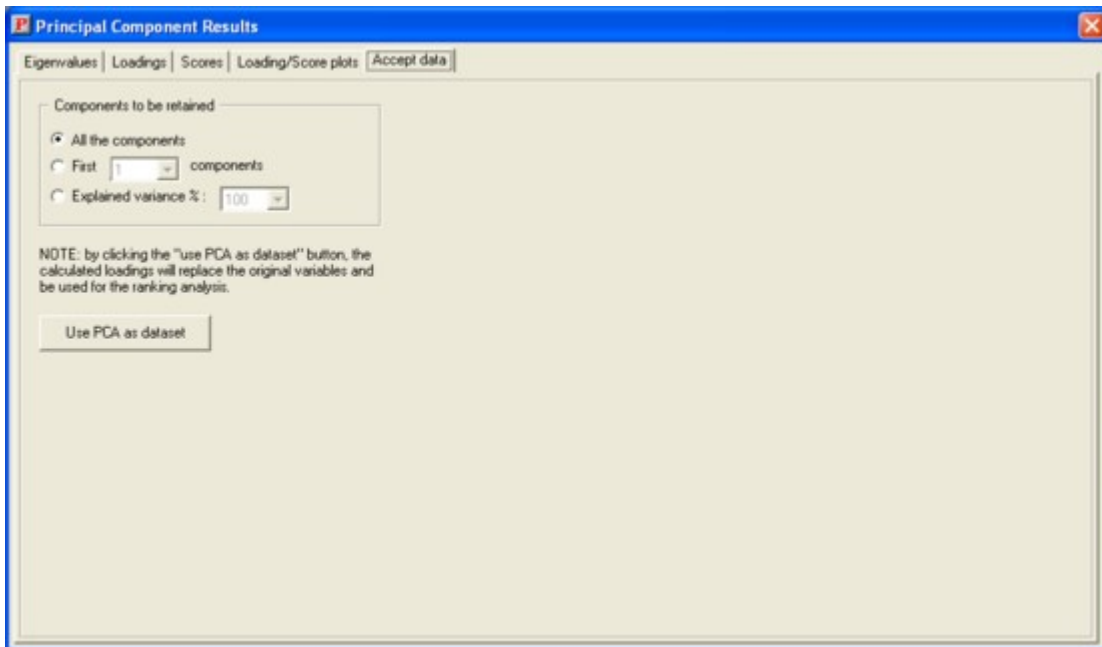
LOADING/SCORE PLOT Tab



The fourth tab reports the loading and the score plots. It is possible to select what components use for the charts. For both charts it is possible to change the visualization options in the upper panels. By double clicking on a chart, the chart itself it is opened in a

new form, making possible to enlarge it and obtaining a better visualization.

[ACCEPT DATA Tab](#)

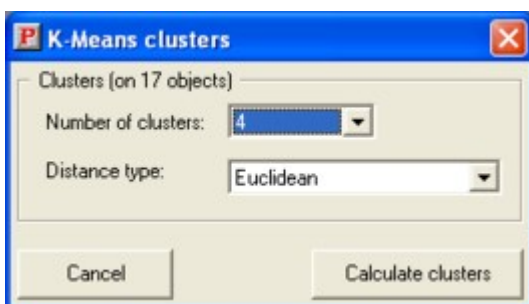


In the fifth tab it is possible to select how many components to use as dataset. By clicking the "Use PCA as dataset", the new dataset will be constituted by the original objects projected in the selected number of principal components.

K-MEANS CLUSTERING menu



The first form shown is the setup of the K-Means clustering, in which you can choose the number of clusters and the method for the calculation of distances.



By clicking the "Calculate clusters" button, the K-Means algorithm is performed and the main form is shown. In this form are reported the centroids for each cluster. If the "Show only centroids" checkbox is not checked, in the grid are also reported all the objects belonging to each cluster just after the representing centroid. By clicking the "Save K-means results" button it is possible to export the results of the clustering in plain text format. By clicking the "Use as dataset" button the new dataset will be constituted only by

the centroids.

ID	Name	LOG(1/EC01)	LOG(1/EC10)	LOG(1/EC20)	LOG(1/EC50)	LOG(1/EC80)
1	Centroid 1 (6 objs)	1.510	0.962	0.787	0.523	0.326
[2]	Cyproconazole	1.729	1.137	0.948	0.662	0.450
[4]	Difenoconazole	1.312	0.769	0.596	0.334	0.139
[5]	Fenbuconazole	1.302	0.698	0.505	0.214	-0.003
[6]	Flusilazole	1.306	0.699	0.505	0.212	-0.006
[8]	Hexaconazole	1.424	0.931	0.774	0.536	0.360
[12]	Prochloraz	1.989	1.539	1.395	1.178	1.017
2	Centroid 2 (4 objs)	0.127	-0.329	-0.475	-0.695	-0.859
[7]	Flutriafol	0.090	-0.391	-0.544	-0.776	-0.948
[9]	Myclobutanil	0.548	-0.035	-0.221	-0.503	-0.712
[10]	Paclobutrazol	-0.498	-0.763	-0.847	-0.975	-1.070
[17]	Triadimenol	0.369	-0.128	-0.287	-0.527	-0.706
3	Centroid 3 (2 objs)	1.019	0.263	0.021	-0.345	-0.616
[1]	Bitertanol	0.852	0.162	-0.059	-0.392	-0.639
[3]	Didlobutrazol	1.186	0.363	0.100	-0.297	-0.592
4	Centroid 4 (5 objs)	0.684	0.192	0.035	-0.202	-0.378
[11]	Penconazole	0.705	0.268	0.128	-0.083	-0.239
[13]	Propiconazole	0.650	0.246	0.117	-0.078	-0.223
[14]	Tebuconazole	0.922	0.343	0.158	-0.121	-0.328
[15]	Tetraconazole	0.632	0.114	-0.051	-0.301	-0.486
[16]	Triadimefon	0.510	-0.012	-0.178	-0.429	-0.616

Introduction



In the CALCULATE menu you can access these voices:

- [TOTAL RANKING](#): opens the form for total ranking calculations.
- [PARTIAL RANKING](#): opens the form for partial ranking calculations.
- [HASSE DIAGRAM](#): opens the Hasse diagram form.

TOTAL RANKING menu

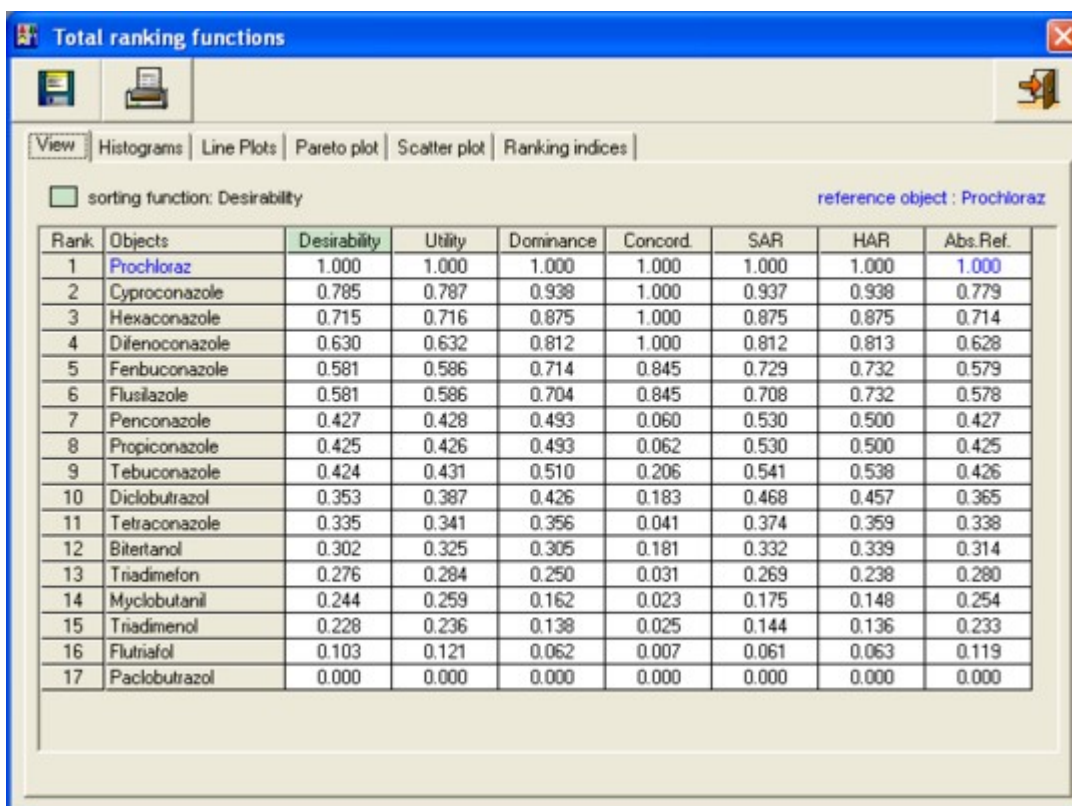


In the Total Ranking form, all calculations related to total ranking are reported, divided in six tabs. By clicking on the save or print icons in the upper part of the form, you can save to a txt file or send to printer the information you will select in a dialog box.

