

Enhancing visual analysis of environmental monitoring data through an information visualization-based prototype

Celmar Guimarães da Silva¹, Gisela de Aragão Umbuzeiro¹

¹Faculdade de Tecnologia – Universidade Estadual de Campinas (Unicamp)
Limeira – São Paulo – Brazil

{celmar,giselau}@ft.unicamp.br

***Abstract.** This paper presents an approach for using interactive visual data representation in the field of environmental monitoring. Monitoring programs of dredged material disposition areas generates lots of spreadsheets, each one containing thousands of data necessary for decision making. This paper proposes using Information Visualization techniques for enhancing visual data analysis, aiming to ease the detection of important patterns and trends in these datasets. An implemented prototype organizes this data in a way that it could be easily available for users' interactive and visual-based queries.*

1. Introduction

The new millennium started with a global concern about the environment and its changes. This concern is related to the increasing human population and the concentration of its activities, and both cause serious problems such as atmosphere, soil and water pollution [Artiola et al., 2004]. Public and private institutions do environmental monitoring activities in order to assess the negative impact of these factors. These activities consist in observing and studying the environment, aiming to get relevant information for decision making related to prevention or remediation of these impacts.

This assessment includes objective observation of distinct characteristics of the environment under study which is composed by physical, chemical and biological aspects. These aspects are interrelated into interconnected and inseparable processes which involve soil, surface waters, atmosphere and live organisms [Artiola et al., 2004]. These characteristics may be analyzed from microscopic spatial scales (such as analyzing the existence of a specific kind of atom at a sample of the environment) to macroscopic ones (which may consider, for example, the Earth as a whole). The observations also have a temporal scale which may vary from a instantaneous time (a moment) to a geological time (which extrapolates ten million years); yearly and seasonal measurements are common in this domain [Artiola et al., 2004]. If analyzed individually, data originated from these environmental measurements provides limited information about the environment under study. However, a holistic analysis may allow for obtaining significant information for more effective decision making.

Artiola et al. (2004) points that these activities require a multidisciplinary approach, which aggregates professionals from Chemistry, Physics, Biology, Mathematics, Statistics and Computer Science. The role of Computer Science on this approach involves data storage and retrieval for decision making. About this role, Brown and Musil (2004) points three computing-related tools for data manipulation: spreadsheet managers (such as OpenOffice and Microsoft® Excel), database

management systems and programming languages. They compare these resources according to easiness of learning, scalability (in relation to the supported amount of data), data searching capability and data manipulation flexibility. Brown and Musil also declare a need of using graphical representation of monitoring data stored by computers. They argue that graphical representation may help users to identify potential problems at the data, such as unexpected outliers and trends, which may be not detected by statistical analysis.

In this sense, this work presents an initial proposal for using Information Visualization (InfoVis) techniques in the context of environmental monitoring. Based on some of these techniques, a software prototype was developed for helping in the data analysis process generated in a monitoring program of an ocean area that received dredged material [Umbuzeiro et al., 2009]. This prototype provides interactive views of these data, which helped analysts to understand data obtained from samples collected from 9 sampling sites in 29 sampling campaigns where several chemical and physical analysis were performed including several metals and organic compounds.

This paper was organized in five sections. In section 2, the environmental problem was presented. Based on theoretical background presented in section 3, a proposed solution is presented in section 4. Section 5 presents the developed prototype used to help data analysis. At section 6 we made final remarks and indicate future possible works.

2. Problem Characterization

Data analyzed in this study was provided by a consultancy agency. The data were obtained from a monitoring program of an area in the ocean that received dredging material for several months. The data was generated in order to obtain data about the quality of the sediment of the disposal area. The monitored area comprised 9 sampling sites where 29 sampling campaigns were performed. One sample of each site was collected in each campaign. Each sample was chemically analyzed for As, Cd, Pb, Cu, Cr, Hg, Zn and Mn. Physical parameters analyzed were pH, temperature and granulometry of the sediment. All these analysis generated a set of measurements of each sample, which were stored in 29 spreadsheets, one for each campaign, as exemplified in Figure 1. Each campaign' spreadsheet presents the values of each analyzed parameter for each sampling site. People in charge of monitoring compare obtained values with environmental quality criteria in order to assess environment quality and to indicate future actions. There are defined sediment criteria in regulations from different countries. Remediation actions may be necessary if the obtained values are not in compliance with the established criteria.

