

SpRay

an R-based visual-analytics platform for large and
high-dimensional datasets

J. Heinrich¹ J. Dietzsch¹ D. Bartz² K. Nieselt¹

¹Center for Bioinformatics, University of Tübingen

²ICCAS/VCM, University of Leipzig

August 12, 2008

Outline

- 1 Introduction
- 2 SpRay
- 3 Discussion
- 4 Future Work

Data Sets Become Increasingly Large

High-Throughput techniques yield a huge amount of data

- Microarrays
- CT scanner
- Simulation data

Many data sets are high-dimensional

- Time series: 100 experiments, 5 replicates, 10000 oligos
- 10000 rows \times 500 columns = $5 \cdot 10^6$ data points

... and complex

- Heterogeneous data (categorical, metric)
- Invalid data (NA, NaN)

Knowledge Discovery Becomes Increasingly Difficult

Effects of Large and High-Dimensional Datasets for the Analysis

- Storage: obvious
- Speed: time to read, locate, compute, render, display the data
- Quality: errors, administration
- Complexity: more variables, more detail, special cases. . .
- Visualization: Dimensionality, Occlusion, Identification

Visual Analytics with R

Analytical Reasoning

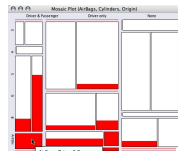
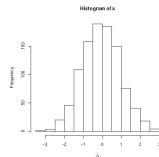
- Gain insight into data
- Reveal underlying structure and model
- Extract information contained