The total ranking window is constituted by six tabs:

- [View](#)
- [Line plot](#)
- [Histograms](#)
- [Pareto plot](#)
- [Scatter plot](#)
- [Ranking indices](#)

[VIEW Tab](#)



In the first tab, all the resulting calculations for total ranking functions (desirability, utility, dominance, concordance, SAR, HAR, absolute reference) are reported in the grid. By clicking on a column header, results will be sorted by the selected function (which will be highlighted in green); by clicking on a row header, the selected object will be set as the absolute reference (which will be highlighted in purple) and the grid will be sorted by the absolute reference function.

[LINE PLOT Tab](#)

In the second tab, the line plots for all the six ranking functions (or it is possible to select only some of them, by clicking on the checkbox in the legend) is shown; results are ordered by the ranking function selected on the first tab (which is also indicated in the footer of the plot, as "Reference function").

[HISTOGRAM Tab](#)

In the third tab, the histograms for all the six ranking functions (or it is possible to select only some of them, by clicking on the checkbox in the legend) it is shown; results are ordered by the original objects order.

[PARETO PLOT Tab](#)

In the fourth tab, the pareto plot for a selected function is shown; results are ordered by the ranking function selected on the first tab (which is also indicated in the footer of the plot, as "Reference function").

[SCATTER PLOT Tab](#)

In the fifth tab, a 2D scatter plot, for two functions that can be selected using the combo boxes on the right panel, is shown.

[RANKING INDICES Tab](#)

In the sixth tab, some mathematical indices are given for each ranking function (number of levels, stability StR index, degeneracy k and D indices, YDbyR index).

PARTIAL RANKING menu

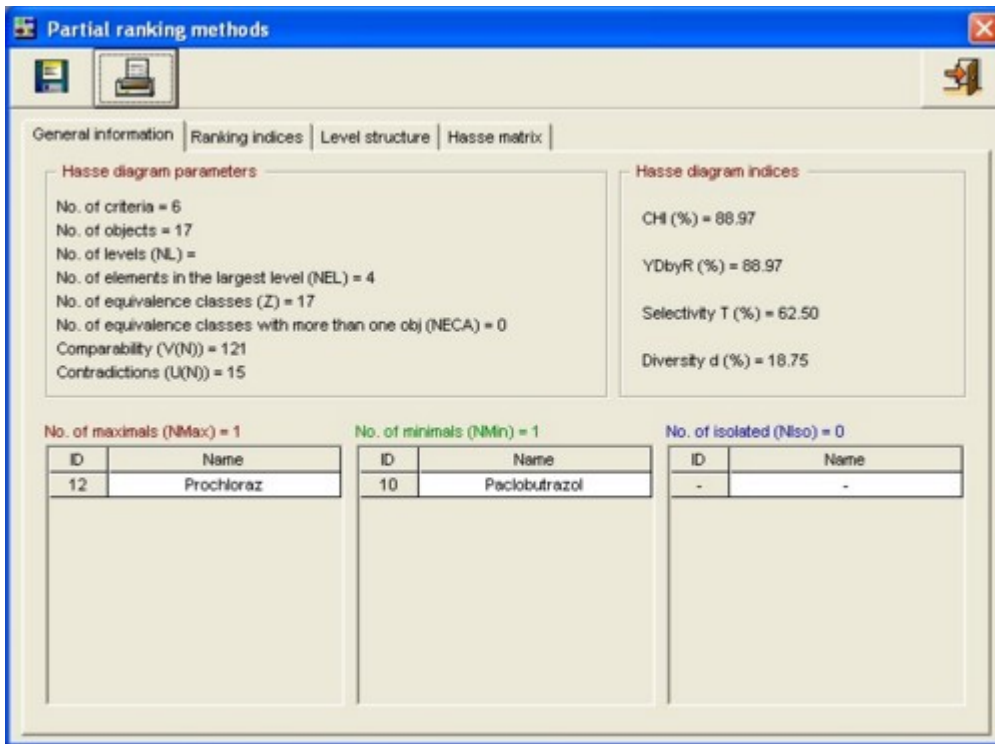


In the Partial Ranking form, all calculations related to partial ranking (i. e. the calculations which lead to the Hasse diagram) are reported, divided in four tabs. By clicking on the save or print icons in the upper part of the form, you can save to a txt file or send to printer the information you will select in a dialog box.

The partial ranking window is constituted by four tabs:

- [General information](#)
- [Ranking indices](#)
- [Level structure](#)
- [Hasse matrix](#)

[GENERAL INFORMATION Tab](#)



In the first tab, all the parameter and indices for the Hasse diagram are reported, together with the number and the complete list of minimal, maximal and isolated objects in the diagram.

[RANKING INDICES Tab](#)

In the second tab, stability and degeneracy indices for the diagram are reported, together with two histograms showing in a graphical way these values.

[LEVEL STRUCTURE Tab](#)

Partial ranking methods

General information | Ranking indices | **Level structure** | Hasse matrix

Level	No. Objs	Objects
11	1	Prochloraz
10	1	Cyproconazole
9	1	Hexaconazole
8	1	Difenoconazole
7	2	Fenbuconazole; Flusilazole
6	4	Diclobutrazol; Penconazole; Propiconazole; Tebuconazole
5	2	Btertanol; Tetraconazole
4	2	Myclobutanil; Triadimefon
3	1	Triadimenol
2	1	Flutriafol
1	1	Pacllobutrazol

In the third tab, the complete structure for the Hasse diagram is reported in a grid, whose rows correspond to diagram's levels. For each level, the number of objects and the list of objects belonging to that level is reported. Objects in equivalence classes are grouped using [parentesi graffe].

[HASSE MATRIX Tab](#)

Partial ranking methods

General information | Ranking indices | Level structure | **Hasse matrix**

ID	Name	1	2	3	4	5	6	7	8	9	10	11	12	13
1	Btertanol	0	-1	-1	-1	-1	-1	1	-1	1	1	0	-1	0
2	Cyproconazole	1	0	1	1	1	1	1	1	1	1	-1	1	
3	Diclobutrazol	1	-1	0	-1	-1	-1	1	-1	1	1	0	-1	0
4	Difenoconazole	1	-1	1	0	1	1	1	-1	1	1	1	-1	1
5	Fenbuconazole	1	-1	1	-1	0	0	1	-1	1	1	1	-1	1
6	Flusilazole	1	-1	1	-1	0	0	1	-1	1	1	1	-1	1
7	Flutriafol	-1	-1	-1	-1	-1	-1	0	-1	-1	1	-1	-1	-1
8	Hexaconazole	1	-1	1	1	1	1	1	0	1	1	1	-1	1
9	Myclobutanil	-1	-1	-1	-1	-1	-1	1	-1	0	1	-1	-1	-1
10	Pacllobutrazol	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	-1	-1	-1
11	Penconazole	0	-1	0	-1	-1	-1	1	-1	1	1	0	-1	0
12	Prochloraz	1	1	1	1	1	1	1	1	1	1	1	0	1
13	Propiconazole	0	-1	0	-1	-1	-1	1	-1	1	1	0	-1	0
14	Tebuconazole	1	-1	0	-1	-1	-1	1	-1	1	1	0	-1	0
15	Tetraconazole	0	-1	0	-1	-1	-1	1	-1	1	1	-1	-1	-1
16	Triadimefon	0	-1	0	-1	-1	-1	1	-1	0	1	-1	-1	-1
17	Triadimenol	-1	-1	-1	-1	-1	-1	1	-1	0	1	-1	-1	-1

Show colors in Hasse matrix

In the fourth tab, the Hasse matrix is reported. By checking the "Show colors in Hasse matrix" voice, the values in the grid will be highlighted with different colours in order to help understanding the matrix structure: light blue for 1 values, light purple for -1 values

and gray for symmetric 1 values (which stands for equivalent objects).