Data spatial patterns in those spreadsheets help analysts to compare the status of all sampling sites according to the analyzed parameters in the same campaign, i. e., at the same time interval. These patterns allow for either the comparison among measures of one parameter at distinct sampling sites, or the comparison among measures of different parameters related to a same sampling site. However, analysts may need to focus at data related to a specific sampling site instead of a specific campaign. In this case, they must analyze all spreadsheets together in order to understand time-related variation of measures, such as how measures of all parameters vary at a single sampling site along time, or how measures of a single parameter change at all sampling sites along the campaigns. Also, different people involved in the analysis process may require distinct data views for their data interpreting and decision making. In all these

situations, it is possible to use a spreadsheet manager for creating subsets of data of interest, but this process is time consuming and error-prone. This suggests that a static table does not fit all needs of viewing distinct data subsets, and that some alternative approach could be used.

In this sense, visual data analysis may be enhanced by techniques and concepts from Information Visualization, a research area which aims to ease the process of obtaining data and understanding information on visual analysis of datasets. Distinct Information Visualization techniques use computing resources for graphically and interactively representing data, trying to optimize the use of human visual capabilities in order to understand phenomena which might not readily lend themselves to visual-spatial representations [Card et al., 1999; Chen, 2002]. Next section presents some techniques which may be important at enhancing environmental monitoring analysis capabilities, and which are the heart of this paper's solution.

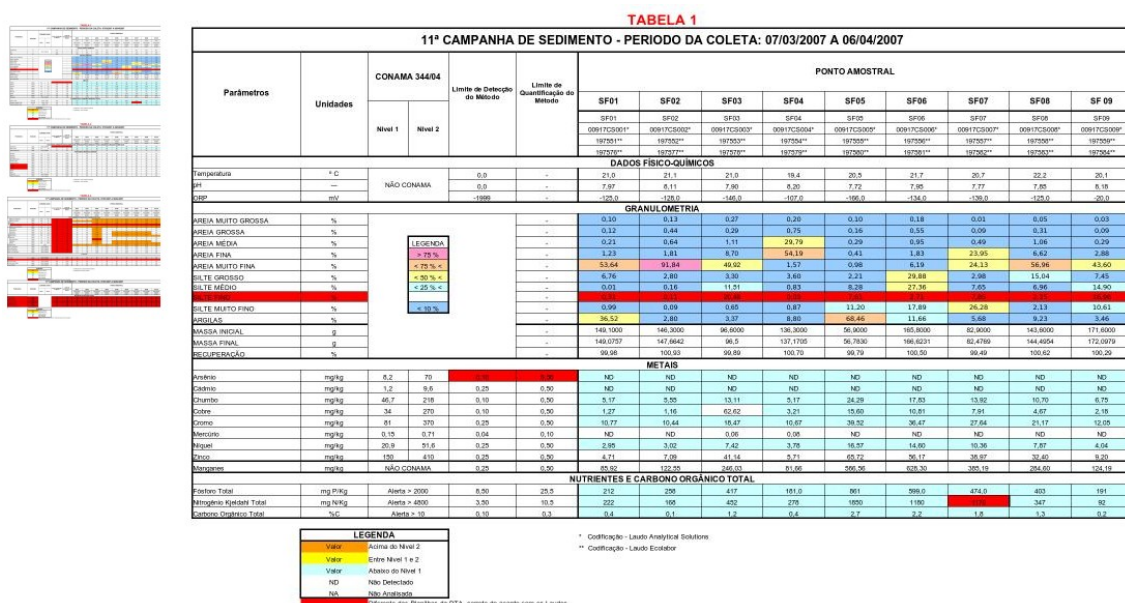


Figure 1. At left, a single campaign spreadsheet. At right, the first part of this spreadsheet, presenting sampling sites (SF01 to SF09), some physical and chemical analysis parameters analyzed and their respective values (Umbuzeiro et al., 2009).