Techniques

- Data Analysis
- Visualization
- Interaction

```

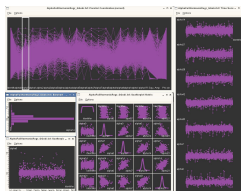
1 # Create a histogram for the variable 'mpg' and showing mean
2 # and standard deviation.
3 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
4 # The title is 'Histogram of mpg'.
5 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
6 # The histogram has 10 bins.
7 # The data is from the 'mtcars' dataset.
8 # The variable 'mpg' is selected.
9 # The histogram is plotted.
10 # The mean and standard deviation are printed.
11 # The plot is titled 'Histogram of mpg'.
12 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
13 # The title is 'Histogram of mpg'.
14 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
15 # The histogram has 10 bins.
16 # The data is from the 'mtcars' dataset.
17 # The variable 'mpg' is selected.
18 # The histogram is plotted.
19 # The mean and standard deviation are printed.
20 # The plot is titled 'Histogram of mpg'.
21 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
22 # The title is 'Histogram of mpg'.
23 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
24 # The histogram has 10 bins.
25 # The data is from the 'mtcars' dataset.
26 # The variable 'mpg' is selected.
27 # The histogram is plotted.
28 # The mean and standard deviation are printed.
29 # The plot is titled 'Histogram of mpg'.
30 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
31 # The title is 'Histogram of mpg'.
32 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
33 # The histogram has 10 bins.
34 # The data is from the 'mtcars' dataset.
35 # The variable 'mpg' is selected.
36 # The histogram is plotted.
37 # The mean and standard deviation are printed.
38 # The plot is titled 'Histogram of mpg'.
39 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
40 # The title is 'Histogram of mpg'.
41 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
42 # The histogram has 10 bins.
43 # The data is from the 'mtcars' dataset.
44 # The variable 'mpg' is selected.
45 # The histogram is plotted.
46 # The mean and standard deviation are printed.
47 # The plot is titled 'Histogram of mpg'.
48 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
49 # The title is 'Histogram of mpg'.
50 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
51 # The histogram has 10 bins.
52 # The data is from the 'mtcars' dataset.
53 # The variable 'mpg' is selected.
54 # The histogram is plotted.
55 # The mean and standard deviation are printed.
56 # The plot is titled 'Histogram of mpg'.
57 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
58 # The title is 'Histogram of mpg'.
59 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
60 # The histogram has 10 bins.
61 # The data is from the 'mtcars' dataset.
62 # The variable 'mpg' is selected.
63 # The histogram is plotted.
64 # The mean and standard deviation are printed.
65 # The plot is titled 'Histogram of mpg'.
66 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
67 # The title is 'Histogram of mpg'.
68 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
69 # The histogram has 10 bins.
70 # The data is from the 'mtcars' dataset.
71 # The variable 'mpg' is selected.
72 # The histogram is plotted.
73 # The mean and standard deviation are printed.
74 # The plot is titled 'Histogram of mpg'.
75 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
76 # The title is 'Histogram of mpg'.
77 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
78 # The histogram has 10 bins.
79 # The data is from the 'mtcars' dataset.
80 # The variable 'mpg' is selected.
81 # The histogram is plotted.
82 # The mean and standard deviation are printed.
83 # The plot is titled 'Histogram of mpg'.
84 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
85 # The title is 'Histogram of mpg'.
86 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
87 # The histogram has 10 bins.
88 # The data is from the 'mtcars' dataset.
89 # The variable 'mpg' is selected.
90 # The histogram is plotted.
91 # The mean and standard deviation are printed.
92 # The plot is titled 'Histogram of mpg'.
93 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
94 # The title is 'Histogram of mpg'.
95 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
96 # The histogram has 10 bins.
97 # The data is from the 'mtcars' dataset.
98 # The variable 'mpg' is selected.
99 # The histogram is plotted.
100 # The mean and standard deviation are printed.
101 # The plot is titled 'Histogram of mpg'.
102 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
103 # The title is 'Histogram of mpg'.
104 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
105 # The histogram has 10 bins.
106 # The data is from the 'mtcars' dataset.
107 # The variable 'mpg' is selected.
108 # The histogram is plotted.
109 # The mean and standard deviation are printed.
110 # The plot is titled 'Histogram of mpg'.
111 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
112 # The title is 'Histogram of mpg'.
113 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
114 # The histogram has 10 bins.
115 # The data is from the 'mtcars' dataset.
116 # The variable 'mpg' is selected.
117 # The histogram is plotted.
118 # The mean and standard deviation are printed.
119 # The plot is titled 'Histogram of mpg'.
120 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
121 # The title is 'Histogram of mpg'.
122 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
123 # The histogram has 10 bins.
124 # The data is from the 'mtcars' dataset.
125 # The variable 'mpg' is selected.
126 # The histogram is plotted.
127 # The mean and standard deviation are printed.
128 # The plot is titled 'Histogram of mpg'.
129 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
130 # The title is 'Histogram of mpg'.
131 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
132 # The histogram has 10 bins.
133 # The data is from the 'mtcars' dataset.
134 # The variable 'mpg' is selected.
135 # The histogram is plotted.
136 # The mean and standard deviation are printed.
137 # The plot is titled 'Histogram of mpg'.
138 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
139 # The title is 'Histogram of mpg'.
140 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
141 # The histogram has 10 bins.
142 # The data is from the 'mtcars' dataset.
143 # The variable 'mpg' is selected.
144 # The histogram is plotted.
145 # The mean and standard deviation are printed.
146 # The plot is titled 'Histogram of mpg'.
147 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
148 # The title is 'Histogram of mpg'.
149 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
150 # The histogram has 10 bins.
151 # The data is from the 'mtcars' dataset.
152 # The variable 'mpg' is selected.
153 # The histogram is plotted.
154 # The mean and standard deviation are printed.
155 # The plot is titled 'Histogram of mpg'.
156 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
157 # The title is 'Histogram of mpg'.
158 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
159 # The histogram has 10 bins.
160 # The data is from the 'mtcars' dataset.
161 # The variable 'mpg' is selected.
162 # The histogram is plotted.
163 # The mean and standard deviation are printed.
164 # The plot is titled 'Histogram of mpg'.
165 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
166 # The title is 'Histogram of mpg'.
167 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
168 # The histogram has 10 bins.
169 # The data is from the 'mtcars' dataset.
170 # The variable 'mpg' is selected.
171 # The histogram is plotted.
172 # The mean and standard deviation are printed.
173 # The plot is titled 'Histogram of mpg'.
174 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
175 # The title is 'Histogram of mpg'.
176 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
177 # The histogram has 10 bins.
178 # The data is from the 'mtcars' dataset.
179 # The variable 'mpg' is selected.
180 # The histogram is plotted.
181 # The mean and standard deviation are printed.
182 # The plot is titled 'Histogram of mpg'.
183 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
184 # The title is 'Histogram of mpg'.
185 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
186 # The histogram has 10 bins.
187 # The data is from the 'mtcars' dataset.
188 # The variable 'mpg' is selected.
189 # The histogram is plotted.
190 # The mean and standard deviation are printed.
191 # The plot is titled 'Histogram of mpg'.
192 # The x-axis is labeled 'mpg' and the y-axis is labeled 'Frequency'.
193 # The title is 'Histogram of mpg'.
194 # The x-axis ranges from -3 to 3 and the y-axis ranges from 0 to 100.
195 # The histogram has 10 bins.
196 # The data is from the 'mtcars' dataset.
197 # The variable 'mpg' is selected.
198 # The histogram is plotted.
199 # The mean and standard deviation are printed.
200 # The plot is titled 'Histogram of mpg'.
  
