**Expanding Frontiers -**

**Challenges and Opportunities in Maritime Education and Training**

# Application of an Internet-based Tool for Visualization of Multi-Dimensional Objects with Missing Data in Maritime Education

*Natalia Nikolova*
*Assoc. Prof. PhD., N. Vaptsarov Naval Academy, V. Drumev 73, Varna 9026, Bulgaria,*
natalianik@gmail.com

*Daniela Toneva Zheynova*
*Assoc. Prof. PhD., Technical University – Varna, 1 Studentska Str., Varna, Bulgaria,*
d_toneva@abv.bg

*Kalin Kalinov*
*Assoc. Prof. DSc, Captain, N. Vaptsarov Naval Academy, V. Drumev 73, Varna 9026,*
*Bulgaria,* va_vms@yahoo.com

*Boyan Mednikarov*
*Prof. DSc, Captain, N. Vaptsarov Naval Academy, V. Drumev 73, Varna 9026, Bulgaria,*
bobmednikarov@abv.bg

*Kiril Tenekedjievld*
*Prof. DSc, N. Vaptsarov Naval Academy, V. Drumev 73, Varna 9026, Bulgaria, Kiril.*
*Tenekedjiev@fulbrightmail.org*

**Abstract:** In certain disciplines of maritime education it is necessary to operate with multi-dimensional data that is difficult for students to comprehend. Such disciplines happen to be environmental monitoring of water, technical diagnostics of ship machines and equipment using statistical pattern recognition, decision analysis of alternatives with multi-dimensional consequences, etc. Worst still, some of the multi-dimensional vectors are not known completely. Even though there are some mathematical algorithms developed to tackle the problem stated, in some countries maritime education suffers certain limitations caused by the high prices of specialized software. Here, an Internet based tool is presented which solves two interconnected problems. The first problem is the visualization of multi-dimension environmental vectors. That problem is addressed with two methods: The Principal Component solution of orthogonal Factor Analysis model and with Multi-Dimensional Scaling procedure. The second problem is the generation (the imputation) of the missing data in multi-dimension environmental vectors. The implemented solution uses modified Roweis algorithm for Expectation Maximization (EM) algorithm for Principal Component solution of orthogonal Factor

Analysis method. The functionality of the system is described along with environmental examples.

**Keywords**: multi-dimensional scaling, principle component solution, Roweis algorithm,

# 1. General Set Up of the Visualization Problem

The records in an environmental database normally contain several environmental variables measured at approximately the same time point. Let us assume that we have an environmental database with $n$ records each of which is represented by a $p$-dimensional column vector

$$\vec{x}_j = \left( x_1^{(j)}, x_2^{(j)}, \ldots, x_p^{(j)} \right)^T,$$

where $T$ stands for the transpose operator and $i=1,2,\ldots,n$. The whole data base can be written in a $p \times n$ matrix $D = \left( \vec{x}_1, \vec{x}_2, \ldots, \vec{x}_n \right)$.

If all the measurements of the environmental variables are known and recorded into the data base $D$ then we can denote this matrix with $D_{full}$. The problem which we face is to compress the $p$-dimensional vectors in $D_{full}$ into $m$-dimensional vectors $\vec{f}_j$ for $j=1, 2, \ldots, n$. If $m=2$, the compressed vectors can be depicted by a scatter plot of the first against the second compressed coordinates. If $m=3$, the compressed vectors can be depicted by three scatter plots which show respectively the first against the second, the second against the third and the first against the third compressed coordinates. We call that a visualization problem. The result of the visualization problem is a compressed database in form of a $m \times n$ matrix $D_{compr} = \left( \vec{f}_1, \vec{f}_2, \ldots, \vec{f}_n \right)$.

Two classical algorithms are employed in the Internet based tool for solution of the visualization problem (that is for the data reduction from the $p$-dimensional to the $m$-dimensional space):

- Multi-Dimensional Scaling Procedure (MDS) [4]
- Principal Component Solutions of Exploratory Orthogonal Factor Analysis Model (PCA_EOFA) [6]

So the visualization module of the Internet-based tool has two sub-modules; the MDS sub-module and the PCA_EOFA sub-module.

The MDS sub-module uses 10 different distances as a measure of dissimilarity between the vectors:

- Manhattan distance: $\delta_{k,j} = \sum\limits_{i=1}^{p} \left| x_i^{(k)} - x_i^{(j)} \right|$

- Euclidean Distance: $\delta_{k,j} = \sqrt{\sum\limits_{i=1}^{p} \left( x_i^{(k)} - x_i^{(j)} \right)^2}$

- Minkovski $r$-distance (with $r=3, 4, 5, 10, 20$) : $\delta_{k,j} = \sum\limits_{i=1}^{p} \left( \left| x_i^{(k)} - x_i^{(j)} \right|^r \right)^{1/r}$

- Hamilton (or Chebishev) distance: $\delta_{k,j} = \max\limits_{i} \left( \left| x_i^{(k)} - x_i^{(j)} \right| \right)$

- Normalized Euclidean Distance: $\delta_{k,j} = \sqrt{\sum_{j=1}^{p} \left( \dfrac{x_i^{(k)} - x_i^{(j)}}{s_i} \right)^2}$ , where $s_i$ is the sample standard deviation of the coordinate $i$.

