

# GGobi Manual

Deborah F. Swayne, AT&T Labs – Research

Di Cook, Iowa State University

Andreas Buja, The Wharton School, University of Pennsylvania

Duncan Temple Lang, University of California at Davis

Hadley Wickham, Iowa State University

Michael Lawrence, Iowa State University

September 2006

GGobi is an open source visualization program for exploring high-dimensional data. It provides highly dynamic and interactive graphics such as tours, as well as familiar graphics such as the scatterplot, barchart and parallel coordinates plots. Plots are interactive and linked with brushing and identification.

It includes 2-D displays of projections of points and edges in high-dimensional spaces, scatterplot matrices, parallel coordinate, time series plots and bar charts. Projection tools include average shifted histograms of single variables, plots of pairs of variables, and grand tours of multiple variables. Points can be labelled and brushed with glyphs and colors. Several displays can be open simultaneously and linked for labelling and brushing. Missing data are accommodated and their patterns can be examined.

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Tutorial</b>	<b>4</b>
<b>3</b>	<b>Layout and functionality</b>	<b>9</b>
3.1	The major functions . . . . .	9
3.2	Graphical displays . . . . .	13
<b>4</b>	<b>Data format</b>	<b>15</b>
4.1	XML . . . . .	15
4.2	CSV files . . . . .	15
<b>5</b>	<b>View modes: projections</b>	<b>16</b>
5.1	1D plots . . . . .	16
5.2	XY plots . . . . .	17
5.3	1D Tour . . . . .	17
5.4	2D Tour . . . . .	18
5.5	Rotation: 2D Tour with Three Variables . . . . .	19
5.6	2x1D Tour . . . . .	19
<b>6</b>	<b>Interaction modes</b>	<b>20</b>
6.1	Scaling of axes . . . . .	20
6.2	Brush: brushing of points and edges . . . . .	20
6.3	Identification . . . . .	23
6.4	Edit edges . . . . .	24
6.5	Move points . . . . .	24
<b>7</b>	<b>Tools</b>	<b>26</b>
7.1	Variable manipulation tool . . . . .	26
7.2	Variable transformation tool . . . . .	26
7.3	Sphering . . . . .	27
7.4	Jittering . . . . .	27
7.5	Color schemes . . . . .	27
7.6	Automatic brushing . . . . .	28
7.7	Color & glyph groups . . . . .	28
7.8	Subsetting . . . . .	29
7.9	Controls for missing data . . . . .	30
<b>8</b>	<b>Multiple datasets</b>	<b>30</b>
<b>9</b>	<b>Edges</b>	<b>31</b>
<b>10</b>	<b>Large data</b>	<b>31</b>
10.1	Large $N$ . . . . .	31

<b>11 Differences from XGobi</b>	<b>33</b>
11.1 Multiple displays	33
11.2 Data format	33
11.3 Integration	34
11.4 Variable selection	34
11.5 The variable manipulation tool	34
11.6 Changes in projection, interaction and tools	35
11.7 On-line help	35

# 1 Introduction

This manual gives an overview of the layout and functionality of GGobi, interactive graphical software for exploratory data analysis. Readers who are familiar with XGobi will find much that is familiar in GGobi's design, and might want to read section 11 first, where key differences between the two programs are described. There are several papers that describe the functionality which is present in both packages [6, 14, 13, 4, 8, 9, 7, 5].

You will find that you can use GGobi for simple tasks with virtually no instruction. All that you need is some cursory knowledge of the developments in interactive statistical graphics of the last 15 years, as well as a willingness to experiment with the sample data provided. In parallel with the hands-on learning process, it is probably useful to acquire a basic understanding of the overall layout and functionality of the system. You will be most successful if you gain experience with the system and combine it with creativity and data analytic sophistication.

We begin with a tutorial, and move on to describe GGobi in detail.

## 2 Tutorial

Several sample data files are included with the GGobi distribution, in a directory called *data*, and there you will find the file *olive.csv*, a dataset on olive oil samples from Italy [10]. The olive oil data consists of the percentage composition of 8 fatty acids (palmitic, palmitoleic, stearic, oleic, linoleic, linolenic, arachidic, eicosenoic) found in the lipid fraction of 572 Italian olive oils. There are 9 collection areas, 4 from southern Italy (North and South Apulia, Calabria, Sicily), two from Sardinia (Inland and Coastal) and 3 from northern Italy (Umbria, East and West Liguria). The data is part of a quality control study of olive oils. It is proposed that the olive oils from different regions have different fatty acid signatures.

Start GGobi. On Windows this will be as simple as double-clicking the GGobi icon on the screen. For linux and Mac users, you'll need to open up a unix terminal and type `ggobi (linux)`.

From the **File** menu, choose **Open** and navigate to the GGobi `data` directory to find the `olive.csv` file.

Two windows will appear, the GGobi console and a scatterplot, as shown in Figure 1.

The console has a panel of controls on the left, labeled **XY Plot**, and a variable selection region on the right. You can see that the scatterplot contains a 2-dimensional projection of the data, a plot of Area vs Region. Move the mouse to one of the variable labels in the variable selection region, and leave it in place until the tooltip appears, explaining how to select new variables for the plot. Begin to get a feeling for the data by looking at several of the 2d plots: Area vs palmitic, Region vs oleic, etc.

Next, get acquainted with the main menubar for the console by exploring each of its menus. Pay particular attention to the Display, View, Interaction and Tools menus.

- The File menu contains the usual open and save data, option setting and close actions.
- The Display menu is the interface for opening new plotting windows.
- The View menu is the interface for specifying the projection (1d, 2d, 3d or higher) for the

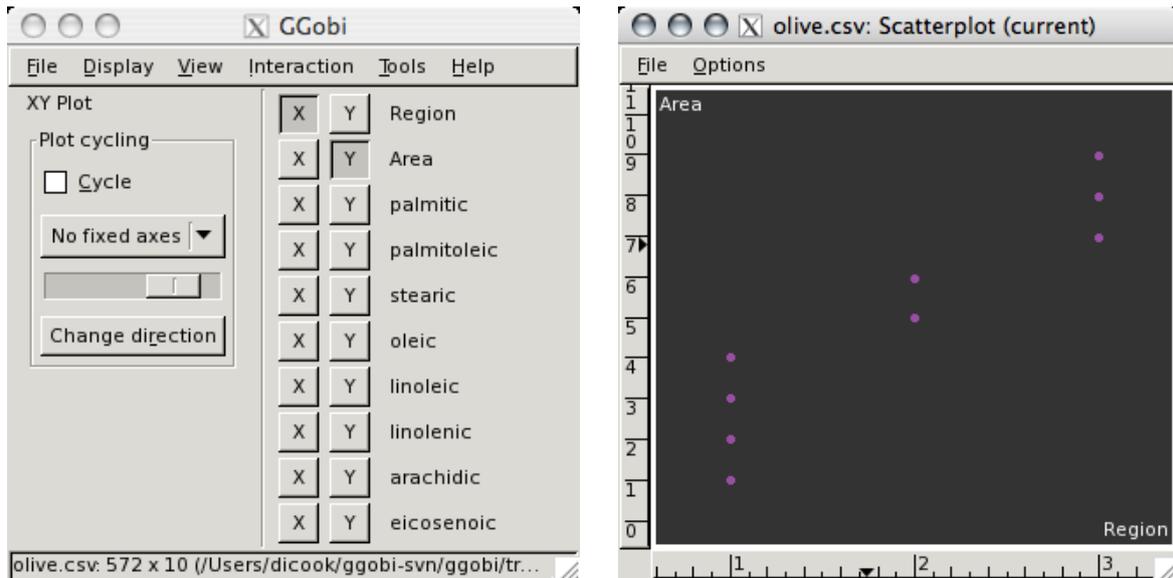


Figure 1: Layout of a GGobi session. The plotting window contains a scatterplot of the Area vs Region, from the olive oils data.

current display. (For some displays, such as the parallel coordinates display, there are no choices in the projection used, so the View menu is not always present.)

- The Interaction menu is the interface for specifying the mouse interactions (for scaling the plot, highlighting points, and so on) for the current plot.
- The Tools menu lets you open other windows to manipulate characteristics of the data and the view.

Using the Interaction menu, choose **Identify** (Figure 2). Look at the buttons inside the leftmost portion of the GGobi console. Notice that they're contained by a frame labeled **Identify**, which is the interaction mode of the current plot. This frame contains most of the row labeling controls, which are described in section 6.3. There are several options for writing labels in the plot window. The default is the record label. Highlight the Region and Area in this panel, then the labels will be the geographic area. Move the cursor around the plot window using the mouse, and observe the labels of the point closest to the cursor. The labels show the geographic area where the sample was taken. Using the variable checkboxes on the main GGobi control window change the variables shown to be linoleic plotted against eicosenoic.

Open a second display window, showing a barchart, using the Display menu. Notice that the new window has a narrow white band drawn around the outside of the plotting area: that means that the new window is now the “current display” and its plot is the “current plot.” Click in the plotting region of the other scatterplot to make it the current plot, and notice what happens in the console when you alternate between the two. Both the controls and the variable selection region correspond to the current plot.

Now set up a plot of linoleic vs eicosenoic in the first scatterplot, and display Region in the barchart, and make the barchart the current plot. Using the Interaction menu, choose **Brush**.

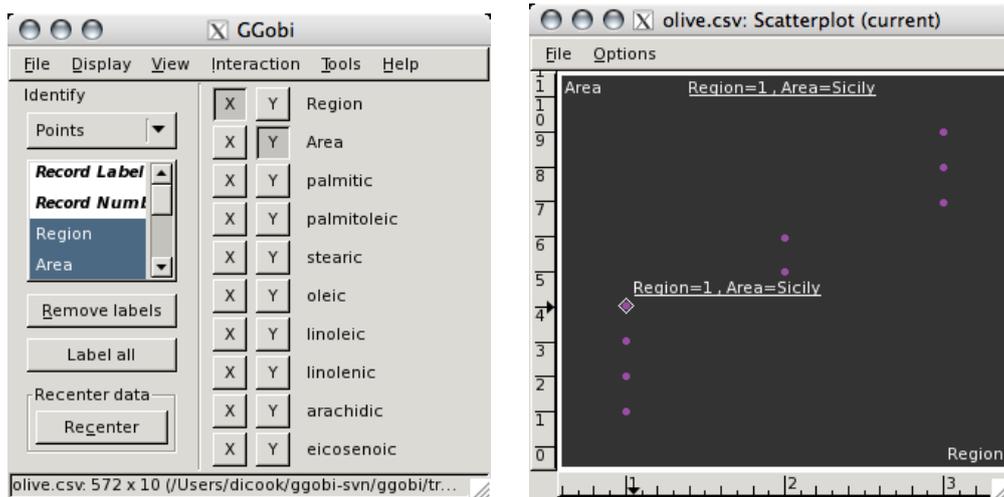


Figure 2: The **Identify** panel and the labels showing in the plot window.