3. Information Visualization concepts and techniques

Information Visualization presents a lot of techniques and concepts related to data organization and characterization, definition of visual data representations and of mapping between data and visual-spatial properties of graphical elements, such as color, shape and position at the screen [Card et al., 1999]. It also points the important role of interaction at the process of understanding data, given that “(...) a graphic is never an end in itself: it is a moment in the process of decision making” [Bertin apud Spence, 2001].

In the context of this work, one of the first considerations to do is related to data dimensionality. Data related to environmental sampling have intrinsically five dimensions (or variables): sampling site, campaign (which may be interpreted as time),

parameter (either physical or chemical), measured value and respective action level. However, only two or three variables may be associated to spatial orthogonal axis at the screen; in fact, it may exist a third axis if a 3D representation is used, but occlusion-related problems may occur, as pointed by Chuah et al. (1995). Hence, in order to represent all dimensions, it is necessary to map them to graphical properties such as color, shape, textures [Card et al., 1999], multidimensional icons [Spence, 2001] or even digits [Tufte, 1990]. Other approaches are related to the reuse of spatial coordinates, such as recursive use of space, based on zooming functions, as used in Pad++ [Bederson e Hollan, 1994]; sharing axis between distinct representations, as presented by Mackinlay (1986); overloading of orthogonal spaces, as used in n-Vision [Feiner e Beshers, 1990]; matrix-organized stamp-like graphics, known as “small multiples” [Tufte, 1990]; and parallel coordinates [Inselberg, 1997; Siirtola, 2000; Novotný & Hauser, 2006].

If not all dimensions are mapped to graphical or spatial properties of visual representation's elements, using interaction is a possible approach for including remaining dimensions on the analysis. Interactive visual representations allow for analyzing data from different viewpoints, which may enhance users' mental model about the analyzed data and, consequently, may help them in decision making activities. As presented by the Information Visualization reference model of Card et al. (1999) and as reinforced by Silva (2007), user interaction may happen when he chooses (or refines) a dataset to be represented, a mapping between dataset dimensions and visual representation's properties, or even some behaviors of this representation, such as showing details when asked by user, or changing user's viewpoint of the representation. Hence, defining how users may interact with graphical representations is also an important issue when defining these representations.

There are many interaction techniques available for Information Visualization-related representations. When defining matrix-based representations, users may be interested in permuting some of its columns and/or rows in order to evidence similar data groups, isolate outliers, and form visual patterns that help them to analyze the presented dataset. Enabling users to make this reorganization – which is called “reorderable matrix” technique by Bertin [apud Siirtola & Mäkinen, 2005] – focus their attention at the permuting actions and results, while hard work of this process is done by the computer [Siirtola & Mäkinen, 2005].

Dynamic queries [Shneiderman, 1994] is also an important interaction technique, which allow for users to control visual query parameters. These parameters are often presented by buttons, check boxes, list boxes, sliders, range-sliders or other interactive controls. This technique requires that users queries be quickly answered by the system. It is derived from a human-computer interaction technique called direct manipulation [Ahlberg & Shneiderman, 1994], which is based on four principles: (1) representing visually the world of action; (2) allowing for rapid, incremental and reversible actions; (3) enabling selection by pointing (and not only by typing); and (4) immediately and continuously displaying results of users interaction. Some of the most relevant advantages of dynamic queries are: users' ease of learning for querying, without need to understand complex syntaxes of typed commands; minimization of user errors, given that visual controls of parameters exposes its own syntax; and a short time between query formulation and its result exhibition, which may ease users construction of a internal model or cognitive map of the data under analysis [Spence, 2001; Shneiderman, 1994].