```



Visual Analytics with R

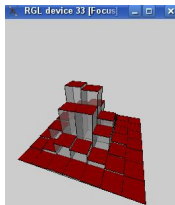
Related Work

GGobi¹



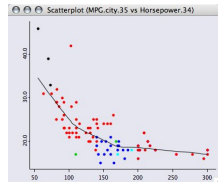
- linked views
- CPU only
- R optional

RGL²



- no linked views
- CPU/GPU
- depends on R

iPlots³



- linked views
- CPU/GPU
- depends on R

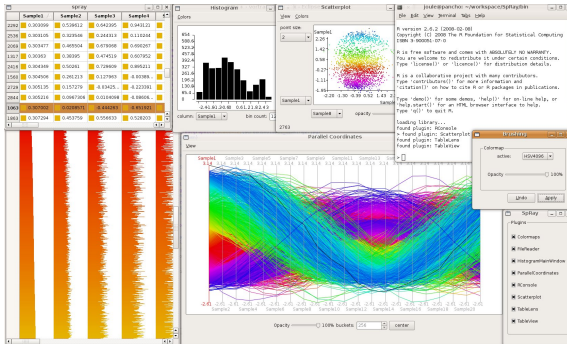
¹[Swayne et al., 2003]

²[Adler and Nenadic, 2003]

³[Urbanek and Theus, 2003]

SpRay

viSual exPlORation and anALYsis of high-dimensional data



- linked views
- CPU/GPU
- R optional

SpRay

Objectives

Objectives

- Extendable
- Interactive
- Portable
- Statistical Backend
- High-Performance



SpRay

Architecture

VisLib

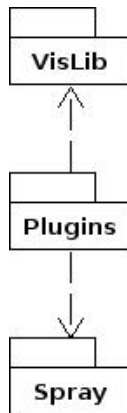
- Independent Visualization Library

Plugins

- Implement the plugin-interface
- Make use of VisLib (optional)

Host Application

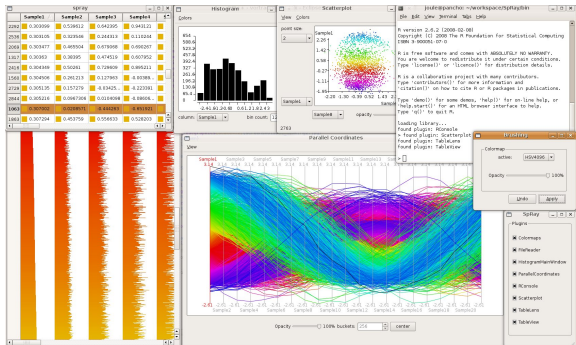
- Defines the plugin-interface
- Organizes communication



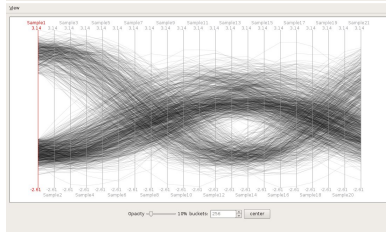
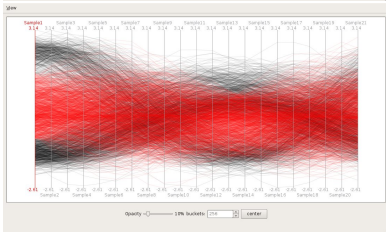
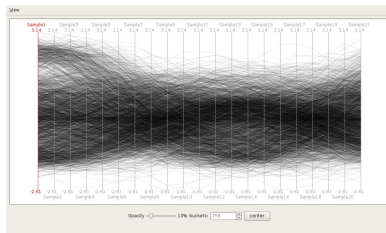
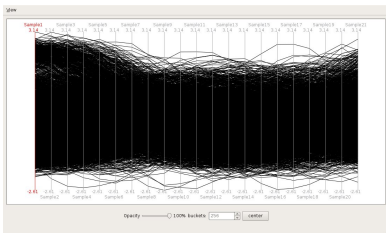
Plugins