Look at the buttons and menus inside the leftmost portion of the GGobi console. Notice that they're contained by a frame labeled **Brush**, which is the view mode of the current plot. This frame contains most of the brushing controls, which are described in the section 6.2. A few of the brushing controls, though, are in the **Brush** menu in the display menubar, and some associated tools are in the **Tools** menu.

The rectangle that appears in the current plot is the “paintbrush,” and dragging it over bars (or in a scatterplot the points) changes their color. Change the color of the brush by opening up the **Choose color & glyph** panel (Figure 3). Hold down the left button and drag the mouse to paint the first region, then the second and third regions, or click on the bars. Since you're brushing in the **Transient** style, the points resume their original color when the paintbrush no longer surrounds them.

While you brush, keep an eye on the plot of linoleic vs eicosenoic, and notice where the painted points fall in that scatterplot (Figure 4). The oils from region 1, south Italy, are separable from the other two regions using eicosenoic acid. Using two linked 2d plots is one way to explore the relationship among three variables.

Change to **Persistent** brushing and paint the three Regions using different colors.

Open the **Variable transformation** tool, select all of the fatty acid variables – either hold down the control button while selecting them one at a time, or select the first one, then hold down the shift button down while you select the last one. Once they're all highlighted, choose **Standardize** in the “Stage 2” transformation panel to standardize all the variables to have mean 0 and variance 1. This tool has numerous transformations, which can be applied in sequence. Stage 0 transformations are used to adjust variables' values so that they're within the range of some stage 1 transformations. Stage 2 transformations allow for post-processing of stage 1 transformations, so that the variables can be logged and then standardized, for example.

But that was just an instructive detour; click “Reset all” to to turn off standardization and any other transformations you've explored.

Now look at the plot of linoleic vs eicosenoic. Open up the **Color & glyph groups ...** tool,

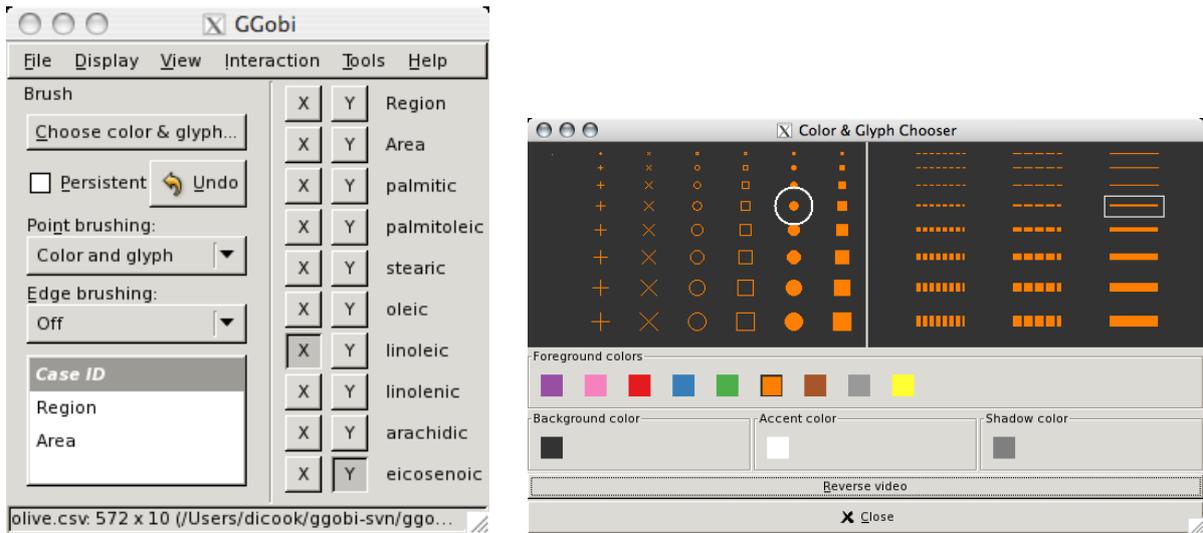


Figure 3: The chooser for selecting color and glyph for painting points, as well as type and thickness for painting lines. It also allows setting the plot background color, and has a full color wheel for tuning the colors.

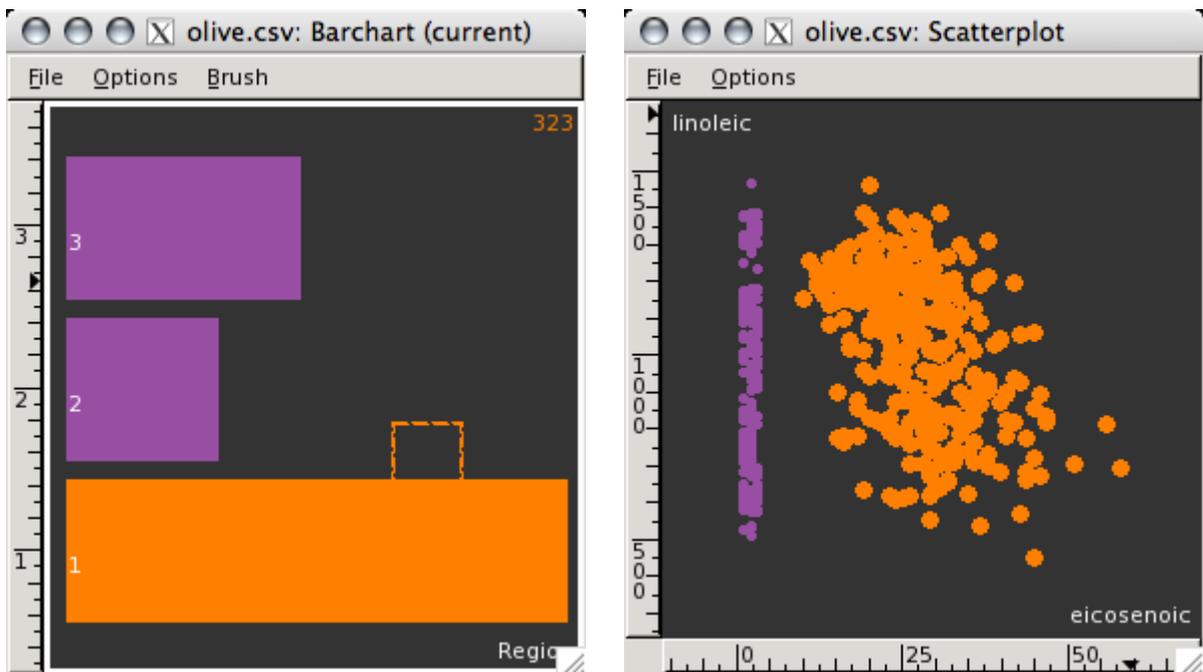


Figure 4: Brushing in a GGobi session. Brushing Region 1 shows that it corresponds to a cohesive cluster in the variables linoleic and eicosenoic acid.

and experiment with it: Use the “Shadow” toggle buttons to have groups of points drawn in a dim color, and then use “Exclude shadows” to exclude them altogether, and “Include shadows” to re-include them. (Figure 5). Finally, use these controls to shadow brush the cases from geographic region 1, and then use **Exclude shadows** to remove them from consideration.

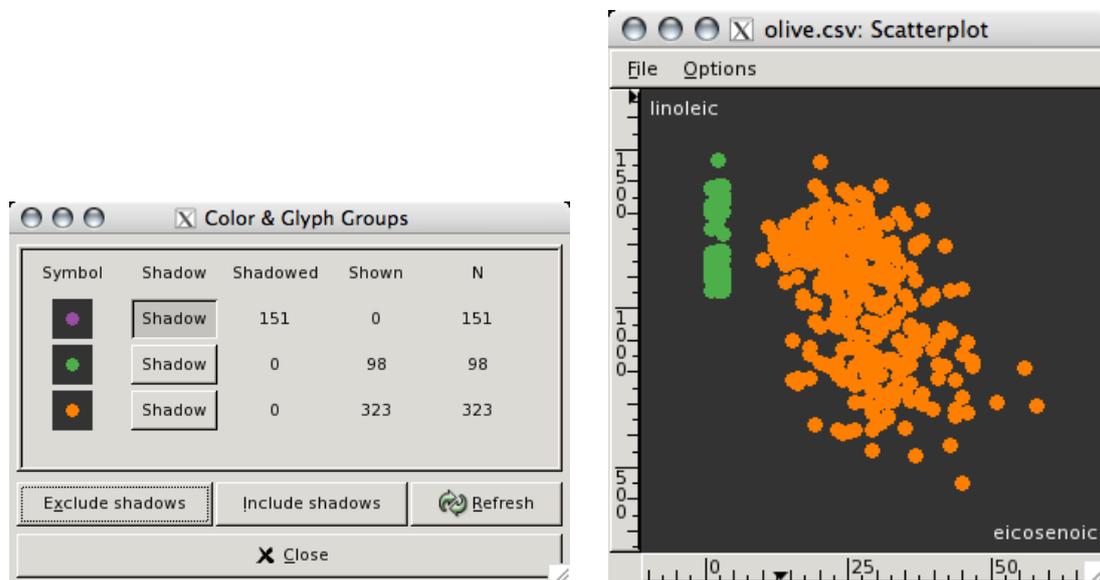


Figure 5: Toggling clusters of points on and off according to color/glyph value.

Switch into **2D Tour** using the **View** menu. A column of variable circles is added to the right of the variable toggles, one for each selected variable. Click on the toggle widgets for the variables palmitoleic through arachidic to make these variables available for touring, and toggle Region and Area out of the tour. Click the **Reinit** button. The tour should now include 7 variables, with one variable circle for each.