In the perspective of a dynamic query-based visual representation, the time available for system feedback may not be dismissed. Spence (2001) points that the time between users interaction and system feedback should be less than 0.1 second, which he called responsive interaction. This property leads user to understand that each screen change presented by the software is a direct consequence of one of his actions, because action and change are stimuli that fuse into a single percept when they happen within this period [Card et al., 1999].

In order to provide some guidance for development of Information Visualization systems, Shneiderman (1996) presented a guideline, called Visual Information Seeking Mantra, which has been useful at distinct application scenarios, as pointed by Craft & Cairns (2005). This guideline is defined as: “Overview first, zoom and filter, then details-on-demand”. In other words: first, present users a overview about data to be visually analyzed; second, give them the possibility of using zooming and filtering resources, in order to focus on relevant datasets and to temporally discard remaining data; last but not least, allow for getting details about data according to users' need.

4. Proposed solution

Given the context of environmental monitoring and the specific problem at hand, this paper proposes the use of Information Visualization techniques in order to enhance visual analysis of data. This work proposes using reorderable matrix-like representations of these data, enhanced by dynamic queries for dealing with data-inherent dimensionality problem. Filtering and detailing capabilities should also be available, in order to provide understanding about general and detailed aspects of data, according to the Visual Information Seeking Mantra.

Besides these characteristics, some data-representation requirements were also defined for the proposed solution. First, the proposed solution must present sample data related to a specific sampling site, which would enable users to compare how measured values for all parameters vary through distinct campaigns (it is, along time). It must represent sample data related to a specific parameter, allowing for visual analysis of how measured values of this parameter vary through distinct campaigns and sampling sites. It must also exhibit sample data from a specific campaign, enabling analysis of measured values at this campaign and its variation for each parameter along all sampling sites.

Using dynamic queries technique, each of these data representations should allow for users to alter, respectively, the sampling site, parameter or campaign they refer to. Providing access to reports from which a selected dataset was obtained is also an important feature related to the “details on demand” part of Shneiderman's mantra.

Other defined requirements are:

- Datasets with distinct units should not be displayed at the same screen, in order to simplify understanding and implementation. For exemplifying distinct units, sediment granulometry is expressed in percentage, while concentration of a chemical element is measured in mg/kg.
- It should be provided a summary statistical measure for analyzing sediment quality related simultaneously to all metal parameters. Each metal has its own action levels, and this measure should be normalized according to action level 1. Hence, consider a campaign c_x , a sampling site p_y , the set of metal parameters $\{m_1, m_2, \dots, m_M\}$, their respective set of action levels of type 1

$\{a_{11}, a_{12}, \dots, a_{1M}\}$, and the measured value $v(c_x, p_y, m_z)$ for a specific campaign, sampling site and metal. The defined summarization value $q(c_x, p_y)$ related to a single campaign and a single sampling site may be calculated as

$$q(c_x, p_y) = \frac{1}{M} \sum_{i=1}^M \frac{m_i}{aI_i},$$
 which will be called *quality quotient* for this

campaign and sampling site. This quotient is, then, a mean of measured values for each metal, which were normalized by action levels of type 1.

- Given that action levels may be distinct for different countries, options for indicating the country whose action levels should be considered at the representation should be available.

Given these requirements, the following section presents a prototype which implements a subset of them.

5. The VisSed Prototype

Based on the elements presented in section 4, a prototype software was implemented in order to facilitate the analysis of the environmental dataset at hand. The software, which was named *VisSed – Sediment-related Data Visualizer (Visualizador de Dados de Sedimentos, in Portuguese)*, is based on the requirements and InfoVis-related techniques presented at the previous section. VisSed uses JInfoVis [Silva, 2006], a prototype of Information Visualization toolkit which was previously used for graphical and interactive representation of learning management systems data. JInfoVis enables VisSed to obtain datasets related to the current environmental problem, which must be available through a previously defined database. It also allows for the representation of queried data through visual representations such as the interactive table-like one used by the prototype.