Currently available

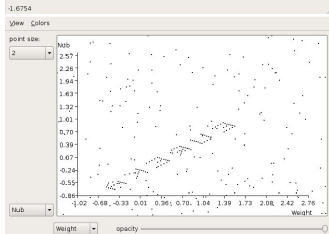
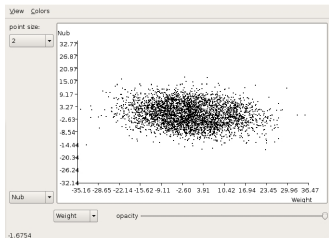
- Parallel Coordinates
- Scatterplot
- Histogram
- Data Table
- TableLens
- R-Console
- Brushing



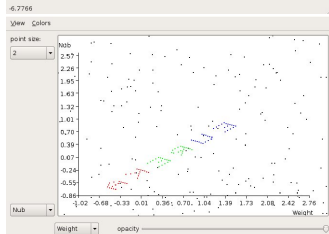
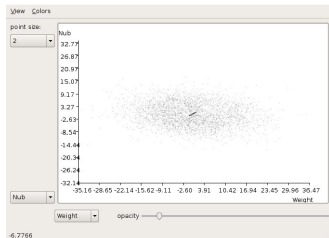
Parallel Coordinates



Scatterplot



0.8207



1.439

Data Table and R-Console

Data Table

	Sample1	Sample2	Sample3	Sample4	Sample5	Sample6
105	2.3415859163...	2.2376690546...	1.8982071874...	1.5555907577...	1.3721192240...	1.0038125142...
106	2.4401413565...	2.5040031651...	2.2656312029...	2.3950072936...	1.9896212244...	1.4697343178...
107	1.3811222629...	1.0465605603...	0.5781103422...	-0.233239960...	-0.621642629...	-0.807906840...
108	2.5384502012...	2.6598043399...	2.3529196256...	2.0231697166...	1.4417099692...	1.1541984958...
109	2.6865113599...	2.5672115712...	2.2875285175...	1.7987117521...	1.6185984467...	0.9827138999...
110	1.5442642534...	1.0196378953...	0.7877514254...	0.1382789851...	-0.403685698...	-0.692011846...
111	2.0488761740...	2.4978078321...	2.4989904982...	2.3638465612...	2.4162556514...	2.3486733609...
112	2.4515904970...	2.4102145324...	2.2785772702...	2.3462009364...	1.6576799128...	1.1253601772...
113	2.8350728203...	2.6358395454...	2.7894649536...	2.4301788917...	2.3818830086...	1.5301226056...
114	2.6457625081...	2.6191339249...	2.2891254100...	2.0749760720...	1.7882662302...	1.2335163649...
115	2.7380323543...	2.8372696299...	2.819595139...	2.4995965980...	2.0364326709...	1.6025039883...
116	2.2689389588...	2.3967752810...	2.2150238929...	1.8784653441...	1.7698521355...	1.3850019018...
117	2.0939331355...	2.1002774908...	1.6969150501...	1.2723382349...	0.7364192866...	0.3901431583...
118	1.9801059987...	1.6898704150...	1.5512342968...	0.9184525149...	0.6300543819...	0.1504721183...
119	2.0147595769...	2.1199225332...	2.2438563815...	2.3026739011...	2.2815192030...	2.0488097932...
120	1.7817203160...	2.0466779923...	2.1316428308...	2.2596791103...	2.2426233845...	2.1510957519...
121	2.6658253103...	2.5174836628...	2.3540442541...	1.8444133751...	1.2396569553...	0.8281764888...
122	2.3363089591...	2.0727448585...	1.6733367610...	1.5319254814...	1.1105856698...	0.6047233089...
123	1.6921383910...	2.0874461781...	2.4729319179...	2.6812288408...	2.3453695767...	2.5820502901...