The scrollbar at the top of the tour controls is used to adjust the speed of the tour. Drag it to the right to speed up the rotation. The circle at the bottom left of the plot window displays the axes for the tour. These can be removed by toggling the “Show axes” button on the display’s **2D Tour** menu. Pause the tour when you see a separation of the two regions (as seen in Figure 6).

Now you are going to use manual tour controls to sharpen up the separation by manually changing the coefficients of the variables in the projection. Click on the **Manip** (magenta) button below the variable selection region of the console, and then click on linoleic acid (a magenta circle will be seen now in the variable circle for linoleic acid). **Oblique** is already selected in the **Manual manipulation** menu. In the plot window, hold down any mouse button while you drag the cursor. The coefficient corresponding to linoleic acid will increase and decrease following the mouse, inducing a rotation of the scattercloud. If the axes are still showing in the scatterplot window, the axis in magenta corresponds to the variable you’re manipulating. Similarly rotate oleic acid and arachidic acid in and out of the plot, with the aim of finding a projection where the two regions are well separated. There are 5 choices of manipulation mode (unconstrained oblique, vertical, horizontal, radial and angular) to explore.

Next use the Tools menu to open the **Variable manipulation** tool (Figure 7). It contains a

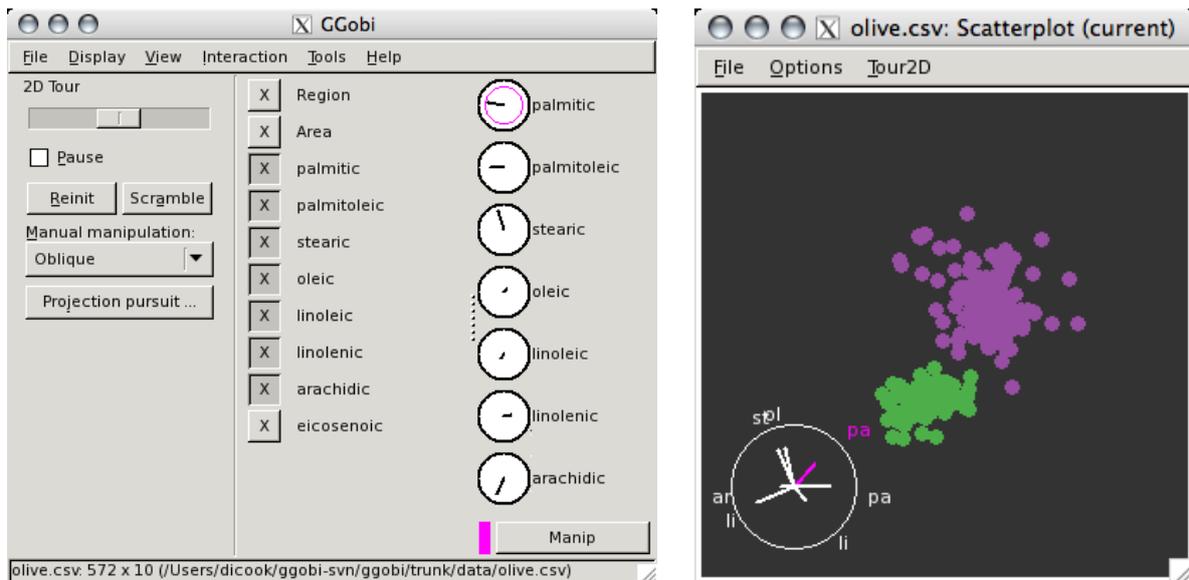


Figure 6: (Left) Tour panel showing the controls for the tour, and (right) a tour projection revealing a difference between regions 2 and 3 of the olive oils.

notebook widget which separates the variables by type, and displays information about them in a set of tables. The buttons below the tables allow you to set variable limits, add new variables, and a few other things. Try clicking the mouse in the table – you can select one row of the table at a time, or use the control and shift keys as modifiers to select more than one. Select all the fatty acid variables, all of which are in the table of real variables.

Open a parallel coordinates display using the **Display** menu. Select the first plot in the parallel coordinates display by clicking on it, and use the **Interaction** menu to switch to **Brush** mode. Choose a new color (say yellow) and large closed glyph, and transiently paint the case with a very low value on palmitic (Figure 8). This case also has a very high value for oleic acid, and low value for linolenic acid.

Using the **Color & glyph groups** tool from the **Brush** controls, click on the appropriate ‘S’ button to bring the Region 1 cases back into the plot. Using the **File** menu, save the data, preserving the colors and glyphs that have been assigned during this session.

This has been a brief introduction to the use of GGobi. The following section contains more detailed information on its functionality.

### 3 Layout and functionality

#### 3.1 The major functions

Across the top of the console, as seen in in Figure 1, stretches a row of menu buttons: **File**, **Display**, **Interaction**, **Tools**, and **Help** are always visible; **View** appears as appropriate.

As expected, the **File** menu contains items for selecting input/output functions and for exiting.

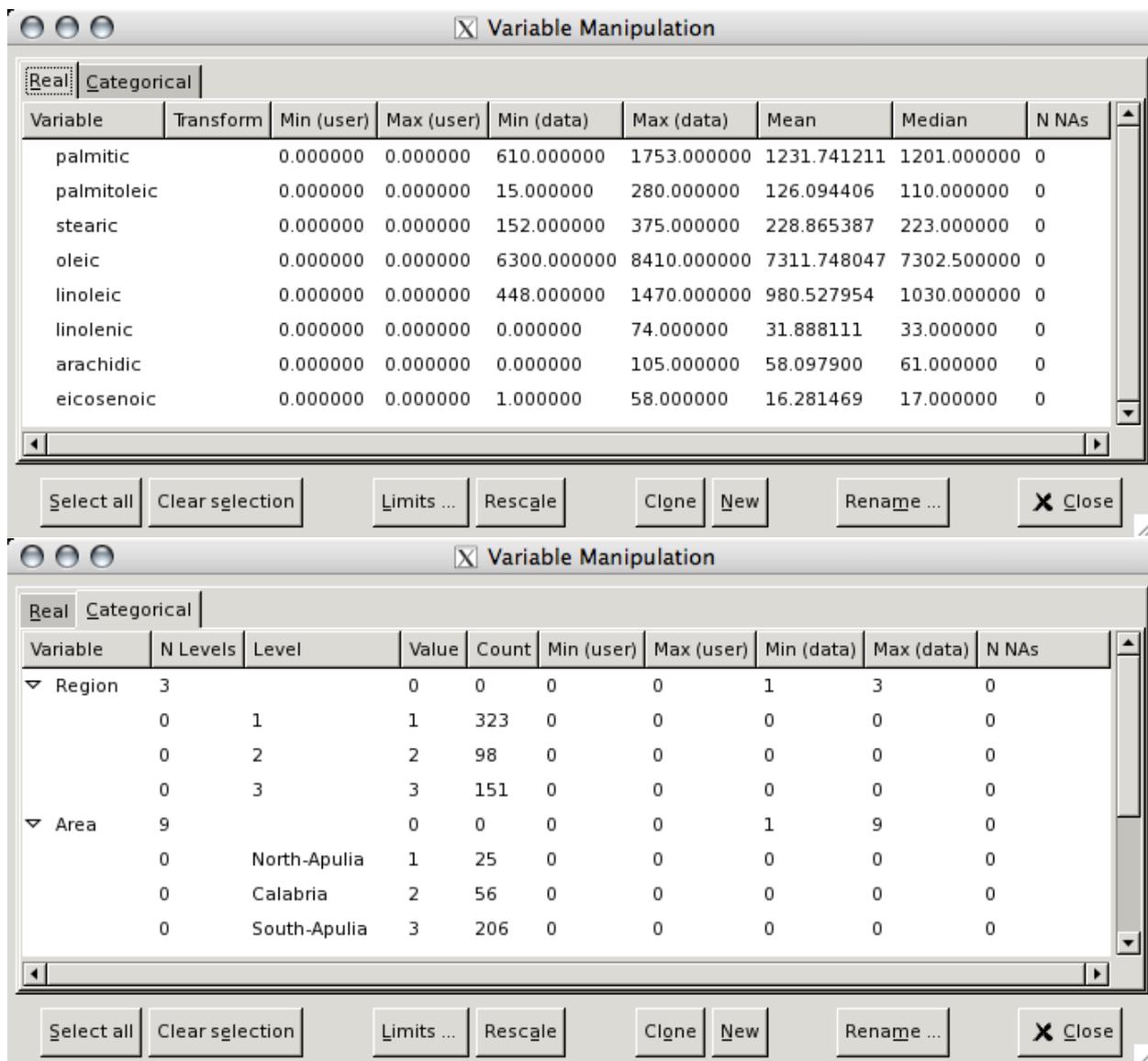


Figure 7: The variable manipulation panel contains basic summary statistics of variables. It can also be used for adding new variables and for variable subset selection before launching multi-plot displays.

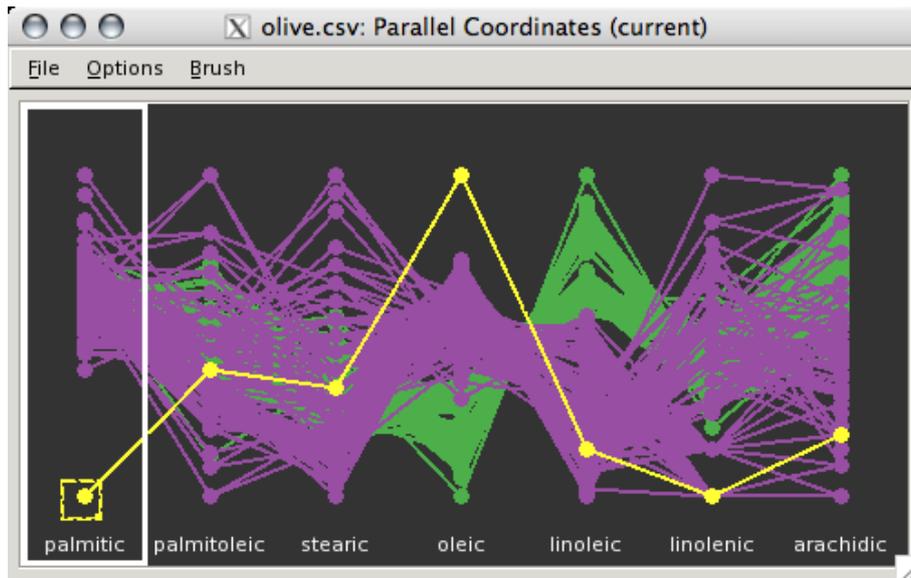


Figure 8: Brushing in a parallel coordinates display reveals an outlier in palmitic, oleic and linolenic acids.