HASSE DIAGRAM menu



In this form, the Hasse diagram is shown. In the lower panel the legend for the diagram is shown; if a class variable has been selected, in the lower right panel the classes legend will be shown. When this form is active, it is possible to access to the ["diagram" menu](#), which contains several options related to the Hasse diagram visualization and manipulation.

In the diagram, classes with more than one object are shown by drawing a gray circle beyond the object, and the equivalence class is reported in the legend.

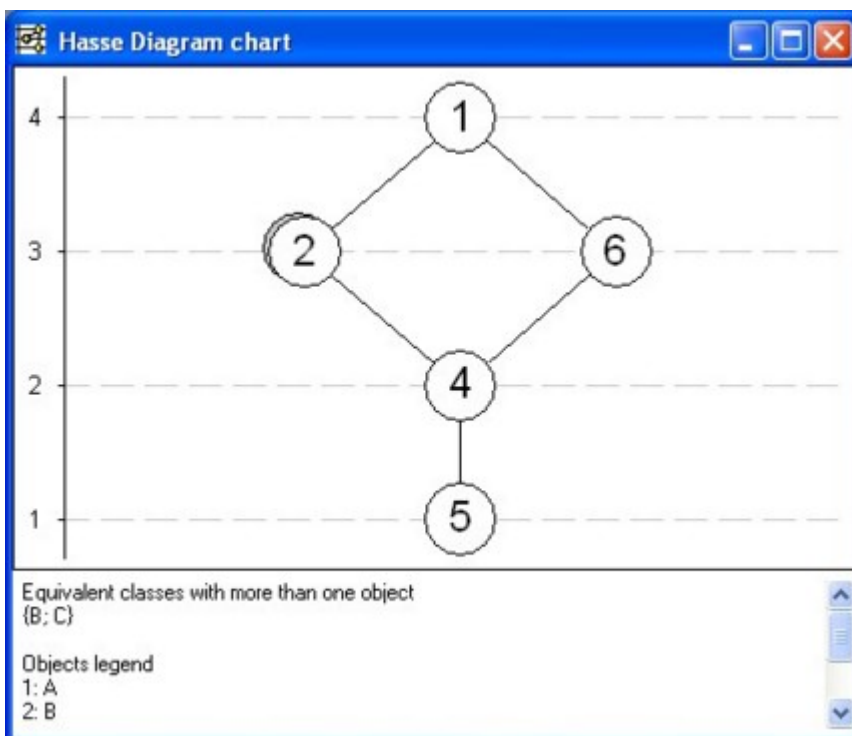


DIAGRAM menu



This menu contains several options related to the visualization and manipulation of the Hasse diagram:

- DRAW MODE: LEFT ALIGN: toggles the left align visualization, in which all objects are aligned to the left of the diagram.
- DRAW MODE: SYMMETRIC: toggles the symmetric visualization, in which all objects are symmetrical centered in the diagram.
- DRAW MODE: GENETIC/SYMMETRIC: toggles the genetic/symmetric visualization, in which all objects are symmetrical centered, after being rearranged using the genetic algorithm (see [Genetic Algorithm](#) section for further details).
- DRAW MODE: GENETIC ALGORITHM: toggles the genetic algorithm visualization, in which all objects are put on the diagram with coordinates calculated by the genetic algorithm (see [Genetic Algorithm](#) section for further details).
- SETUP GENETIC ALGORITHM: opens the genetic algorithm setup for, in which is possible to adjust the parameters used by the algorithm (explained in the [Genetic Algorithm](#) section).
- SHOW GRID: toggles the grid on the diagram, i.e. the level lines.
- SHOW OBJECT NAMES: toggles whether the names or the id number should be shown inside objects on the diagram.
- DRAW LINES OVER OBJECTS: toggles whether the lines joining objects should be drawn over the objects or not.
- SHOW CLASSES COLOR: toggles wheter the objects should be coloured with the class color or not (this voice is active only if a class variable has been defined).
- EXPORT DIAGRAM TO FILE: saves the current diagram as a bitmap file.
- COPY DIAGRAM TO CLIPBOARD: copies the current diagram to clipboard, so that it can be used by other Windows applications.
- PRINT DIAGRAM: opens the print dialog box, in which it is possible to set the desired options before sending the current diagram to printer.

Genetic algorithm for Hasse diagram



As often the Hasse diagram get quite complex, when several objects are included, the problem of drawing the diagram in the "best" way possible is raised. It is not easy to determine an optimal endpoint for which a diagram is the most easy-to-be-read, anyway as the idea is that the best diagram should have its objects rearranged so that as few lines as possible cross each other, a genetic algorithm that optimizes this condition has been developed.

The genetic algorithm is used in two different ways. As sometimes it's not easy to find the right adjustment for algorithm's parameters, together with the straight algorithm (explained below) it is given also a symmetric/genetic algorithm, which means that after getting a suitable solution from the genetic algorithm, objects in the diagram are disposed as in the symmetric view option.

[ALGORITHM](#)

The population of this algorithm is made by the possible diagrams, and the fitness function is:

$$F = w1*Y + w2*M - w3*C$$

where Y is the average squared horizontal distance between objects on the same level, M is the minimum squared distance between objects on the same level, C is the average squared horizontal distance between all couples of objects, and w1, w2, w3 are the weights of the three parameters. Other parameters for the algorithm are E, the size of the elite population, and N, the number of cycles to be performed.

The genetic algorithm performs N cycles on E populations, made by E*E organism, and for each cycle only the E best organisms for each population mate each other to build the population for the next cycle; after N cycles, the best organism from each population is taken in order to build an elite population, on which N further cycles are executed; at the end, the best organism of this elite population is taken as the final result.

Missing values



DART can handle dataset with missing values, as all the total ranking and partial ranking calculations can be performed even if some values in the dataset are missing.

Objects/samples with missing values are indicated with a * character at the end of their names, to remind that calculations for that object are somehow approximate. Also the preprocessing techniques can be performed on datasets containing missing values, except for the Principal Component Analysis.

Anyway, when missing values are found, DART gives you the possibility of filling them with some desired values, such as minimum, maximum, mean or random values (see [data Setup](#) section for further details).

Introduction



A way to perform data exploration is by rank methods which analyse the order relationship among elements. The different kinds of order methods available can be roughly classified as total (called even-scoring) and partial-order ranking methods, according to the specific order they provide. These methods are the ones needed to support and solve decision problems, setting priorities. Besides sophisticated multivariate statistics, used mostly in pre-processing and modelling data, priority setting makes use of quite simple methodologies. However the increasing of problem complexity leads to the decision processes becoming more complex, requiring the support of new tools. Thus there has been increased interest in decision making strategies and several techniques have been proposed. The intrinsic complexity of the systems analysed in chemistry research, and the multiplicity of objectives involved like economic efficiency, environmental quality and availability of resources, has led to complex multicriteria decision problems.

In the complex systems evaluated by ranking strategies, elements (chemical substances, chemical processes, regions,...) are described by several attributes, referred to also as the criteria; thus the system must be analysed by more than one criterion, and decisions must be made by taking several criteria into account contemporaneously.

The criteria are any set of attributes which must reliably represent the system required properties and which must be orientable, i.e. for each criterion it is necessary to explicitly ascertain whether the best condition is satisfied by a minimum or maximum value of the criterion.