Given that VisSed is a prototype, only a subset of the complete dataset was inserted into the VisSed database. This subset contains all sampling sites and campaigns, but only the following parameters: percentage of sand and of silt and clay; quantities of As, Cd, Pb, Cu, Cr, Hg, Ni, Zn, Mn and total PAH - Polycyclic Aromatic Hydrocarbons. The quality quotient presented in the previous section was also included at VisSed; however, it is not stored at the database, but it is calculated when needed.

Figure 2 presents a conceptual model of the problem's dataset. According to this model, a measured value is related to a sampling site, a campaign and a parameter. Also, a given parameter has associated action levels, which are defined by each country. VisSed must represent these concepts, and, hence it must treat them as problem's dimensions. The following presentation of VisSed snapshots will exhibit how these dimensions were incorporated as graphical, spatial and interactive elements.

VisSed's initial screen is presented at Figure 3. It is organized in four parts: (A) a chart or visual representation, (B) a value selection list, (C) a view selection list, and (D) a country selection list. The prototype maintains this organization through its entire use.

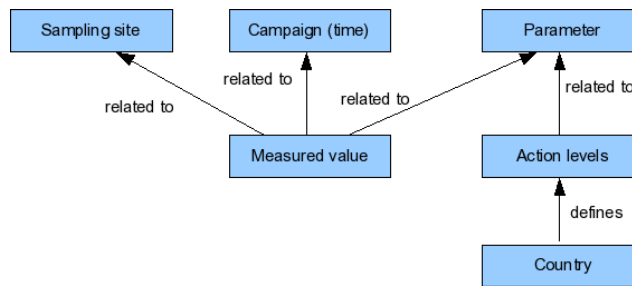


Figure 2: Conceptual model of the problem of sediment analysis.