It also contains **Options** allowing users to set options for the main control window: whether tooltips are displayed and whether the control panel is shown.

The **Display** menu allows a new plotting window to be opened. The display types include

- scatterplot,
- scatterplot matrix,
- parallel coordinates plot,
- time series plots, and
- barchart.

Each display type is discussed in section [3.2.2](#).

The **DisplayTree** button at the bottom of the menu allows users to open a tree listing all the currently open display windows, each of which may contain several plots.

The **View** menu contains items to set the projection; the full menu is available only for scatterplot displays:

- 1DPlot: 1-D dotplots and average shifted histograms,
- XYPlot: 2-D scatterplots,
- 1D Tour: 1-D tour,
- Rotation: 2-D tour constrained to use exactly 3 variables,

- 2D Tour: 2-D tour,
- 2x1D Tour: a correlation tour; that is, independent 1-D tours on horizontal and vertical axes,

When you choose a new view, the **Interaction** mode (described next) is also changed. Each view is discussed in section 5.

The **Interaction** menu contains items to set the interaction type. Each **View** also has its own set of interactions, and appears on the Interaction menu when selected. The Interaction modes themselves are:

- Scale: axis scaling,
- Brush: setting point glyphs, edge types, and point and edge colors,
- Identify: labeling points,
- Edit Edges: add points or edges.
- Move Points: direct manipulation of point positions,

When you choose a new interaction mode, the controls at the left of the main window will change correspondingly: each mode has its own parameters and its own rules for responding to mouse actions in the plotting windows. The interaction modes are discussed in section 6.

The **Tools** menu gives access to

- a variable manipulation table,
- a variable transformation pipeline,
- a variable sphering panel,
- jittering controls,
- a panel for selecting a new color scheme for drawing, and then for automatically brushing points and edges by mapping the color scheme onto a selected variable.
- a panel for brushing and excluding groups of cases,
- subsetting functions for systematic and random subsampling,
- a tool for managing missing values.

Each of these tools is discussed in section 7.

## 3.2 Graphical displays

### 3.2.1 Current display, current plot

Since there are multiple displays, some of which contain multiple plots, the question arises: Which plot in which display window corresponds to the console? If you select the **Scale** mode, how can you tell which plot is going to respond?

There is in GGobi a notion of the “current display” and the “current plot.” (We need both because some displays, like the scatterplot matrix, contain multiple plots.) The current plot is the one which is outlined with a thick white border; the current display is the one which contains the current plot.

To reset the current plot and display, just click once (left, right or middle) in the plot you wish to address. To understand the effects of this selection, open a few displays and set them in different interaction modes, then click on different plots and see what happens. The white border follows your actions, and the console updates so that its control panel corresponds to the current display type and view mode.

### 3.2.2 Display types

Each display type is briefly described here. As mentioned earlier, there are presently five main display types:

The **scatterplot** display is a window containing a single scatterplot. It has the largest number of view (projection) modes of any display, and each projection has its own rules for variable selection. The main variable selection interface for the XYPlot modes uses toggle buttons: clicking on the **X** button selects the variable to plot horizontally, and clicking on **Y** selects the variable to plot vertically. (There’s a less obvious variable selection interface as a shortcut: clicking left (right) on the variable label selects the **X** (**Y**) variable.) The interface for the 1D plot is similar: it only shows one column of toggle buttons, so clicking on **X** selects the variable to be plotted irrespective of the plot orientation. (Using the variable label shortcut, you can change the plot orientation as you make a variable selection.)

The tour modes, including rotation, all use toggle buttons to select a subset of variables to be available for touring, and a column of variable circles, one for each variable in the subset. The variable circles can be used to further refine the selection of variables that are actively touring, and they provide some feedback about current projection. The variable selection behavior for the tour modes will be described in the section for each mode.

The **scatterplot matrix** is a window containing a symmetric matrix of scatterplots for the chosen variables. The plots along the diagonal are ASHes (Average Shifted Histograms). The matrix is required to be symmetric, and that constraint affects its variable selection behavior.

- **Replace:** First select one of the ASH plots along the diagonal to tell GGobi uniquely which variable to replace, then click on one of the toggle buttons in the variable selection region.
- **Insert:** First select one of the plots along the diagonal to tell GGobi uniquely where to insert the new plot, then click on one of the buttons. (GGobi will not add a variable that’s already plotted.)

- Append: Click on one of the buttons to append a new plot after all other plots. (GGobi will not add a variable that's already plotted.)
- Delete: No plot selection is required; just select the variable you want to delete.

For the scatterplot matrix display, the variable label “shortcut” works, but it's simply redundant: that is, it makes no difference which button you press.

The **parallel coordinates** display contains a single parallel coordinates plot, which can be arranged horizontally or vertically. (To understand this plot if you are encountering it for the first time, imagine deconstructing a high-dimensional scatterplot and arranging its axes in parallel instead of orthogonally. To represent case  $i$ , think of drawing a dot on each axis, with the point on axis  $j$  being the value of  $x[i][j]$ , and then connecting the dots into one set of connected line segments [11, 18].) The line segments are drawn by default, but you can turn them off using the display's **Options** menu.

By default, the plots are simple dotplots, but they can also be drawn using one of the two methods for 1D plots: as a textured dot strip or an ASH.

The variable selection behavior works as follows:

- Replace: First select one of the plots, then click on one of the toggle buttons to replace its plotted variable.
- Insert: First select one of the plots, then click on one of the toggle buttons to insert a new plot before the current plot. (GGobi will not add a variable that's already plotted.)
- Append: Click on one of the toggle buttons to append a new plot after all other plots. (GGobi will not add a variable that's already plotted.)
- Delete: No plot selection is required; just select the variable whose plot you want to delete.

For the parallel coordinates display, the variable label “shortcut” works, but it's simply redundant: that is, it makes no difference which button you press.

The **time series display** contains a row or column of 2-variable plots with a common axis, usually a time variable. By default, the points are connected with line segments.

The behavior of the toggle buttons and labels depends on the state of the **Selection mode** option menu in the console.

- Replace: If you want to replace the horizontal (time) variable, no plot selection is required; simply click the X toggle button for the variable you want. To replace a vertical variable, first select the plot you want to change, and then click the Y toggle for the variable you want. (GGobi won't let you plot the same variable twice.)
- Insert: Select a plot in the display, and then click on a the Y toggle to add a new plot before the current plot.
- Append: Click on one of the Y toggles to append a new plot after all other plots. (GGobi will not add a variable that's already plotted.)

- Delete: No plot selection is required; just click the Y toggle for the variable in the plot you want to delete.

The variable label shortcuts can be used in the Time Series display. In general, clicking left (right) corresponds to selecting the X (Y) toggle button.

The **barchart display** contains a single plot, a barchart if the variable is categorical, and a histogram if it is real. Variable selection is simple: one variable at a time. An option menu can be used to switch to a spineplot display style, where all the bars are the same height, and it is their width that varies instead.

## 4 Data format

### 4.1 XML

The richest file format supported by GGobi is an XML format, which is described in *The GGobi XML Input Format* (XML.pdf), available from the web site. This format allows a great deal of detailed specification, such as

- multiple datasets within a single XML file, all available within a single GGobi process;
- multiple variable types, such as categorical, integer, and random uniform;
- rules for linking between datasets;
- rules for specifying edges: line segments connecting pairs of points.

### 4.2 CSV files

GGobi reads CSV (Comma-Separated Variables) files, with “,” as a field delimiter and carriage return as a record delimiter. This is a format that has been made popular by Excel. The file extension should be `.csv`.

An example file is:

```
,Var1,Var2,Var3
A,1,F,9
C,2,G,NA
D,-1,NA,3
E,.,F,7
F,1,F,-5
```

The first row contains column labels, and note that it begins with “,” indicating that the first column contains the row labels. Non-numeric variables are treated as categorical. Missing values can be denoted with “NA” or “.”. There are a few example CSV files included with the distribution.

Note that CSV files are extremely basic: they don’t allow the specification of color or glyph, for example, nor do they allow row ids, so they can’t be used for linking across data sets and the display of graphs.

## 5 View modes: projections

Selecting any item on the **View** menu changes the projection of the current display. This causes two changes: the control panel for the new projection appears to the left of the console window, and defines the appearance and behavior of the variable selection panel.

As an example, start GGobi with some data, and watch what happens in the main GGobi window and a single scatterplot display. When GGobi starts, it's in **XYPlot** mode by default, so the scatterplot window shows a 2-dimensional projection of the data.

Now select **2D Tour** in the View menu.

- The control panel at the left of the console changes, because the **2D Tour** has its own set of interactions.
- The variable selection panel changes, because the variable selection behavior for high-dimensional projection types is quite different than that for low-dimensional projection types.
- The plot in the scatterplot display changes, because it's now showing a projection of 3 variables instead of 2. Furthermore, it's moving, because a grand tour process is running.

Now we'll describe each view mode in more detail, starting at the top of the View menu.

### 5.1 1D plots

The 1D plot can be displayed in two ways: as a textured dot plot or an average shifted histogram, or ASH.

The textured dot plot uses a method described in [17]. This method spreads the data laterally by amounts that are partly constrained and partly random, resulting in a fairly smooth spreading of the points and minimizing artifacts of the plotting method, such as stripes, clusters, or gaps.

The ASH is due to Scott ([12]), and the code is also his. In this method, several histograms are calculated using the same bin width but different origins, and the averaged results are plotted. His algorithm has two key parameters: the number of bins, which controls the bin width, and the number of histograms to be computed. In GGobi, the number of bins is held constant (at 200), while the smoothing parameter available on the console controls the number of histograms (which ranges from 1 to 50). The effect is a smoothed histogram – a histogram that allows us to retain case identity so that the plots can be linked case by case to other scatterplots.

Line segments can be added that run between the plotted point and the baseline of the ASH. This is helpful when the smoothing parameter is low, because it helps your eye make out the shape of the ASH.