Total and Partial order ranking strategies, which from a mathematical point of view are based on elementary methods of Discrete Mathematics, appear an attractive and simple tool to perform data analysis. A complete evaluation by the ranking technique requires a pre-processing phase to establish an adequate data matrix, and a post-processing phase to extract information and decisions on the system investigated. Obviously both pre-processing and post-processing may influence the results significantly. Pre-processing statistical techniques like Clustering, Principal Component analysis and broad order statistics seems to be very suitable tools, providing a satisfactory solution to those drawbacks related to noise and measurement error.

Total and partial order rankings can be analysed to establish the quality of the result obtained. As is usual for regression and classification strategies, the quality of a ranking procedure has to be evaluated by a deep analysis and by several indices, i.e. scalar functions which describe features of an ordered set and allow comparison among different rankings. Thus, the post-processing phase mainly consists in evaluating the quality of the ranking procedure by calculating ranking indices.

Preprocessing techniques



Preprocessing

Preprocessing techniques used in DART consists in various well-known statistical techniques to be applied to the dataset before any further ranking analysis, with the aim of obtaining better results from the desired analysis. Usually, preprocessing appears to be especially useful for partial ranking, as in this case one of the most relevant problem is the growth of incomparable objects; preprocessing the dataset can lead to a reduction of incomparabilities, thus to a better ranking result.

The preprocessing techniques used in DART are briefly explained here.

Significant digits

The most simply technique, it consists in rounding all values to a given decimal digit. Rounding can help to reduce incomparabilities in the Hasse diagram.

Bins partition

The method of bins partition consists in dividing the variable range into n bins, i.e. regular intervals, consequently each element's value is set to the average value of the bin in which it falls. This operation makes the dataset more homogeneous, as similar values are set to an exact and equal value, thus resulting very useful for Hasse diagram.

Principal Component Analysis

From a mathematical point of view the aim of principal component analysis is to transform p -correlated variables into a set of orthogonal variables which reproduce the original variance/covariance structure. This means rotating a p -th dimensional space to achieve independence between variables. The new variables, called principal components, are linear combinations of the original variables along the direction of maximum variance in the multivariate space, and each linear combination explains a part of the total variance of the data. Being orthogonal the information contained in each PC is unique. A maximum of p principal axes can be derived from the original data containing p variables. The new variables are defined by calculating eigenvalues and eigenvectors of the correlation matrix C (or the covariance matrix S) obtained from the data matrix X . Because of their properties, PCs can be used to summarize, in a few dimensions, most of the variability of a dispersion matrix of a large number of variables, providing a measure of the amount of variance explained by a few independent principal axes. This means that it is possible to choose only few of the first components and to project the original elements in their new space, obtaining a dataset that has less variables than the original one while it represents the same amount of information.

K-means clustering

Clustering consists in a partition of the elements into k clusters, then representatives (centroids) of each cluster are defined and used to represent all cluster's belonging elements. In this way clusters are considered equivalence classes. Thus, elements assignment to a cluster is based on some measurement of similarity.

K-Means clustering algorithm is an iterative procedure for the division of elements into k arbitrary clusters, starting from a random partition and then moving each element, during each iteration, to the cluster which centroid is more near to the element itself (measurement of distance can be made using different methods: euclidean, mahalanobis, chebyshev).

Total ranking



Total ranking

Total order ranking methods are multicriteria decision making techniques used for the ranking of various alternatives on the basis of more than one criterion. A criterion is a standard by which the elements of the system are judged. Criteria are not always in agreement, they can be conflicting, motivating the need to find an overall optimum that can deviate from the optima of one or more of the single criteria.

Total order ranking methods are based on an aggregation of the criteria y_r , where $r = 1, \dots, R$:

$$\Gamma \equiv f(y_1, y_2, \dots, y_R)$$

Thus, if an element is characterised by R criteria, then a comparison of different elements needs a scalar function, i.e. an order or ranking index, to sort them according to the numerical value of Γ . Several evaluation methods which define a ranking parameter generating a total order ranking have been proposed in the literature [Keller and Massart, 1991; Hendriks et al., 1992; Lewi et al., 1992]; those used in DART are [Desirability functions](#), [Utility functions](#), [Dominance functions](#), [Concordance Analysis](#), [Simple Additive Ranking](#), [Hasse Average Ranking](#) and [Absolute Reference](#) method.

These methods require the definition of the values and situations of optimum, i.e. for each criterion it is necessary to ascertain explicitly if the best condition is satisfied by a minimum or a maximum criterion value, and the trend from the minimum to the maximum must also be established.

The attribute setting is a crucial point in ranking methods since it requires the mathematization of decision criteria which are often not completely defined or explicit.

Total order ranking results are strictly dependent on the criteria setting and thus can be completely different for different settings.

[Hendriks et al., 1992] Hendriks, M.M.W.B., Boer, J.H., Smilde, A.K., Doorbos, D.A. (1992). Multicriteria Decision Making. Chemom.Intell.Lab.Syst., 16, 175-191.

[Keller and Massart, 1991] Keller, H.R., Massart, D.L. (1991). Multicriteria Decision Making: a case study. Chemom.Intell.Lab.Syst., 11, 175-189.

Desirability functions



Desirability functions

Desirability functions are a well-known multicriteria decision making method. The approach is based on the definition of a desirability function for each criterion in order to transform values of the criteria to the same scale. Different kinds of functions can be used, the more common ones being linear, sigmoid, logarithmic, exponential, step, normal, parabolic, Laplace, triangular and box.

Each criterion is independently transformed into a desirability d_{ir} by an arbitrary function which transforms the actual value of each element into a value between 0 and 1:

$$d_{ir} = f_r(y_{ir}) \quad 0 \leq d_{ir} \leq 1 \quad r = 1, 2, \dots, R.$$

r being the selected criterion, f the function chosen and y_{ir} the actual value of the i -th element for the r -th criterion.

Once the kind of function and its trend for each criterion is defined, the global desirability D of each i -th element can be evaluated as follows:

$$D_i = \sqrt[R]{d_{i1} \cdot d_{i2} \cdot \dots \cdot d_{iR}} \quad 0 \leq D_i \leq 1$$

The overall desirability is calculated combining all the desirabilities through a geometrical mean. It must be highlighted that the desirability product is very strict: if an element is poor with respect to one criterion, its overall desirability will be poor. If any desirability d_i is equal to 0 the overall desirability D_i will be zero, whereas the D_i will be equal to one only if all the desirabilities have the maximum value of one.

In addition each criterion can be weighted in order to take into account criterion importance in the decision rule. In the case of weighted desirability functions the overall desirability of the i -th element is defined as follows:

$$D_i = d_{i1}^{w_1} \cdot d_{i2}^{w_2} \cdot \dots \cdot d_{iR}^{w_R} \quad 0 \leq D_i \leq 1$$

w_r being the weight of the r -th criterion and

$$\sum_{r=1}^R W_r = 1$$

Once D for each element has been calculated, all the elements can be ranked according to their D value and the element with the highest D can be selected as the best one, if its D value is acceptable.

A Desirability scale, shown in Table 1.1, was developed by Harrington [Harrington, 1965] :

Scale of D	Quality evaluation
1.00	Improvement beyond this point has no preference
1.00 – 0.80	Acceptable and excellent
0.80 – 0.63	Acceptable and good
0.63 – 0.40	Acceptable but poor
0.40 – 0.30	Borderline
0.30 – 0.00	Unacceptable
0.00	Completely unacceptable

Table 1.1 – Harrington qualitative definition of the Desirability scale.

The critical feature of this approach to multicriteria decision making problems is the establishment of the relation between criteria and desirability values which must be performed by the decision maker.

[Harrington, 1965] Harrington, E.C. (1965). The Desirability Function. Industrial Quality Control., 21, 494-498.

Utility functions



Utility functions

The approach is very similar to the desirability functions; each criterion is independently transformed into a utility u_r by a function which transforms the actual value of each element into a value between 0 and 1.

$$u_{ir} = f_r(y_{ir}) \quad 0 \leq u_{ir} \leq 1$$

r being the selected criterion, f the function selected and y_{ir} the actual value of the i -th element for the r -th criterion.