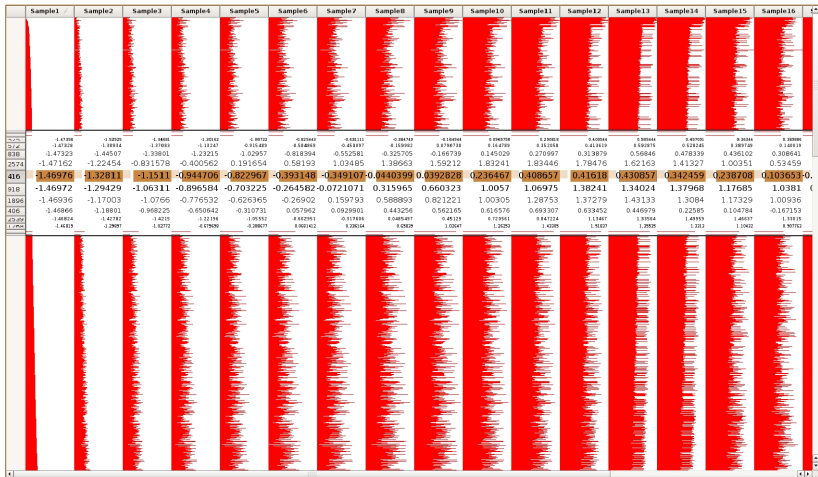
R-Console

The R-Console window displays the following R code and output:

```
Sitzung Bearbeiten Ansicht Le  
[[261]] *2761* *2762* *2763*  
0*  
[[2771]] *2771* *2772* *2773*  
0*  
[[2781]] *2781* *2782* *2783*  
0*  
[[2791]] *2791* *2792* *2793*  
0*  
[[2801]] *2801* *2802* *2803*  
0*  
[[2811]] *2811* *2812* *2813*  
0*  
[[2821]] *2821* *2822* *2823*  
0*  
[[2831]] *2831* *2832* *2833*  
0*  
[[2841]] *2841* *2842* *2843*  
0*  
[[2]]  
[1] *Sample1* *Sample2* *Sample3* *Sample4* *Sample5* *Sample6*  
[7] *Sample7* *Sample8* *Sample9* *Sample10* *Sample11* *Sample12*  
[13] *Sample13* *Sample14* *Sample15* *Sample16* *Sample17* *Sample18*  
[19] *Sample19* *Sample20* *Sample21*  
  
> plot(m)  
> help.search("scatterplot")  
> pairs(m)  
> pairs(m[,1:4])  
> pairs(m[,1:4])  
> []
```

The scatterplot matrix shows the relationship between variables Sample1 through Sample6. The diagonal elements are labeled 'sample' and show a strong positive correlation. The off-diagonal elements show varying degrees of correlation between pairs of samples.

TableLens



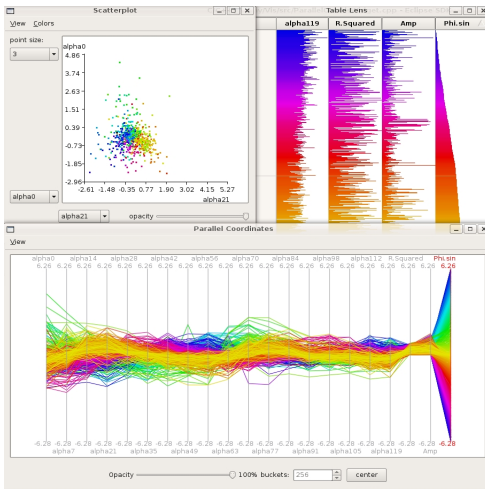
[Rao and Card, 1994]

Linking and Brushing

Colormap

active:

Opacity



Performance

Depends on

- Size of the data set
- Number of plugins loaded
- Operation in progress
- Available hardware (GPU?)

Results

- Lower response times than GGobi/iPlots/RGL/Mondrian
- Good performance for middle-sized datasets

Discussion

Objectives achieved

- Extendable Visual-Analytics-Framework
- Independent Visualization Library
- Hardware-accelerated Graphics
- Statistical Backend using R
- Interactivity
- Good performance / Low response times

Problems

- Redundancy in frequently used calculations
- Very basic interface to R
- categorical data only supported via the R-plugin

Future Work

Future Work

- Incorporate meta-information into datamodel to avoid redundancy (e.g. maxima)
- Add/Improve plugins (Heatmap, 3D Plots, ...)
- Extend interface to R (hot-linking, selections)
- Improve GPU-usage (textures, framebufferobjects ...)

Thank You!

References I



Adler, D. and Nenadic, O. (2003).

A Framework for an R to OpenGL Interface for Interactive 3D graphics.

In Proc. of the 3rd International Workshop on Distributed Statistical Computing.



Rao, R. and Card, S. K. (1994).

The table lens: merging graphical and symbolic representations in an interactive focus + context visualization for tabular information.

In Proc. of SIGCHI conference on Human factors in computing systems, pages 318–322, New York, NY, USA. ACM.



Swayne, D. F., Lang, D. T., Buja, A., and Cook, D. (2003).

GGobi: evolving from XGobi into an extensible framework for interactive data visualization.

Computational Statistics and Data Analysis, 43(4):423–444.



Urbanek, S. and Theus, M. (2003).

iPlots - High Interaction Graphics for R.

In Proc. of the 3rd International Workshop on Distributed Statistical Computing.