The 1D plot will be arranged horizontally if you select a variable with a left click, and vertically if you use a middle or right click.

The cycling controls can be used to make GGobi step through the plots automatically, one after another.

To activate this view mode from the keyboard, type *d* or *D* with the focus in the plot window, or type Control-*d* or Control-*D* with the focus in the console.

## 5.2 XY plots

The XY plots are the rudimentary 2 variable scatterplot (or draughtsman plot) displays. Two variables are chosen, one to be plotted horizontally and the other vertically. The cycling controls can be used to make GGobi step automatically from one pairwise plot to the next.

To activate this view mode from the keyboard, type  $x$  or  $X$  with the focus in the plot window, or type Control- $x$  or Control- $X$  with the focus in the console.

## 5.3 1D Tour

The 1D tour generates a continuous sequence of 1-D projections of the active variable space. The projected data are displayed as an average shifted histogram (ASH), horizontally or vertically. The scrollbar at the top of the controls allows the speed of rotation to be adjusted. The pause checkbox stops and starts the tour. **Reinit** initializes the tour to the projection (the ASH) of the first active variable. **Scramble** sets the view to a random projection.

Variables can be toggled into and out of the subset of active variables by clicking on the the toggle buttons. The variable will be immediately removed from the tour.

The variable circles on the right hand side of the control panel add further control for adding or removing variables from the current tour. The active variable space is the subset of variables currently selected, and their variable circles are drawn with a bold outline. When a variable is de-selected on the variable circles the variable fades out gradually, to maintain continuity of motion.

The reason for the two displays is to make handling of large numbers of variables more convenient. It might not make sense to include all the variables into a tour at once. The toggle buttons provide an efficient way to interact with the variable list. The variable circles provide information about the variables in the tour, how they are project to give the current projection, and information on which variable is the current manip variable. It is possible to select and de-select these variables to fine tune the tour.

The variable projection coefficients can be manually manipulated using manual controls. To select the variable to manipulate, click on the purple **Manip** button and then click on the variable circle. Horizontal mouse motions in the plot window then alter the coefficient for the manipulation variable, constrained by the values of the coefficients of other active variables (which may also change).

The variable axes and projection coefficient values can be toggled on or off using the display's **1D Tour** menu.

To activate this view mode from the keyboard, type  $t$  or  $T$  with the focus in the plot window, or type Control- $t$  or Control- $T$  with the focus in the console.

### 5.3.1 Projection Pursuit

A guided tour is available when the **Projection Pursuit** button is selected. It is controlled through a separate pop-up projection pursuit window, which contains a plot of the projection pursuit index. When **Optimize** is selected, the tour is guided by the index rather than proceeding randomly. The numbers displayed to the right of the **PP index** label are the minimum,

current value, and maximum of the index. A selection of indices is available.

### 5.3.2 1D Tour Options

The display's **1D Tour** menu contains controls for laying out variable circles in the variable selection panel in different ways, and also for toggling variable fading on or off. Variable fading means a variable smoothly fades out when it is de-selected. The alternative is to zero the variable out of view immediately, which creates a discontinuity in the tour motion, but is desirable for some situations.

## 5.4 2D Tour

The 2D tour generates a continuous sequence of 2-D projections of the active variable space. The projected data are displayed as a scatterplot. The scrollbar at the top of the controls allows the speed of rotation to be adjusted. The pause checkbox stops and starts the tour. **Reinit** initializes the tour to the projection of the first two active variables. **Scramble** sets the view to a random projection.

Variables can be toggled into and out of the subset of active variables by clicking on the the toggle buttons. The variable will be immediately removed from the tour.

The variable circles on the right hand side of the control panel add further control for adding or removing variables from the current tour. The active variable space is the subset of variables currently selected, and their variable circles are drawn with a bold outline. When a variable is de-selected on the variable circles the variable fades out gradually, to maintain continuity of motion.

The variable projection coefficients can be manually manipulated using manual controls. To select the variable to manipulate, click on the purple **Manip** button and then click on the variable circle. Once a manipulation mode has been selected, horizontal mouse motions in the plot window alter the coefficient for the manipulation variable, constrained by the values of the coefficients of other active variables (which may also change).

There are 5 manipulation modes: *oblique* allows unconstrained manipulation, *horizontal* and *vertical* constrain manipulation along the axes, *radial* constrains manipulation to the current direction of the variable keeping angle fixed, and *angular* manipulation allows rotating the variable axis in the plane of the plot window, keeping the length of the axis fixed.

The variable axes can be toggled on or off using the **2D Tour** menu at the top of the display window.

To activate this view mode from the keyboard, type *g* or *G* with the focus in the plot window, or type Control-*g* or Control-*G* with the focus in the console.

### 5.4.1 Projection Pursuit

A guided tour is available when the **Projection Pursuit** button is selected. It is controlled through a separate pop-up projection pursuit window, which contains a plot of the projection pursuit index. When **Optimize** is selected, the tour is guided by the index rather than pro-

ceeding randomly. The numbers displayed to the right of the **PP index** label are the minimum, current value, and maximum of the index. A selection of indices is available.

Often, sphering the data ahead of time provides more interesting results with the 2D guided tour, especially for the holes and central mass indices. Use the **Tools** menu and use the **Sphering...** tool to clone sphered counterparts of the currently active variables to do this.

#### 5.4.2 2D Tour Options

The display's **2D Tour** menu contains controls for laying out variable circles in the variable selection panel in different ways, and also for toggling variable fading on or off. Variable fading means a variable smoothly fades out when it is de-selected. The alternative is to zero the variable out of view immediately, which creates a discontinuity in the tour motion, but is desirable for some situations.

### 5.5 Rotation: 2D Tour with Three Variables

The rotation mode is essentially a 2D tour that is restricted to use three variables. Its graphical user interface is a subset of the 2D tour interface, with the exception that the three axes are individually represented by toggle buttons labelled X, Y and Z, principally so that it's possible to unambiguously specify the variable to be replaced when selecting a new one.

### 5.6 2x1D Tour

The 1x1D tour generates 2 independent continuous sequences of 1D projections of 2 active variable spaces, plotting the results horizontally and vertically generating a scatterplot. The scrollbar at the top of the controls allows the speed of rotation to be adjusted. The pause checkbox stops and starts the tour. **Reinit** initializes the tour to the projection of the first two active variables. **Scramble** sets the view to a random projection.

Variables can be toggled into the tour by clicking on the variable circles. A click with the left mouse toggles a variable in the horizontal direction, and a click with the middle mouse toggles a variable in the vertical direction. The active variable space is the subset of variables currently selected, and their variables circles are drawn with a bold outline. When a variable is toggled out of the tour it fades out gradually, to maintain continuity of motion.

The variable projection coefficients can be manually manipulated using manual controls. To select the variable to manipulate, click on the purple **Manip** button and then click left or right on the variable circle. Once a manipulation mode has been selected, mouse motions in the plot window alter the coefficients for the manipulation variable or variables, constrained by the values of the coefficients of other active variables (which may also change).

There are 4 manipulation modes: *combined* changes both horizontal and vertical manipulation variable coefficients, *equal combined* constrains the horizontal and vertical changes to be equal, *horizontal* and *vertical* constrain manipulation in the corresponding direction.

The **2x1D Tour** menu on the menu bar in the display window contains controls for laying out variable circles in the variable selection panel in different ways, and also for toggling variable fading on or off. Variable fading means a variable smoothly fades out when it is de-selected.

The alternative is to zero the variable out of view immediately, which creates a discontinuity in the tour motion, but is desirable for some situations.

The variable axes can be toggled on or off using the same menu.

To activate this view mode from the keyboard, type *c* or *C* with the focus in the plot window, or type Control-*c* or Control-*C* with the focus in the console.

## 6 Interaction modes

Selecting an item on the **Interaction** menu changes the interactions available for manipulating the current display. These new interactions are visible in the control panel that appears to the left of the console window. (There may also be cues added to the plot to tell you the interaction of the plot.

### 6.1 Scaling of axes

You can think of view scaling as if you're operating a camera and looking at a projection of the data in the viewfinder. That is, you aren't transforming the data itself, just your view of it. There are two ways to perform view scaling: by using the sliders on the console or by direct manipulation in the plot window.

To activate this interaction from the keyboard, type *s* or *S* with the focus in the plot window, or type Control-*s* or Control-*S* with the focus in the console.

#### 6.1.1 Sliders

There are two Zoom sliders, one for zooming in and out horizontally and the other for zooming vertically. Drag a slider, or click in the trough, to zoom. To manipulate both sliders at once, select **Fixed aspect**.

#### 6.1.2 Direct manipulation

Use the left button to pan move the data freely around the window, and the middle or right button to zoom. Once again, selecting **Fixed aspect** will hold the aspect ratio constant during zooming.

#### 6.1.3 Reset

To reset the plot, use the **Scale** menu in the display's menubar: it has two entries, allowing pan and zoom to be reset separately.

## 6.2 Brush: brushing of points and edges

Brushing is often performed when only a single display is visible, but it is most interesting and useful to perform brushing with more than one linked display showing different views of the

same data.

To interactively paint points, drag left to move the “brush” within the plotting window, or drag middle to change the size or shape of the brush while you paint. (If you lose the brush by pulling it outside the plotting window, you can grab it again if you press the left or middle button while the cursor is inside the display window.)

To activate this interaction from the keyboard, type *b* or *B* with the focus in the plot window, or type Control-*b* or Control-*B* with the focus in the console.

### 6.2.1 Point and edge brushing

If **Point brushing** is in any state but *Off*, the brush has a rectangular outline. As the brush is moved across the points, any points contained by the brush are affected. You may be changing the color, glyph shape, glyph size, or the “visible” state of the point, depending on the menu setting. The brushing style called “Shadow” deserves a few words: When it is selected, the brushed points are drawn in a color that’s very close to the background color [1]. The points are de-emphasized but they provide context for the rest of the data.

Similarly, if **Edge brushing** is in any state but **Off**, the brush includes a crosshair, and as the brush is moved in the window, any edges (line segments) intersecting either the vertical or horizontal “hair” are affected. You may be changing the color, line type, line thickness, or visible state of the edge.