Once the kind of function and its trend for each criterion has been defined, the overall Utility U of each i -th element is defined as:

$$U_i = \frac{\sum_{r=1}^R u_{ir}}{R} \quad 0 \leq U_i \leq 1$$

In the case of weighted utility functions the overall utility is calculated as:

$$U_i = \sum_{r=1}^R w_r \cdot u_{ir} \quad 0 \leq U_i \leq 1$$

with

$$\sum_{r=1}^R w_r = 1$$

In this case the overall utility is calculated less severely: in fact the overall quality of an element can be high even if a single utility function is zero.

Like the desirability functions, the utility functions are affected by arbitrariness related to the a priori selection of the functions and corresponding upper and lower limits. Both desirability and utility functions are very easy to calculate, thus specific software is not required.

Dominance functions



Dominance functions

The dominance function method is based on the comparison of the state of the different criteria for each pair of elements. This approach does not require the transformation of each criterion into a quantitative function, it has only to be established whether the best condition is satisfied by a minimum or maximum value of the selected criterion.

For each pair of elements (i, j) three sets of criteria are determined:

$R^+(i,j)$ is the set of criteria w^+ where i dominates j, i.e. where i is better than j, $R^0(i,j)$ is the one where i and j are equal, and $R^-(i,j)$ is the set of criteria w^- where i is dominated by j.

The dominance function between two elements i and j is calculated considering the weights of these two sets, considering that the total sum of all of them is always equal to one, as follows:

$$C_{ij} = \frac{1 + \sum_{R^+} w^+}{1 + \sum_{R^-} w^-} \quad 0.5 \leq C_{ij} \leq 2$$

with

$$\sum_{r=1}^R w_r = 1$$

A C_{ij} value equal to 1 means equivalence of the two elements; $C_{ij} > 1$ means that the element i is, on the

whole, superior to the element j, whereas $C_{ij} < 1$ means that the element i is, on the whole, inferior to the element j. The obtained values can be normalised according to:

$$C'_{ij} = \frac{C_{ij} - 0.5}{2 - 0.5} \quad 0 \leq C'_{ij} \leq 1$$

A global score of the i-th element is then calculated as:

$$\Phi_i = \sum_j C'_{ij} \quad 0 \leq \Phi_i \leq N - 1$$

and the corresponding i-th scaled value is:

$$\Phi'_i = \frac{\Phi_i}{N - 1} \quad 0 \leq \Phi'_i \leq 1$$

Elements with higher values of Φ' are the optimal points.

Concordance analysis



Concordance analysis

The use of Concordance Analysis was introduced by Opperhuizen and Hutzinger as a multicriteria decision making method for the priority setting of chemicals [Opperhuizen and Hutzinger, 1982]. The main difference between Concordance Analysis and Desirability, Utility and Dominance functions is the introduction of a reference element to which each element is compared. The reference element can be a real element or a fictitious one: the centroid, i.e. the vector of the means, is frequently used as the fictitious reference element.

Because of the different dimensions of the criteria, each criterion first undergoes normalisation, and each is weighted according to its importance in the decision process. For each criterion the normalised value is compared with the normalised value of the reference element.

For each element Concordance and Discordance sets are defined. The Concordance set $ConSet_i$, related to the i-th element, is composed by those criteria for which the i-th element has values higher than those of the reference element i^* :

$$ConSet_i = \left\{ \forall r \mid (y_{ir} > y_{i^*r}) \cdot w_r \right\}$$

The Discordance set $DiscSet_i$, related to the i-th element, is composed by those criteria for which the i-th element has values lower than or equal to those of the reference element i^* :

$$DiscSet_i = \left\{ \forall r \mid (y_{ir} \leq y_{i^*r}) \cdot w_r \right\}$$

For each element a Concordance Indicator Cl_i , which measures the number of criteria for which the i-th element is preferred to the reference element, is calculated as the sum of the weights belonging to the criteria of the Concordance set, $ConSet_i$:

$$Cl_i = \sum_{r \in ConSet_i} w_r \quad 0 \leq Cl_i \leq 1$$

Similarly a Discordance Indicator DI_i , which quantifies not only the number of criteria with a

worse i-th element than the reference element but also how much worse it is, is calculated as the weighted maximum difference between the criteria of the Discordance set and those of the reference element:

$$DI_i = \max_{r \in \text{DiscSet}_i} |(y_{ir} - y_{i^*r}) \cdot w_r|$$

The maximum is taken over all the criteria of the Discordance set. The elements are ranked according to the global score Γ_i :

$$\Gamma_i = CI_i - DI_i$$

Since both CI_i and DI_i range from 0 to 1, the global scaled score Γ of the i-th element is calculated as: i

$$\Gamma_i = \frac{CI_i - DI_i + 1}{2} \quad 0 \leq \Gamma_i \leq 1$$

It must be pointed out that the Concordance Indicator, as defined in the classical Concordance Analysis proposed by Opperhuizen and Hutzinger, is a measure of the number of criteria for which each element is preferred to the reference element, since the Indicator is defined as the sum of the weights belonging to the criteria of the Concordance set, however no account is taken of the real quantitative distance between the two elements.

[Opperhuizen and Hutzinger, 1982] Opperhuizen, A., Hutzinger, O. (1982). Multicriteria Analysis and Risk Assesment. Chemosphere., 11, 675-678.

Simple Additive Ranking (SAR)



Simple Additive Ranking (SAR)

The Simple Additive Ranking method is based on the ranking of the objects with respect to each criterion separately, and the subsequent aggregation of the weighted ranks by arithmetic mean. After all the ranking r_{ij} of the i-th object for the j-th criterion are calculated, the index value is calculated as

$$S_i = \frac{\sum_{j=1}^p w_j \cdot r_{ij}}{n}$$

Then it is normalized as

$$S'_i = \frac{S_i - 1/n}{1 - 1/n}$$

Hasse Average Ranking (HAR)



Hasse Average Ranking (HAR)

The Hasse Average Ranking method is calculated as an empirical relation based on linear extensions of the partial order. A linear extension is a projection of the partial order into a total order that comply with all the relations in the partial order. Thus, from a partial order ranking it is possible to obtain all its linear extensions, and then to evaluate an average ranking for each object based on its ranking frequencies in all the linear extensions. The main problem is that the number of linear extensions increases dramatically as the number of incomparable pairs increases, so a simple empirical relation has been proposed by Brüggemann et al. [Brüggemann et al., 2004] and Carlsen [Carlsen, 2005] to evaluate the average rank of each object:

$$AR(i) = (n+1) - \frac{(S_i+1) \cdot (n+1)}{(n+1-U_i)}$$

where n is the total number of objects, S_i the number of objects ranked below the i-th in the Hasse diagram, and U_i the number of objects incomparable with the i-th.

[Brüggemann et al., 2004] Rainer Brüggemann, Peter B. Sørensen, Dorte Lerche, and Lars Carlsen, Estimation of Averaged Ranks by a Local Partial Order Model, J. Chem. Inf. Comput. Sci. 2004, 44, 618-625
[Carlsen, 2005] L. Carlsen, A QSAR Approach to Physico-Chemical Data for Organophosphates with Special Focus on Known and Potential Nerve Agents, Internet Electron. J. Mod Des. 2005, 4, 355-366

Absolute reference



Absolute reference

The absolute reference method is based measuring the distance between each element and a reference element, which is supposed to represent the overall optimum of all the considered criteria. This method require the definition of the values and situations of optimum, i.e. for each criterion it is necessary to explicitly ascertain not only whether the best condition is satisfied with a minimum value or a maximum value of the criterion, but also the specific optimum values. To get rid of different criterion dimensions, each criterion first undergoes normalisation and weighting to account for its importance. Once a distance measure has been selected, the Absolute reference method calculates the entire N distances between the elements and the reference element. If the Euclidean distance is selected, the distance of the i-th element from the reference element (i^*) is defined as:

$$d_{i^*} = \sqrt{\sum_{r=1}^R (y_{ir} - y_{i^*r})^2 \cdot W_r}$$

For each element a measure of its similarity with the reference element is derived from the Euclidean distance according to the following expression:

$$S_i = 1 - d_{i^*} \quad 0 \leq S_i \leq 1$$

This similarity measure is used to rank the elements. It ranges from 0 (no similarity exists between the considered element and the reference one) and 1 (there is complete similarity

between the considered element and the reference one).