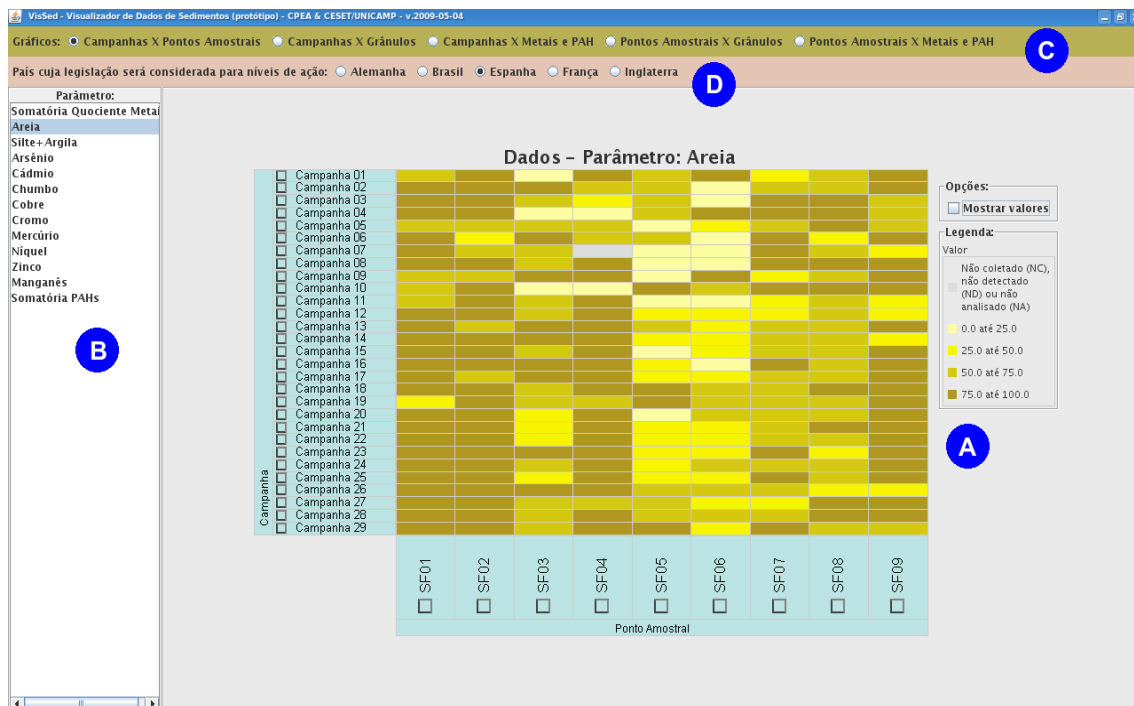


Figure 3. VisSed's initial screen.

The matrix-based chart (A) is provided by JInfoVis. The representation has two axes, X and Y, and it uses colors for classifying cells' values according to a caption. At Figure 3, each cell represents values of the “sand” parameter for each sampling site and campaign; the darker the cells, the higher the percentage of sand at the analyzed sediment. A check box named “Show values” (“Mostrar valores”), over the caption, allows for users to indicate if they want to see the exact value of each cell with digits, instead of only presenting color. Besides, columns and lines may be manually reordered by users according to their needs. The color set defined by chart caption is different according to the objective of the presented chart: yellow to brown sequences for showing sand concentrations, different intensities of red for representing metal and HPA concentrations, and blue and red for representing quality quotients.

Taking into consideration that a set {parameter, campaign, sampling site} has a single measured value to be presented, and that two of these three dimensions are mapped to the chart axes, values of the remaining dimension of this set are presented at a list box (B). Users may select on this list a single value which defines what subset of the measured data must be presented. At Figure 3, the parameter “sand” is selected, and it defines that all values from the chart are related to percentages of sand at the sample.

If this parameter is changed, the chart is immediately updated.

Different views of the dataset may be selected through a view selection list (C). Each view presents the following possible combinations of axes: Campaigns X Sampling sites, Campaigns X Grains, Campaigns X Metals and total PAH, Sampling sites X Grains, and Sampling sites X Metals and total PAH. Each view aims to ease user analysis of the variation of measured values according to axes-related dimensions. Figure 4 shows an example of the “Sampling sites X Metals and total PAH” view.

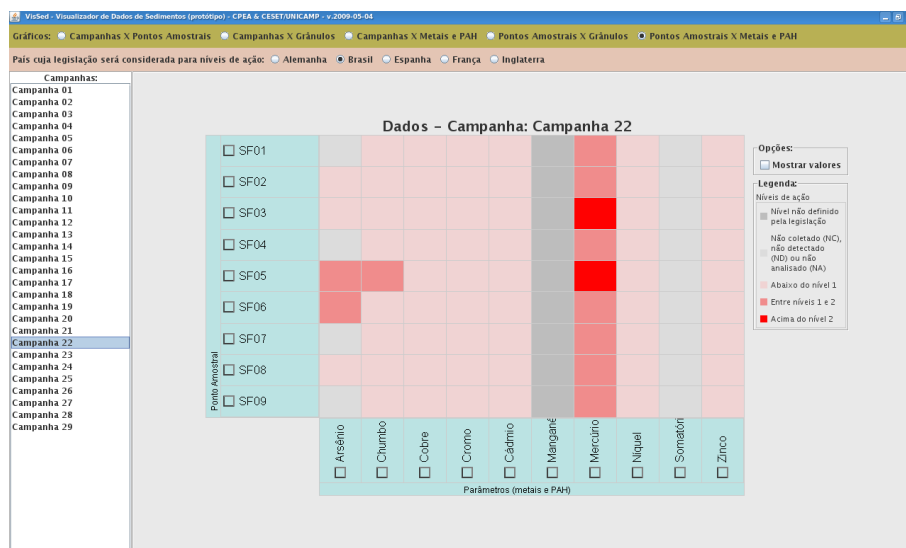


Figure 4. Measured values from a specific campaign, grouped by metals and HPA and by sampling sites.