If both **Point brushing** and **Edge brushing** are on, the brush is drawn as a crosshair inside a rectangle, and both point and edge brushing are performed.

### 6.2.2 Persistent brushing

If **Persistent** is selected, then a brushed point or edge retains its new characteristics after the brush has moved on.

### 6.2.3 Undo

Clicking on the **Undo** button restores the characteristics of all points and edges persistently painted between the last mouse-down and mouse-up.

### 6.2.4 Choose color & glyph ...

Clicking on this button opens the **Choose color & glyph** panel, which can be used to choose the point color and glyph as well as the edge type for brushing. At the top of the panel, there is a table of all possible point glyphs and another table of all possible edge types. Clicking on a point glyph sets both the glyph and edge type.

The reason that glyphs and edge types are linked is that points in one display may be linked to edges in another, and then brushing a point with a new glyph may cause a linked edge to acquire a new edge type at the same time. (Since there are more point types than edge types for now, it’s clear which edge type to select if a new glyph is chosen, but it isn’t clear which glyph to

select if a new edge type is chosen. For that reason, the edge symbols don't yet respond to button clicks.)

Below those tables is a row of rectangles of color which represent the current color scale. Clicking on one of these sets the brushing color. Double-clicking on one of these rectangles, or on the two rectangles just below them for the background and accent colors, opens a color selection widget with access to the full color map. The **Reverse video** button allows you to swap the background and accent colors.

### 6.2.5 Linking rules

This list is used to define which of two linking rules is to be used. The default rule, linking by case identifier, dictates that points representing records that have the same *id*, as specified in the XML description (or in the API), will respond identically to brushing events. Ids are unique within a dataset, so this rule has no effect when only a single dataset is being studied.

If instead the data includes categorical variables, and one of those is selected, the linking logic uses the levels of that variable to link points. When a case is brushed in one display, all cases with the same value of the categorical variable will change accordingly, in this and all other displays.

For example, look at the *algal-bloom.xml* data supplied with GGobi. It contains four datasets, one of which contains the levels of a factorial experimental, while another represents the response. Both include the same categorical variables. Open two scatterplots, one for each of those two datasets. In the measurements display, plot algal count against day. In the plot of experimental conditions, plot the level of phosphorus against the level of carbon. Prepare to brush in the plot of experimental conditions. Choose **Link by Carbon** as the linking rule: Note that you have to select brushing and set the linking rule twice, once when each scatterplot is the current plot. Highlighting the low values of carbon, note that all of the points highlighted in the measurements display are among the lowest in algal count.

### 6.2.6 Tools of importance for brushing

See section 7.5 for a description of the selection of color schemes, section 7.6 for a description of **Automatic Brushing** and section 7.7 for a description of a tool for showing and hiding clusters.

### 6.2.7 Brush menu in display

When the brushing mode is active, the **Brush** menu in the display's menu bar contains these items:

- Exclude shadowed points: Excluded points aren't drawn, and the views are scaled without them. Excluding a lot of points from large data sets can improve the performance of many operations.
- Include shadowed points: Redraw these points as shadows, and include them in view scaling.

- Un-shadow all points: Restore the points to their usual colors.
- Exclude shadowed edges: As is the case with points, excluding a lot of edges from large data sets can improve performance.
- Include shadowed edges: Redraw these edges as shadows.
- Un-shadow all edges: Restore the edges to their usual colors.
- Reset brush size: Reset the brush to its default size and position.
- Brush on: When this is not checked, the brush can be freely moved across the plotting window and it does not change the points. This is useful if you are plotting a very large number of points, and you want to position the brush before painting, because you can move it much more quickly across the plot.
- Update brushing continuously: Update linked brushing with every mouse motion. The alternative is to update linked views only when the mouse is released, which is more efficient when there are a great many points in the plot, or a great many plots on the screen.

The first few items in that list may affect more than the current display. Because they really operate on the *data* in the current display, all other displays showing the same data will also respond.

### 6.3 Identification

This mode is used to display labels near points in the plotting window. To see these labels, simply move the cursor inside the plotting window. The label of the point nearest the cursor is displayed. The possible labels are

- the record (case) label supplied by the user either in ASCII or XML (the default),
- the record number (the backup default in case no record label was supplied),
- a list of variables and variable values, where the variables are specified in the list widget above, or
- the record id.

Identification in one window is instantly reflected in all linked windows. [Some thought is required before deciding how or whether this interaction mode should reflect the linking rules used in brushing.]

To cause a label to become “sticky,” click left when the target label is displayed. The printing style changes and the label now remains printed as the cursor moves off, and even remains printed as you leave the **Identify** mode. It is possible to rescale or rotate data, and the sticky labels will continue to be displayed next to their associated points.

To cause a label to become “unsticky,” return to the **Identify** mode and click left again when the target point is nearest the cursor. It is also possible to restore all labels to unsticky status

by clicking on the **Remove labels** button. You can also see all the labels at once by clicking on **Label all**.

Notice that once a point's label is sticky, you can click **Recenter** to make it the center for the rotation and tour modes.

To activate this interaction mode from the keyboard, type *i* or *I* with the focus in the plot window, or type Control-*i* or Control-*I* with the focus in the console.

## 6.4 Edit edges

Here we add points or edges to the datasets by adding them to the displays.

Adding edges: Press the mouse button when the cursor is near the source node, and drag it around the window. You'll see a temporary edge between the source and the nearest node, drawn using the current color and "glyph." When you release the button, one of two things will happen: If you pressed the middle or right mouse button, a dialog window will appear with default values describing the new edge; if you pressed the left button, the edge will simply be added with those same default values.

Adding points: Simply click a mouse button when the cursor is where you want the new point to be located. As above, the middle or right button raises a dialog, and the left button simply adds the point.

The default values that are assigned are the record label and record id (often the same), and the variable values (if the dataset being augmented has variables). By default, the record label and record id are simply the new record number (represented as a string). If this string already exists as a record id, the new record will not be added.

The default variable values are assigned based on where you clicked on the screen and on the current projection. For any variables not part of the current projection, the default value is 0.

Edge and point deletion have not been implemented. For now at least, you'll have to shadow brush any unwanted elements to get rid of them.

To activate this mode from the keyboard, type *e* or *E* with the focus in the plot window, or type Control-*e* or Control-*E* with the focus in the console.

Note: This mode and the next, "Move Points," may seem like peculiar, even dangerous, additions to data analysis software. They were initially added for the use of another community of xgobi users: discrete mathematicians use xgobi an GGobi to visualize graphs. In that context, moving and editing graphs is quite natural – as it sometimes turns out to be in the context of data analysis, too.

## 6.5 Move points

In this mode, points or groups of points can be moved. Move the cursor in the window until it's nearest to the point you want to move, then press any mouse button and drag until the point is where you want.

The **Direction of motion** menu allows the movement to be constrained. If the **Move brush group** checkbox is checked, then all points with the same glyph and color as the selected point

will be moved with it.

The **Undo last** and **Reset all** buttons allow movement to be reversed.

To activate this interaction mode from the keyboard, type *m* or *M* with the focus in the plot window, or type Control-*m* or Control-*M* with the focus in the console.

## 7 Tools

### 7.1 Variable manipulation tool

This powerful tool is opened by selecting the first entry on the **Tools** menu. It has several important functions:

- the display of variable statistics,
- variable subset selection for launching multi-plot displays,
- setting variable ranges,
- cloning variables, and
- adding other new variables.

Its first purpose is to report information about each variable. It begins by separating the variables by type: variables are currently classified as categorical or real, though more types can be specified in XML. Categorical variables are displayed hierarchically, and the information reported includes the number of records for each level. For “real” variables, GGobi report the current variable transformation (if any); the minimum, maximum, mean, and median of the raw data; the number of missing values per variable.

Its second purpose is to specify subsets of variables to be plotted when launching a parallel coordinates or scatterplot matrix display. Variables are selected by highlighting rows, and the control and shift keys are modifiers that allow multiple rows to be highlighted.

These selected variables will also respond to operations contained within the panel: you can reset the variable ranges that are used for projecting the data into the plotting window. This allows variables with the same units or potential range (such as percentages) to use the same range, and facilitates visual comparisons.

The selected variables can also be cloned, and the new variables you create will be added to the table as well as the console.

There’s another way to add new variables, and that relies on the **New ...** button, which brings up a small panel. Use that panel to specify the variable’s name and to set its values: either the row numbers or a set of integers reflecting the assignment of a group identifier to each combination of point color and glyph.

### 7.2 Variable transformation tool

The first step in variable transformation is to specify the variables you want to transform.

There are three stages in the transformation pipeline, with a transformation function in each stage operating on the output of the previous stage. It’s equally acceptable to use any or all of them.

You can think of stage 0 as a domain adjustment stage: if a variable has negative values, for instance, many transformation functions can’t be applied to it, so you may need to add an increment to each value.

Stage 1 transformations include the Box-Cox family of linear transformations  $T(X) = (X^\lambda - 1)/\lambda$  [2], and you can either type the Box-Cox parameter into the text box and hit return, or use the spin button to gradually increase or decrease the parameter.

Many of the stage 2 transformations are not linear; they include sorting and ranking.

### 7.3 Sphering

To sphere one or more variables, first select the variables in the list at the top of the window, then click on **Update scree plot**.

It's common to standardize the variables before sphering, that is, use the correlation matrix instead of the variance-covariance matrix. The check box **Use correlation matrix** allows for this option.

Now you're ready to sphere the selected variables. Working your way down the panel, use your visual interpretation of the scree plot together with the information in the labeled section "Prepare to sphere" to decide how many principal components you want to create. By default, all the selected variables will be sphered, but you decide that the first few principal components account for a sufficiently high proportion of the variance. In that case, you can use the spin button to the right of the label "Set number of PCs" to decrease the number of principal components you're going to generate. The variance and condition number are displayed to help you make that choice.

Once you're satisfied with the selected variables and the number of principal components, proceed to the last step. Click on **Apply sphering** to create new variables and add them GGobi's variable selection panels. The names of the selected variables will be added to the "sphered variables" to help you remember which variables you sphered.

### 7.4 Jittering

Select **Variable jittering ...** to open a panel that allows random noise to be added to selected variables. This ameliorates overplotting in scatterplots of data with many ties.