It must be pointed out that this approach measures how far each element is from the reference element, the distance being calculated in only one direction, i.e. on the assumption that none of the investigated elements can be better than the one selected as the reference.

Indices for total ranking



Degeneracy indices

A degeneracy index $k(N)$ was proposed by Bruggemann [Bruggemann and Halfon, 1999] to measure the degeneracy of an order ranking. First proposed for partial ordered rankings, but easily applicable to total ordered rankings, it is defined as:

$$k = \sum_{c=1}^C n_c \cdot (n_c - 1)$$

n_c being the number of elements of the c -th equivalence class and C the total number of equivalence classes.

The corresponding standardized index is:

$$k_{std} = \frac{\sum_{c=1}^C n_c \cdot (n_c - 1)}{N(N-1)} \quad 0 \leq k_{std} \leq 1$$

Note that in the case of two equivalence classes, one containing five elements and the other only one element, the k_{std} index takes a value equal to 0.67, whereas in the case of two equivalence classes, each containing three elements, the k_{std} index takes a value equal to 0.40. Thus the more the degeneracy is shared among the equivalence classes, the less is the numerical value of ; thus this index depends not only on the degeneracy degree but also on the degeneracy distribution in the equivalent classes. k_{std}

To avoid degeneracy distribution dependency an absolute degeneracy degree (D) of a ranking is defined as:

$$D = \frac{\sum_{c=1}^C \left(\frac{n_c}{N} - \frac{1}{N} \right)}{\frac{(N-1)}{N}} \quad 0 \leq D \leq 1$$

The numerator represents the difference between the amount of degeneracy of each equivalence class and the case of total absence of degeneracy (uniform distribution); the denominator corresponds to the maximum value reached by the numerator and is used to scale the values between 0 and 1. Degeneracy D allows the evaluation of the non-uniformity or diversity of the element distribution; D takes a value of 1 when all the elements have the same value as the ranking parameter Γ , in which case the degeneracy is maximum and the total ranking method used is not able to differentiate the elements, i.e. the elements are correlated and only one equivalence class exists. On the other hand D takes the value of 0 for minimum degeneracy when all the elements differ from each other, and N equivalence

classes exist, each with only one element. The greater the degeneracy, the less the diversity of the elements.

Discrimination power by ranking

In most of the cases, total ordered ranking methods are used in multicriteria decision-making problems with the aim of defining priorities. For this purpose, of particular relevance is the method capability of differentiating the elements with different values of the ranking parameter. The quality of an order set can be quantified by the index proposed here, called Discrimination power by Ranking (DbyR), which measures the capability of discriminating elements by a ranking according to the following expression:

$$DbyR = 1 - \frac{D}{L} \quad 0 \leq DbyR \leq 1$$

D being the absolute degeneracy degree and L the number of Levels, i.e. the number of different values of the ranking parameter Γ .

This index ranges from value 0 for the case of all elements equal to each other, i.e. only one equivalent class ($D = L$; $L = 1$), to 1 for the case of totally ordered sequence with no degeneracy ($D = 0$); and increases with the decreasing of the degeneracy index.

Stability index

A total ranking, performed by any whatever total ranking method, is strictly determined by the set of criteria used to describe the system, thus by changing the criteria different rankings arise. The set of criteria used may vary, and an additional criterion may be used. Thus it is of interest to forecast the effect on the ranking of increasing the number of considered criteria, i.e. to evaluate ranking stability. The stability ranking index for a total ordered sequence is defined as:

$$StR = \frac{1 - \sqrt{D}}{L} \quad 0 \leq StR \leq 1/N$$

where D is the degeneracy index and L the number of levels.

This index allows the distinguishing of the case of totally ordered sequence with no degeneracy from the case of full degeneracy, in fact it ranges from 0 for full degeneracy to $1/N$, which is assumed as the stability of an ordered sequence of N elements.

[Bruggemann and Halfon, 1999] Bruggemann, R. and Halfon, E. Introduction to the General Principles of the Partial Order Ranking Theory. In Order Theoretical Tools in Environmental Sciences; Proceedings of the Second Workshop on Order Theoretical Tools in Environmental Sciences, 7-43.

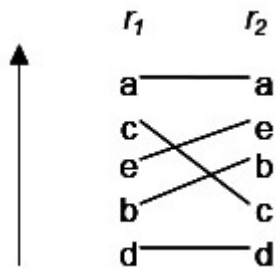
Partial ranking



Partial ranking

Ordering is one of the possible ways to analyse data and to get an overview over the elements of a

system. The elements are commonly characterised by more than one quantity, i.e. they are described by several variables. As a consequence of the multivariate property of the elements, their ordering requires specific techniques as “conflict” among the criteria is bound to exist. Total order ranking methods, being scalar methods, combine the different criteria values into an index, the ranking index Γ , and element comparison and ordering is performed according to the numerical value of Γ . In this way the elements are always ranked in a total or linear ordered sequence, but the information on conflict among criteria is inevitably lost. Partial order ranking is a vectorial approach that recognizes that not all elements can be directly compared with all other elements because, when many criteria are used, contradictions in the ranking can be present. An example could help to better understand what criteria conflict is. The system is made up of five, not perfectly correlated, elements (a, b, c, d, e), each described by two criteria r_1 and r_2 , and the aim is to discover which element is better than the other with respect to all the criteria. The elements are sorted, arranging them according to r_1 and r_2 in the permutation diagram or by parallel coordinates with a vertical orientation, as shown in figure:



This representation highlights the inversions between the two criteria. Elements mutually exchange their position according to the criterion used to sort them. Obviously the higher the number of criteria, the higher the probability that contradictions in the ranking exist. The partial ranking approach not only ranks elements but also identifies contradictions in the criteria used for ranking: some "residual order" remains when many criteria are considered and this motivates the term "partial order". Thus the more known concept of order is the one demanding that all elements be comparable i.e. linear or total order, while partial order is the one in which elements can be “not comparable”. If many elements are to be investigated, and especially if many criteria are to be considered, the parallel coordinates become complex and confusing. The [Hasse diagram technique](#) is a useful tool to perform partial order rankings with an easy visualisation of the obtained results.

Hasse Diagram Technique



Hasse Diagram Technique (HDT)

The Hasse diagram technique is a partial order ranking technique introduced in environmental sciences by Halfon [Halfon and Reggiani, 1986] and refined by Bruggemann Bruggemann [Bruggemann and Bartel, 1999c]. It is based on a specific order relation, named product order, and it provides a diagram, which visualises the results of the sorting.

In this approach the basis for ranking is the information collected in the full set of criteria, called even attributes, E , which is called the "information basis" of the comparative evaluation of elements.

The processed data matrix Q ($N \times R$) contains N elements (rows) and R attributes (columns). The entry y_{ir} of Q is the numerical value of the r -th attribute of the i -th element. According to the product order relation, which the Hasse diagram technique is based on, IB being the information basis of evaluation and E the set of N elements, the two elements s and t are comparable if for all $y_r \in IB$ either $y_r(s) \leq y_r(t)$ or $y_r(t) \geq y_r(s)$. If $y_r(s) \leq y_r(t)$ for all $y_r \in IB$ then

$s \leq t$.

The request "for all" is very important and is called the generality principle:

$$s, t \in E; s \leq t \Leftrightarrow y(s) \leq y(t)$$

$$y(s) \leq y(t) \Leftrightarrow y_r(s) \leq y_r(t) \text{ for all } y_r \in IB$$

If there are some y_r , for which $y_r(s) < y_r(t)$ and some others for which $y_r(s) > y_r(t)$ then s and t are not comparable, and the common notation is st . If only one attribute is used or all the attributes are perfectly correlated then total order is obtained, and all the elements are comparable.

Partial order is determined by the actual information base, thus by changing the information base (IB) different partial orders arise. Partial order sets can be developed easily with the Hasse diagram technique, comparing each pair of elements and storing this information in the Hasse matrix which is a $(N \times N)$ antisymmetric matrix. For each pair of elements s and t the entry h_{st} of this matrix is:

$$h_{st} \begin{cases} +1 & \text{if } y_r(s) \geq y_r(t) \text{ for all } y_r \in IB \\ -1 & \text{if } y_r(s) < y_r(t) \text{ for all } y_r \in IB \\ 0 & \text{otherwise} \end{cases}$$

Thus according to the so-called cover-relation, if there is no element "a" of E , for which $s \leq a \leq t$, $a \neq s$, t and $s \neq t$, then s is covered by t , and t covers s .