A last chart control is a country selection list (D). It allows for users to select the country whose definitions of action levels will be considered for defining chart's caption colors (when presented data is related to PAH or metals). Captions related to metal or total PAH are defined by action levels, which are differently defined by each country. When a user changes the selected country, the chart instantaneously presents how dataset is interpreted according to this country's legislation. For example, Figure 5 presents how a same dataset (quality quotients for all campaigns and all sampling sites) is differently represented according to Brazilian and German action levels. This comparison should be viewed with caution because the basis for the calculation of the action levels are different.

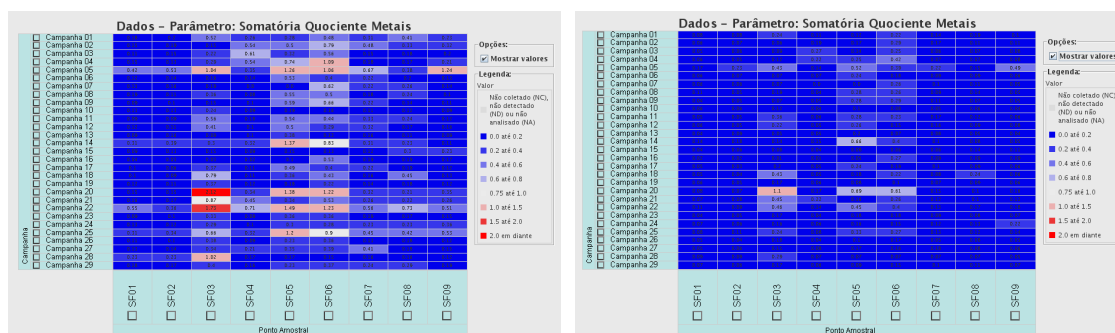


Figure 5. Quality quotients according to Brazilian (at left) German (at right) action levels.

The effectiveness of using VisSed was experimented when the authors used it for analyzing the complete dataset of dredged sediments monitoring presented at section 2. While VisSed was under development, one of the authors of this paper executed the analysis using regular spreadsheet manager to generate graphs to enhance data visualization. It was necessary a task force of a group of people to generate the various graphs required for a good overview of the dataset. After finishing the current version of VisSed, it was possible to use it for analyzing patterns and trends on the dataset in much less time.

6. Final Remarks and Future Works

Using software which allow for interactive visual data exploration may ease analysis and interpretation of datasets. Therefore, the main contribution of this paper is to present an Information Visualization-based approach for multivariate analysis of environmental monitoring data. Based on the use of visual and interactive representation of dredged sediments' measures, VisSed allows for data observation and interpretation under distinct viewpoints, which enable analysts to form a better mental model about these data and to make better decisions than they could do if they use only static spreadsheets. Besides these advantages, VisSed may be adapted for other similar environmental monitoring tasks, such as air pollutant or water programs, given that they have similar conceptual models and requires fast responses that could be provided by the software on a proof-of-concept analysis.

VisSed is currently at a prototyping stage and demands enhancements and user evaluations. Future versions of VisSed may include enhancements such as: distinct types of visual representations of data, with distinct levels of details about it; statistics-related functionality; capabilities for saving a dataset to a new file (such as CSV or XLS format files), in order to be used at other software; enabling users to access the original reports from which presented data was originated; and the definition of file patterns for information interchange between data suppliers (which are responsible for managing data collection) and data consumers (which are responsible for interpreting the collected data, such as the presented prototype).

Acknowledgement

This work was funded by CPEA – Consultoria Paulista de Estudos Ambientais. The authors want to thank Raquel Carnivale Silva for helping in the data organization and interpretation, and Sylvia Lima for valuable suggestions.