First specify the variables you wish to jitter. Choose between uniform and normal random jitter, and then set the degree of jitter using the slider. To rejitter without changing the degree of jitter, simply click on the **Jitter** button.

### 7.5 Color schemes

The **Color schemes** tool has two purposes: to select a new color scheme, and to automatically color points using the current color scheme.

Preview a new color scheme by selecting from the tree at the left. If you would like to apply it, click on the button "Apply color scheme to brushing colors." That replaces the current set of colors visible in the "Choose color & glyph" menu and used in all the displays.

If you are currently using  $n$  colors and the color scheme you have selected has fewer than  $n$  colors, you'll be prohibited from applying the new scheme. If you want to use that scheme, you must use brushing to reduce of colors in use.

The sample file presently contains 250 or so color schemes, of different types and sizes, and they're based on the work of Brewer ([3], [www.personal.psu.edu/cab38](http://www.personal.psu.edu/cab38)). The four types represented are

- diverging: used when the range of the coloring variable has a meaningful midpoint;
- sequential: used to highlight a continuous progression of values;
- qualitative: used when the coloring variable is categorical;
- spectral: Brewer's modifications of the popular spectral scale to reduce its drawbacks. She has made the perceptual steps more uniform, and made the scales more friendly for people with color vision impairments.

## 7.6 Automatic brushing

To paint the data according to the values of a variable, first select a variable in the list at the top of the tool window. This adds two sets of numbers to the display: along the bottom of the display, at the center of each stripe of color, is shown the number of points that will be drawn with this color once **Apply** is pressed. Along the top of the display, at the boundaries between the stripes of color, appear the values of the chosen variable that define the boundaries between colors.

There are two methods available for defining the bin boundaries, and a menu for choosing between them. The "constant bin width" method simply partitions the range of the selected variable into equal-sized sub-ranges, and maps points into those bins. The "constant bin count" method attempts to map the values into  $n$  bins of equal size. Since it also tries to assign all equal values into the same bin, it usually doesn't produce uniform bins if the variable has many equal values.

To adjust the boundaries between stripes of color, grab one of the sliders, and notice that both the values and the counts adjust as you move it. The displays will respond as the sliders are moved: either continuously, or only when you release the mouse. If you your data set has a large number of cases, continuous updating will lag behind the mouse, so it is probably more effective to update only on mouse release.

Try it with the olive oil data used in the tutorial, and select the *Area* variable. It can be a real time-saver.

It's not clear that this tool knows the best way to handle categorical variables yet.

## 7.7 Color & glyph groups

The **Color & glyph groups** tool displays a table. Each unique symbol in the data (combination of color and glyph) occupies a row, and for each row has a toggle button that shadow brushes all points drawn in the corresponding symbol. The three remaining columns report the number of cases shadow brushed, the rest of the cases, and the total number of cases in the group.

Clicking on one of the little symbol displays will assign the currently selected color and/or glyph (depending on the state of the Point Brushing menu on the Brush control panel) to all the points

in the corresponding cluster. If your selection would collapse two clusters into one, it will refuse to go ahead: it seems highly likely that someone might do that by mistake, and unlikely that they would do it deliberately. (That choice should be resolved by a dialog, probably, and an 'undo' button should be added.)

The **Exclude shadows** button at the bottom of the window will exclude all shadow-brushed points from consideration in the displays: the views will be scaled without those points, and they won't be drawn. (Excluding a lot of points from large data sets can improve the performance of many operations.) The **Include shadows** button will bring those points back into the plot, and they'll be drawn in the shadow color as before.

The **Update** button updates the contents of the table in case it isn't responding properly to changes in the displays as you continue to brush.

See [7.1](#) to read how you can add a new variable to the existing data which serves as an indicator for these "clusters."

## 7.8 Subsetting

Select **Case subsetting and sampling ...** to open a panel that allows subsets to be specified in one of six ways.

- **Random sample without replacement:** Specify the number of cases to be in the sample.
- **Consecutive block:** Specify the first and last row of the block. (The two controls in the second row control the increment used in the control in the first row.)
- **Limits:** Use the limits defined by the user in the variable manipulation table to define the subset.
- **Every nth case:** Specify the interval and the first row.
- **Sticky labels:** All cases with a "sticky" label will be in the subset. If no points have a sticky label, this will have no effect. (See [section 6.3](#) for a description of sticky labels.)
- **Row labels:** Type in a string, and specify where it should fall (or not fall) in the row labels, and whether case should be ignored. Matches will be included in the subset.

Select one of those six, then click on **Subset** in the bottom row of the panel. If you want to re-include all rows, select the **Include all** button in the bottom row.

Click the **Rescale** button to rescale all plots excluding the points not in the subset.

One purpose of subsetting is to allow the use of GGobi on data matrices that are so large that dynamic and interactive operations begin to become painfully slow. By selecting a smaller subset, a user can work on that subset at a comfortable speed of tour motion and interaction. Another purpose is to do graphical cross-validation: if the feature you see is still there in repeated subsamples, there's a good chance it's not just an artifact of visualization.

## 7.9 Controls for missing data

When your data includes missing values, you can add a new dataset to the current GGobi using the button at the top of the **Missing data** tool. The resulting dataset has the same number of rows as the original, but it may have fewer variables, since it only uses those variables which actually have missing values. It has the same variable names, row labels, glyphs and colors as the original. The difference is this: if the original data in position  $i, j$  is missing, the new dataset's value is 1.0; otherwise it's 0.0.

Once the new dataset has been created, plots of missingness information can be launched and linked to plots of the data. The missings plots are pre-jittered to spread the points; the degree of jittering can be later adjusted with the jittering tool.

Linking missings plots to displays of the data allows us to explore the joint distribution of missing values across variables.

### 7.9.1 Imputation

By default, missing values are assigned a value 10% lower than the minimum non-missing value for that variable, but you may sometimes find that to be an inconvenient choice. Use the **Missing values** panel to assign alternative numbers.

Select the variable or variables whose imputed value you wish to change, then select a method and fill in the text window if necessary, and click **Impute**.

Below the list of variables, a notebook widget allows you to specify the type of imputation you'd like to use.

**Random imputation:** Sample from the present values for each variable to populate the missing values. If you have done some brushing to partition the cases, you can specify that you want the sampling to be done using only cases brushed with the same color and glyph.

**Fixed value:** Specify any value to use instead of the default.

**Percentage below minimum:** Specify a value that is  $x$  percent below the minimum value. For example, if the variable ranges from 40 to 80, specifying 10 will assign the missings the value  $40 - (10\% * 40) = 36$ .

**Percentage above minimum:** Specify a value that is  $x$  percent above the maximum.

Click the resale button when you want to update the view to respond to the new values.

For more complex imputation, consider using `rggobi`

## 8 Multiple datasets

Several datasets can be open in GGobi simultaneously. The **File** menu interfaces reading in additional data sets. Data tabs corresponding to each open data set appear above the variable selection window on the main controls window, and in various tool windows. The rules for linking these data sets is described in Section 6.2.5.

Several of the sample datasets included with GGobi use multiple datasets. In most cases, the additional datasets are used to add edges, but *algal-bloom.xml* contains four datasets, one of

which contains the levels of a factorial experimental, while another represents the response. These two can be linked by variable.

## 9 Edges

A special case of the use of multiple datasets is use of edges (line segments). There are many reasons one would want to display line segments in a scatterplot.

There are several datasets with edges distributed with GGobi: *algal-bloom.xml*, in which edges are used to structure the display of an analysis of variance; *buckyball.xml*, data describing a graph – a geometrical object; *eies.xml*, social network data; *pigs.xml*, a dataset which includes several time variables; *prim7.xml*, in which line segments are used to illustrate structure in the data which was found during extensive exploratory data analysis.

We specify line segments (edges) for GGobi by the addition of a second dataset in an XML file (or through the API) in which each record has tags for *source* and *destination*. These edge datasets can still have variables, of course, which might represent variables measured on transactions or interactions.

A scatterplot which has corresponding edge datasets (which we might call edge sets) will have an **Edges** menu in its menubar. If there's a choice of edge sets, you'll see a cascading menu showing their names. Edges can have "arrowheads" added, to indicate the edge's direction; it's also possible to see the arrowheads alone.

Edges can be brushed directly, as described in [6.2](#).

If the records that define each edge also have variables, then displays of the variables in those edge datasets are also possible. A point in the scatterplot of an edge dataset corresponds to the same record as an edge in another scatterplot. So brushing an edge in one is the virtually the same action as brushing a point in another.

Two plugins, GraphLayout and ggvis, offer methods for laying out graphs. They are documented elsewhere.

## 10 Large data

There exist at least two meanings of "large" in data analysis: large  $N$  (number of cases) and large  $P$  (number of variables).

### 10.1 Large $N$

We won't attempt to define "large," because it depends so much on your computing hardware and software, but we do know something about the sources of sluggishness, and some effective workarounds.

First of all, there are inefficiencies in the Windows implementation of `gdk`, the drawing library underlying the `gtk+` toolkit in which GGobi is written. Windows users should restrict themselves to rectangular glyphs or point glyphs to get the best possible performance. They certainly should avoid the use of circles. (We have tried to persuade the folks who port `gdk` to Windows to do

a better job, but we haven't yet managed to convince them it's important.) Even in dialects of UNIX, single-pixel points are drawn faster than other glyphs.

Use subsets when possible, particularly with the rotation and tour methods. Use the Tools-Case subsetting and sampling panel to select a sample of the data.

If the touring and rotation methods become too slow, pause the movement and use the **Scramble** button to see different views.

Avoid the 1-D plotting methods in scatterplots; they both become slow as the data gets very large. The barchart/histogram display is of course much less affected by  $N$ , and is probably easier to read because it isn't affected by overplotting.

There are a few ways to improve brushing performance, most of which involve postponing linked updates.

- In Brush mode, open the Reset menu in the main menubar and turn off "Update brushing continuously." In that case, linked windows will only be redrawn when you lift your finger off the mouse button. Sometimes you can go even further and perform all your brushing operations with only a single display open, and open the other displays afterwards.
- Turn off the brush until you have it shaped and positioned.
- Use the API instead of the GUI. Using the `rggobi` package, for example, you can brush a region in a single command if you can describe it in terms of variable values or indices.
- Working with edges slows down the drawing routines, so you might remove edges during point brushing (or rotation), and then restore them.