The results of the partial order ranking is visualised in a diagram which is constructed as follows:

1. 1. each element is represented by a small circle
2. 2. within each circle the element name, or the equivalence class, is given. Equivalent elements are different elements that have the same numerical values with respect to a given set of attributes. The equality according to a set of attributes defines an equivalence relation
3. 3. if an order or cover relation exists then a line between the corresponding pairs of elements is drawn, the elements belonging to an order relation are "comparable"
4. 4. if $s \leq t$ then s is drawn below t , therefore the diagram has orientation, consequently a sequence of lines can only be read in one direction either upwards or downwards
5. 5. if $s \leq t$ and $t \leq z$ then $s \leq z$ according to the transitivity rule; however a line between s and z is not drawn because this connection can be deduced from the lines between s and t and t and z
6. 6. if either $s \leq t$ or $t \leq s$ then s and t are not connected by a line; thus they are called "incomparable"
7. 7. 'incomparable' elements are located at the same geometrical height and as high as possible in the diagram, resulting in a structure of levels. Elements belonging to a given level are 'incomparable'. Note, however, that a location of elements at different levels does not imply comparability.

In the Hasse diagram, the elements at the top of the diagram are called maximals and there are no elements above them; instead elements which have no elements below are called minimals and they do not cover any further element. If there is only one minimal element, then this is called the least element and if there is only one maximal element, it is called the greatest element. In the environmental field, where the Hasse technique was first applied, the criteria describe the elements in terms of environmental hazard. The main assumption is that the lower the numerical value the lower the hazard. If a high numerical value of an attribute corresponds to low hazard the attribute values must be multiplied by -1 to invert their order.

Therefore, by this convention, the maximal elements are the most hazardous, and are selected to form the set of priority elements. Elements that are not comparable with any other element are called isolated elements, and can be seen as maximals and minimals at once: according to the caution principle they are located at the top of diagram within those elements that require priority attention.

A chain is a set of comparable elements, therefore levels can be defined as the longest chain within the diagram. An antichain is a set of mutually incomparable elements. The height (longest chain) and width (longest antichain) of an order set are indicators of the relative number of comparable pairs of elements compared to the total number of pairs. An example is provided to understand the Hasse diagram interpretation. Let E be the set of 10 elements, and IB the information basis of four attributes describing the elements then the data matrix processed is:

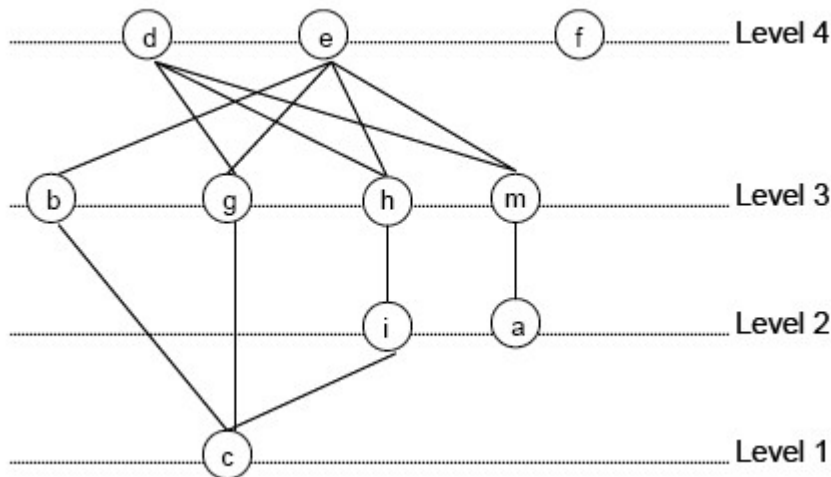
Element	r1	r2	r3	r4
a	15	4	6	8

b	12	22	57	31
c	3	5	6	8
d	44	33	54	33
e	22	38	66	35
f	11	2	69	27
g	6	29	44	28
h	14	31	32	22
i	13	18	20	21
m	18	19	23	28

The corresponding Hasse matrix is:

	a	b	c	d	e	f	g	h	i	m
a	-	0	0	-1	-1	0	0	0	0	-1
b	0	-	1	0	-1	0	0	0	0	0
c	0	-1	-	-1	-1	0	-1	-1	-1	-1
d	1	0	1	-	0	0	1	1	1	1
e	1	1	1	0	-	0	1	1	1	1
f	0	0	0	0	0	-	0	0	0	0
g	0	0	1	-1	-1	0	-	0	0	0
h	0	0	1	-1	-1	0	0	-	1	0
i	0	0	1	-1	-1	0	0	-1	-	-1
m	1	0	1	-1	-1	0	0	0	1	-

The corresponding Hasse diagram is:



In this Hasse diagram there are no equivalence classes; the elements are arranged in four levels, elements d and e are maximals, and are not covered by any other element. The element f is an isolated element since it is not comparable with any of the other elements; c is a minimal and, especially, a least element. Several chains arise: for example $d \geq g \geq c$ and $e \geq h \geq i \geq c$. Obviously maximals are mutually incomparable. At level 3, b and g are not comparable. Incomparability is due to contradictory attributes: for each incomparable pair of elements there must be at least one pair of attributes of counteracting values. Such attributes are called antagonistic. The key diagram interpretation is provided by the meaning of chain and antichain. A chain indicates that the values of the attributes increase synchronously, whereas antichains correspond to diverse patterns. Thus if attributes describe the hazard caused by chemicals which are toxic to different species, then maximals are those elements of highest priority, the most toxic ones, whilst incomparability expresses a diverse pattern of toxicity e.g. toxicity to different species. In this case maximal elements are, in the same way, of priority attention, being toxic but in a different way.

In accordance with the literature [Bruggemann et al., 1993b], the Hasse diagram technique has some relevant advantages:

- evaluation can be represented as a graph
- the mathematics is very simple
- it can easily manage criteria of different scales (linguistic, ordinal and ratio-scaled criteria) since it does not perform any numerical aggregation of the criteria.

Nevertheless there are some severe drawbacks:

- It is strictly dependent on the clarity of the graphical diagram: diagrams that are too complex or too poorly structured, with more isolated elements than comparable elements because of conflict, are of little use.
- if there are too many contradictions criterion reduction must be performed by preliminary multivariate statistic techniques, like Principal Component Analysis (PCA) Multidimensional Scaling.
- if many elements are to be evaluated, preliminary multivariate statistic techniques, like Cluster analysis, are needed to get a readable diagram
- the generality principle is very restrictive and requires appropriate data handling. In fact it must be ensured that any two elements ordered by ">" can be considered as physically and numerically significantly different, i.e. they should have numerically significant data differences. Differences within statistical noise, numerical uncertainty and experimental error are considered physically meaningless, but the Hasse diagram technique considers such elements as different.

Data, and especially environmental data, are often associated with a significant degree of uncertainty inherent in ranking analysis. The comparison of two elements (comparable/incomparable) and thus of the ranking can obviously be affected by this uncertainty. There are two main sources of ranking uncertainty: the relationship assumed between

the attributes and the phenomenon described by the ranking, and the input uncertainty. The first type of uncertainty can be minimized by increasing the number of attributes so that a large number of different aspects are taken into account; nevertheless the greater the number of attributes, the higher the probability that contradictions will occur in ranking the attributes (incomparabilities), and thus the greater the uncertainty in ranking the elements.

Uncertainty from input is the uncertainty induced from variability in the input parameters, which may be due to true variability or to errors in the procedure used to determine the values.

[Bruggemann et al., 1993b] Bruggemann, R., Bucherl, C., Pudenz, S., Steinberg, C.E.W. (1993b). Application of the Concept of Partial Order on Comparative Evaluation of Environmental Chemicals. *Acta Hydrochim. Hydrobiol.*, 27, 170-178.

[Bruggemann and Bartel, 1999c] Bruggemann, R., Bartel, H-G. (1999c). A Theoretical Concept to Rank Environmentally Significant Chemicals. *J.Chem.Inf.Comput.Sci.*, 39, 211-217.