References

- Ahlberg, C. and Shneiderman, B. (1994). “Visual Information Seeking: Tight Coupling of Dynamic Query Filters with Starfield Displays”. *Proceedings of ACM Conference on Human Factors in Computing Systems*, New York, pp. 313-317.
- Artiola, J.F., Pepper, I.L. and Brusseau, M.L. (2004). “Monitoring and Characterization of the Environment”. In: Artiola, J.F., Pepper, I.L. and Brusseau, M.L. (2004). *Environmental Monitoring and Characterization*. Elsevier, USA, pp. 1-9.
- Brown, P. and Musil, S.A. (2004). “Automated Data Acquisition and Processing”. In: Artiola, J.F., Pepper, I.L. and Brusseau, M.L. (2004). *Environmental Monitoring and Characterization*. Elsevier, USA, pp. 49-67.

- Card, S.K., Mackinlay, J.D. and Shneiderman, B. (1999). *Readings in Information Visualization – Using Vision to Think*. Morgan-Kaufmann Publishers, San Francisco.
- Chen, C. (2002). “Editorial – Information Visualization”. *Information Visualization 1*, pp. 1-4, Palgrave Macmillan.
- Chuah, M.C., Roth, S.F., Mattis, J. and Kolojejchick, J.A. (1995). “SDM: Selective Dynamic Manipulation of Visualizations”. *Proceedings of ACM Symposium on User Interface Software and Technology*, pp. 61-70.
- Craft, B. and Cairns, P. (2005). “Beyond Guidelines: What Can We Learn from the Visual Information Seeking Mantra?”. *Proceedings of the Ninth International Conference on Information Visualization*, pp. 110-118.
- Feiner, S.K. and Beshers, C. (1990). “Worlds within Worlds: Metaphors for Exploring n-Dimensional Virtual Worlds”. *Proceedings of ACM Symposium on User Interface Software and Technology*, pp. 76-83.
- Inselberg, A. (1997). “Multidimensional Detective”. *Proceedings of IEEE Information Visualization'97*, pp. 100-107.
- Mackinlay, J.D. (1986). “Automating the Design of Graphical Presentations of Relational Information”. *ACM Transactions on Graphics*, 5(2), pp. 110-141.
- Novotný, M. and Hauser, H. (2006). “Outlier-preserving Focus+Context Visualization in Parallel Coordinates”. *IEEE Transactions of Visualization and Computer Graphics* 12(5), pp. 893-900.
- Shneiderman, B. (1994). “Dynamic Queries for Visual Information Seeking”. *IEEE Software*, 11(6), pp. 70-77.
- Shneiderman, B. (1996). “The Eyes Have It: A Task by Data Type Taxonomy for Information Visualizations”. *Proceedings of IEEE Symposium on Visual Languages*, pp. 336-343.
- Siirtola, H. and Mäkinen, E. (2005). “Constructing and reconstructing the reorderable matrix”. *Information Visualization 4*, pp. 32-48.
- Siirtola, H. (2000). “Direct Manipulation of Parallel Coordinates”. *Proceedings of CHI 2000*, pp. 119-120.
- Silva, C.G. (2006). *Learning Management Systems' database exploration by means of Information Visualization-based query tools*. Doctoral thesis. Institute of Computing, University of Campinas, Brazil. (In Portuguese).
- Silva, C.G. (2007). “Considerations about using Information Visualization for helping information management”. *Proceedings of XXXIV Integrated Seminary about Software and Hardware* (In Portuguese).
- Spence, R. (2001). *Information Visualization*. Addison-Wesley.
- Tufte, E. R. (1990). *Envisioning Information*. Graphics Press.
- Umbuzeiro, G. A., Silva, C. G. and Silva, R. C. (2009). *Assessment of main contamination data related to the monitoring of CODESP's oceanic disposition area for dredged material and adjacent regions – Santos, SP, Brazil*. Technical report (In Portuguese).