## 11 Differences from XGobi

In this section, we summarize the key differences between GGobi and XGobi for those readers who are already familiar with XGobi.

### 11.1 Multiple displays

The first thing you'll notice when you look at a GGobi display is that the plotting window has become separated from the control panels. The main reason for that change is so that a GGobi process can have multiple display windows, of the same type or of different types. In addition to the basic scatterplot, GGobi currently has scatterplot matrices, parallel coordinates plots, and time series plots. (See [3.2](#) for more detail.)

This design change has had far-reaching effects.

First, user interactions available for the simple scatterplot display can now be made available for other display types. In `xgobi`, for instance, there is a parallel coordinates display, but it's not possible to brush it – in GGobi, it is.

Unfortunately, having multiple displays introduces a new source of ambiguity: you now have to tell GGobi which display, and which plot within a display, you want to address. Do that by simply clicking inside the target plot. GGobi will draw a thick white outline around that plot so that you can check which plot your actions will be addressing. The interaction mode control panel in the GGobi console should always correspond to the state of the current display, too. We need more user experience before we can tell whether this approach is satisfactory.

The basis for linking has changed, too. All displays of the same data are now linked by default, so it's no longer necessary to run multiple processes in order to achieve linked displays.

One of the most interesting implications of using multiple displays is that a GGobi process is no longer restricted to a single data set. The XML file format makes it easy to specify two or more data matrices in a single file, as described in [section 4.1](#). In addition, it's possible to add data matrices using the Read button on the File menu or using R. (The R - GGobi interface will be introduced below, and it's more fully described in [section ??](#) and in other documentation.)

The rules for linking in XGobi had evolved into a rather complicated hodge-podge with special handling of “row groups,” the “nlinkable” notion to exclude points from linking, and linking points to line segments (edges). This has been replaced with a single set of rules that can be specified in the XML file. See [section 6.2.5](#) for details.

With multiple displays, too, we no longer have to launch a new process to open missing value plots – the plots of 1's and 0's which represent the presence and absence of data in each cell.

A convenient side effect of multiple displays is that, since each display now sits in a window of its own, it's now very simple to adjust the aspect ratio of a plot; this simple operation is very awkward in XGobi.

### 11.2 Data format

Before we describe other key changes in the visible design, we'll introduce the changes in data format. The XGobi format, in which a set of files with a common base name is used, is no

longer supported. (See section 4 for more detail.) The functions served by the discarded file types are usually served now in a different way: for example, XGobi uses the *.vgroups* file to force a group of variables to use the same axis ranges. In GGobi, that's accomplished by setting limits in the variable manipulation tool.

The new format, in which all the data lives in one file, is written in XML. XML (Extensible Markup Language) is a widely used language for specifying structured documents and data to be viewed and exchanged. It was initially intended to be read by browsers, but it is also used to define documents that are read by other software. XML files can be validated automatically, and XML specifications can be easily extended, too, by adding to the set of tags in use.

We strongly favor the XML format, and we won't try to keep the old format up to date. For instance, it's no longer possible to specify edges in the *ascii* data format, but they can be specified in XML – and they can have associated data values.

Several XML data files are included in the sample GGobi data, and the details of their format is described in *The GGobi XML Input Format* (XML.pdf) which is included as part of the GGobi distribution.

### 11.3 Integration

Many xgobi users are also users of the SPlus or R statistics software, and have used the S function which launches an xgobi process viewing S data. Once that launch occurred, the resulting process was utterly independent of its parent. The XGobi authors did some experiments in the early 90s to achieve a more intimate connection, but made little headway with S, though interprocess communication was used successfully with ArcView [15, 16].

With GGobi, that problem has been solved, and it's now possible to have real-time integration between GGobi and a variety of other software environments. An example of this integration is the embedding of GGobi into R, with the addition of a set of R functions that manipulate GGobi data and displays.

For more details, see section ??.

### 11.4 Variable selection

There have been a few changes in variable selection. The familiar variable circles are still used for high-dimensional projections, but we've switched to a simple checkbox interface for plots where only one or two variables can be selected simultaneously. In these plots, the rich feedback provided by the variable circles is not needed, and may just be confusing to novices.

The basics of the user interface haven't changed, though: click with the left button to select a variable to be plotted horizontally and the middle (or right) button to plot vertically.

### 11.5 The variable manipulation tool

The **Variable manipulation** tool presents summary statistics for each variable. If a variable has been described in the XML data file as categorical, it also shows the name and count for each level. The tables can also be used to select variables – not for plotting, but for specifying variables subsets when launching a parallel coordinates or scatterplot matrix display.

This tool can also be used to specify scaling ranges for variables or groups of variables, and so we have dispensed with the *.vgroups* functionality in XGobi.

Variable cloning is another new feature: it appears as a button on the variable selection and statistics table.

For more detail on this table, see section 7.1.

## 11.6 Changes in projection, interaction and tools

**Brush** may be the mode that has changed the most, because GGobi has a much richer notion of color selection than XGobi. Open the **Choose color & glyph** panel, and then double-click on any element of the color palette to bring up a color selection widget with access to the full color map. Notice that you can change the background color as well as any of the foreground colors, and notice that changing the glyph also changes the line type to be used in edge brushing.

The **Color schemes** tool has extended brushing considerably, by allowing a choice of color schemes and enabling a form of automatic brushing. See 7.5.

**Scale** has also changed. The direct manipulation shifting and scaling methods work as they did in XGobi, but we've added what we call "Click-style interaction" for more precise control. See section 6.1.

The tour methods in GGobi are still evolving: many things are not yet implemented, but new things are beginning to appear. In the **1DTour**, a projection of several variables is viewed as an average shifted histogram, as described in section 5.3.

The redesign of the tour methods is reflected in the **Sphering** tool, which also appears in the most recent versions of XGobi. Instead of automatically sphering variables in projection pursuit, GGobi requires you to specify which variables to sphere. Since this method makes use of variable cloning, it creates new variables, allowing you to look at plots of principal components against the original data.

See section 7.3 for details.

## 11.7 On-line help

The on-line help system used in XGobi has been replaced with "tooltips," so leaving the mouse over a widget for a couple of seconds brings up a phrase describing the function of that widget. If the tooltips annoy you, you can turn them off using a checkbox on the **File-;Options** menu.

## Web Links

The GGobi web site, <http://www.ggobi.org>, contains details on downloading and installing software, related documentation, a picture and video gallery.

## References

- [1] Richard A. Becker and William S. Cleveland. Brushing scatterplots. *Technometrics*, 29:127–142, 1987.
- [2] George E P Box and David R Cox. An analysis of transformations. *Journal of the Royal Statistical Society*, B-26:211–243, 1964.
- [3] Cynthia A. Brewer. Color use guidelines for data representation. In *Proceedings of the Section on Statistical Graphics*, pages 55–60, Baltimore, 1999. American Statistical Association.
- [4] A. Buja, D. Asimov, C. Hurley, and J. A. McDonald. Elements of a Viewing Pipeline for Data Analysis. In W. S. Cleveland and M. E. McGill, editors, *Dynamic Graphics for Statistics*, pages 277–308. Wadsworth, Monterey, CA, 1988.
- [5] A. Buja, D. Cook, D. Asimov, and C. Hurley. Dynamic Projections in High-Dimensional Visualization: Theory and Computational Methods. Technical report, AT&T Labs, Florham Park, NJ, 1997.
- [6] A. Buja, D. Cook, and D. Swayne. Interactive High-Dimensional Data Visualization. *Journal of Computational and Graphical Statistics*, 5(1):78–99, 1996. See also [www.research.att.com/~andreas/xgobi/heidel/](http://www.research.att.com/~andreas/xgobi/heidel/).
- [7] D. Cook and A. Buja. Manual Controls For High-Dimensional Data Projections. *Journal of Computational and Graphical Statistics*, 6(4):464–480, 1997. Also see [www.public.iastate.edu/~dicook/research/papers/manip.html](http://www.public.iastate.edu/~dicook/research/papers/manip.html).
- [8] D. Cook, A. Buja, and J. Cabrera. Projection Pursuit Indexes Based on Orthonormal Function Expansions. *Journal of Computational and Graphical Statistics*, 2(3):225–250, 1993.
- [9] D. Cook, A. Buja, J. Cabrera, and C. Hurley. Grand Tour and Projection Pursuit. *Journal of Computational and Graphical Statistics*, 4(3):155–172, 1995.
- [10] M. Forina, C. Armanino, S. Lanteri, and E. Tiscornia. Classification of olive oils from their fatty acid composition. In H. Martens and H. Russwurm Jr., editors, *Food Research and Data Analysis*, pages 189–214. Applied Science Publishers, London, 1983.
- [11] A. Inselberg. The Plane with Parallel Coordinates. *The Visual Computer*, 1:69–91, 1985.
- [12] David W Scott. *Multivariate Density Estimation: Theory, Practice, and Visualization*. John Wiley, New York, 1992.
- [13] D. Swayne and A. Buja. Missing Data in Interactive High-Dimensional Data Visualization. *Computational Statistics*, 13(1):15–26, 1998.
- [14] D. F. Swayne, D. Cook, and A. Buja. XGobi: Interactive Dynamic Graphics in the X Window System. *Journal of Computational and Graphical Statistics*, 7(1):113–130, 1998. See also [www.research.att.com/areas/stat/xgobi/](http://www.research.att.com/areas/stat/xgobi/).

- [15] Deborah F. Swayne, Andreas Buja, and Nancy Hubbell. XGobi meets S: Integrating software for data analysis. In *Computing Science and Statistics: Proceedings of the 23rd Symposium on the Interface*, pages 430–434, Fairfax Station, VA, 1991. Interface Foundation of North America, Inc.
- [16] J. Symanzik, D. Cook, N. Lewin-Koh, J. J. Majure, and I. Megretskaja. Linking ArcView 3.0 and XGobi: Insight Behind the Front End. *Journal of Computational and Graphical Statistics*, 9(3):470–490, 1999.
- [17] John Tukey and Paul Tukey. Strips displaying empirical distributions: I. textured dot strips. Bellcore Technical Memorandum, 1990.
- [18] E. Wegman. Hyperdimensional Data Analysis Using Parallel Coordinates. *Journal of American Statistics Association*, 85:664–675, 1990.