[Halfon and Reggiani, 1986] Halfon, E., Reggiani, M.G. (1986). On Ranking Chemicals for Environmental Hazard. *Environ. Sci. Technol.*, 20, 1173-1179.

Indices for partial ranking



Degeneracy indices

A degeneracy index $k(N)$ was proposed by Bruggemann to measure the degeneracy of an order ranking. First proposed for partial ordered rankings, but easily applicable to total ordered rankings, it is defined as:

$$k = \sum_{c=1}^C n_c \cdot (n_c - 1)$$

n_c being the number of elements of the c -th equivalence class and C the total number of

equivalence classes.

The corresponding standardized index is:

$$k_{std} = \frac{\sum_{c=1}^c n_c \cdot (n_c - 1)}{N(N-1)} \quad 0 \leq k_{std} \leq 1$$

Note that in the case of two equivalence classes, one containing five elements and the other only one element, the kstd index takes a value equal to 0.67, whereas in the case of two equivalence classes, each containing three elements, the kstd index takes a value equal to 0.40. Thus the more the degeneracy is shared among the equivalence classes, the less is the numerical value of ; thus this index depends not only on the degeneracy degree but also on the degeneracy distribution in the equivalent classes. kstd

To avoid degeneracy distribution dependency an absolute degeneracy degree (D) of a ranking is defined as:

$$D = \frac{\sum_{c=1}^c \left(\frac{n_c}{N} - \frac{1}{N} \right)}{\frac{(N-1)}{N}} \quad 0 \leq D \leq 1$$

The numerator represents the difference between the amount of degeneracy of each equivalence class and the case of total absence of degeneracy (uniform distribution); the denominator corresponds to the maximum value reached by the numerator and is used to scale the values between 0 and 1. Degeneracy D allows the evaluation of the non-uniformity or diversity of the element distribution; D takes a value of 1 when all the elements have the same value as the ranking parameter Γ , in which case the degeneracy is maximum and the total ranking method used is not able to differentiate the elements, i.e. the elements are correlated and only one equivalence class exists. On the other hand D takes the value of 0 for minimum degeneracy when all the elements differ from each other, and N equivalence classes exist, each with only one element. The greater the degeneracy, the less the diversity of the elements.

Comparability degree

Peculiar information encoded in a partial ordered set is that related to the comparability degree which can be quantified by a simple comparability index (χ). Taking into account the number of comparabilities in the ranking, it is defined as:

$$\chi = \frac{V(N,R)}{N(N-1)/2} \quad 0 \leq \chi \leq 1$$

where the numerator $V(N,R)$ represents the number of comparable pairs of elements counted in only one direction and the denominator corresponds to the maximum theoretical value and is used to scale the values between 0 and 1.

This index assumes value 1 for the chain case, i.e. total order, which represents the maximum comparability, whereas, value 0 is assumed for the antichain case where no comparabilities exist. It must be observed that this index assumes value 1 for both the one chain case and the theoretical case of all elements equal each other, as in both these cases the comparability is maximum.

Discrimination power by ranking

When partial order ranking is performed for priority settings it is of great relevance to evaluate the ranking procedure capability of discriminating elements according to different ranks; the discrimination power by ranking (DbyR) of an order set, proposed in a similar formula even for totally ordered rankings, can be calculated as:

$$D_{byR} = \chi - \frac{D}{L} \quad 0 \leq D_{byR} \leq 1$$

where χ is the comparability degree, D the degeneracy degree and L the number of levels. It can be observed that in the case of a chain, total order ranking, the comparability degree takes value 1, thus this expression is equal to the one defined for total ranking. The discrimination power index ranges from value 0 for the case of one antichain to 1 for the case of one chain, and increases with the increasing of the comparability degree and the decreasing of the degeneracy index: it can be observed that D_{byR} , by taking into account both the number of comparabilities in the ranking and the amount of degeneracy of each equivalence class, permits the distinguishing of the case of one chain with no degeneracy, for which $\chi = 1$, $D = 0$ and thus $D_{byR} = 1$, from the case of all elements equal to each other for which $\chi = 1$, $D = 1$ and thus $D_{byR} = 0$.

Stability indices

Partial order ranking is determined by the criteria considered in the ranking procedure, the actual information base, thus by changing the information base (IB), different orders arise. The set of criteria used may vary, and an additional criterion may be used in the information basis. Thus it is of interest to forecast the effect on the ranking of increasing the number of considered criteria, i.e. evaluate the ranking stability. The stability index proposed in the literature is defined as follows:

$$P(N,R) = U(N,R) / S(N)$$

with

$$S(N) = U(N,R) + 2 \cdot V(N,R) - k(N,R)$$

where $V(N,R)$ is the number of comparabilities, $U(N,R)$ the number of incomparabilities (counted in both the directions) and $k(N,R)$ the Bruggemann degeneracy index. This stability index ranges from 0 to 1: when $P(N,R)$ is near zero, then $U(N,R)$ must be near zero and in such a case adding an attribute may have quite a big influence on the ranking, in fact the higher the number of criteria, the greater the probability that contradictions (incomparabilities) in ranking exist among criteria.

Conversely, when $P(N,R)$ is near 1, then $U(N,R)$ must be near $S(N)$, and adding an attribute may have a little influence on the ranking.

The quantity $P(N,R)$ does not differentiate between full degeneracy and the one chain case because, in both cases, no incomparability appears ($U(N,R)=0$).

Nevertheless the stability of the case of full degeneracy is different from the one chain case stability. An example may be useful to better understand this concept: let E be the set constituted of two elements (s,t) and IB the actual information base of R attributes

In the case of full degeneracy:

$$(E, IB) = \{(s,t)\}$$

on adding an attribute, full degeneracy may still exist or a chain may arise.

In the case of a chain:

$$(E, IB) = \{s,t\}$$

on adding an attribute, the chain may still exist or an antichain may arise. Thus, assuming the antichain case to be the case of maximum stability, because adding an attribute changes nothing as the number of incomparabilities is already maximum, from the moment that the one-chain is nearer the antichain case than the full degeneracy, the one-chain case should be more stable than the full degeneracy case.

To take account of the differing stability of the one chain case and the case of full degeneracy, a new ranking stability index is proposed:

$$StR = \left(\frac{1 - \sqrt{D}}{L} \right)^\chi \quad 0 \leq StR \leq 1$$

where D is the absolute degeneracy index, χ the comparability degree defined above and L the number of levels. This index ranges from 0 for full degeneracy to 1 for the case of one antichain, and increases with decreasing degeneracy, with the comparability decreasing and with the decreasing of the number of levels. It can be observed that for a chain with no degeneracy, where $L = N$ and $D = 0$, StR takes the value of $1/N$, which is assumed as the stability of an ordered chain with N elements.

Diversity index

Other useful information encoded in partial ordered ranking is the diversity existing among the elements. A ranking characterised by many incomparabilities between elements, indicates that the elements analysed are of high diversity as far as concerns the criteria they are described with. Therefore antichain corresponds to maximum diversity which can be measured as:

$$div = \frac{NEL(N,R) - 1}{N - 1} \quad 0 \leq div \leq 1$$

where $NEL(N,R)$ is the number of elements in the level, which contains the most elements, and N is the total number of elements. In an antichain $NEL(N,R) = N$ and $div = 1$; whereas for a chain $NEL(N,R) = 1$ and $div = 0$.

If equivalent classes with more than one element exist, the diversity is calculated as:

$$div = \frac{NEL(N,R) - 1}{Z - 1} \quad 0 \leq div \leq 1$$

where Z is the number of equivalent classes.

Selectivity index

The selectivity of a partial ordered ranking is a measure of its capability to providing a unique orientation from "good" to "bad" and therefore it is assumed maximum in a total chain and minimum in a total antichain, as in this case all the elements are incomparable with each other and no orientation is founded. A selectivity index is defined as:

$$T = \frac{L - 1}{N - 1} \quad 0 \leq T \leq 1$$

and for equivalent classes with more than one element, it is computed as:

$$T = \frac{L - 1}{Z - 1} \quad 0 \leq T \leq 1$$

where Z is the number of equivalent classes.