- Mahalonobis distance: $\delta_{k,j} = \left( \vec{x}_k - \vec{x}_j \right)^T K^{-1} \left( \vec{x}_k - \vec{x}_j \right)$, where $K$ is the sample covariance matrix of the data.

The Java realization of the MDS sub-module has been based on [1] where the 'stress' criterion (formula (8) in the accompanying paper [8]) is optimized. The mathematical algorithms of PCA_EOFA sub-module are explained in the accompanying paper [8].

Sometimes the record $\vec{x}_j$ of the data matrix $D$ is not a full vector. Instead, part of the coordinates $x_1^{(j)}, x_2^{(j)}, \ldots, x_p^{(j)}$ are missing due to various reasons. In this case the same data matrix is called $D_{miss} = \left( \vec{x}_1, \vec{x}_2, \ldots, \vec{x}_n \right)$. For example, let $p=5$ and for the third observation $\vec{x}_3$, the second and the fifth coordinates are missing (denoted as NaN), whereas the first the third and the forth coordinates are -17, 24 and 6 respectively. Then:

$$\vec{x}_3 = \begin{pmatrix} x_1^{(3)} \\ x_2^{(3)} \\ x_3^{(3)} \\ x_4^{(3)} \\ x_5^{(3)} \end{pmatrix} = \begin{pmatrix} -17 \\ NaN \\ 24 \\ 6 \\ NaN \end{pmatrix}$$

There are many reasons for the missing data, especially when social surveys are involved, because people do not always want to answer all questions. As long as the environmental data is gathered by direct measurements with some measurement devices we have accepted the hypothesis that the data is missing at random [7]. The problem that we face is to restore the values of the missing part in $D_{miss}$. Of course the missing part of $D_{miss}$ is an unobservable quantity which has to be estimated from the observable quantities (the known data) in the data base. We call that the Missing Data Imputation Problem. The result of the missing data imputation problem is an imputed database in the form of a $p$ x $n$ matrix $D_{imp} = \left( \vec{x}_1^{(imp)}, \vec{x}_2^{(imp)}, \ldots, \vec{x}_n^{(imp)} \right)$. $D_{imp}$ has the same values as $D_{miss}$ but the missing values are substituted with values generated by a mathematical algorithm called imputed values.

A powerful expectation maximization algorithm for estimation of unobservable quantities was proposed in the seminal paper [3]. Six years later that algorithm was properly proven [10]. In [9], the algorithm was used for imputation of missing data in the context of a factor analysis method. In the accompanying paper [8] a modification of the Roweis algorithm is proposed which quickly converges for medium sized data points. In the developed missing data imputation module of the Internet based tool the linear system $T_j \vec{z}_j^{(nimp)} = \vec{t}_j$ with $p$ equations and $N_j$ unknowns in point 7) of the algorithm in [8] is solved by QR decomposition of the matrix $T_j$ [5].

Once the missing data imputation problem is solved and $D_{imp}$ is generated then (having in mind that there are no more missing position in $D_{imp}$) it can be put that $D_{full} = D_{imp}$ and the imputed data can be passed to the visualization problem for data reduction.

## 2. Examples and demonstration

An Internet based system is developed to realize procedures for visualization of multi-dimensional data and for imputing missing data into multi-dimensional data sets. The system is available through the webpage www.ubss-tuv.com and it has been developed under the activities of the UPGRADE Black Sea Scientific Network (financed by the Seventh Framework Program of the EU). The main window of the platform is shown in fig. 1. Here we focus on the options "Visualization" and "Missing Values", the others being described in [3]. The procedures to be commented on may run on three file formats: text document (*.txt), Excel document (*.xls) and comma delimited text (*.csv).

### 2.1 Visualizing multi-dimensional data

A test file is created, containing 70 numbers of 7-dimensional test measurements. The structure of the file is shown in fig. 2 (for the sake of all examples, only the operation with Excel files shall be discussed here, all the others being very similar). The names of the coordinates (i.e. the environmental variables) are given on the second row, columns from 2 to $p+1=8$ of the file. The third row (same columns) contains the variable dimensions. The fourth row, columns from 2 to $p+1=8$ contains the

values of $\vec{x}_1$, and the other vectors are written consecutively on the next rows. The last one is on row $n+3=70+3=73$. The first column of the file contains the consecutive number of the vectors.

The visualization section of the UBSS system is given in fig. 3. Let us use factor analysis to perform the visualization. After selecting the file (fig. 4) the user might choose from the panel shown in fig. 5:

1.  whether to use standardized data or not (the first option is selected in this

    example), i.e. whether to replace the original data $\vec{x}_j$ in the input by the data normalized in the first place in formulae (30) and (31) in [8];

2.  the output dimension to be 2 or to be 3 ( $m=3$ is selected in this example);

3.  the type of compressed data that is ordinary least square factor score using formula (28) or weighted least square factor score (29) of [8] (the first option is selected in this example).

Choosing to proceed, the system generates output as follows:

1.  scatter plot of the first two coordinates of the compressed data (see fig. 6 for one of the three scatter plots, because the compressed data is three-dimensional);

2.  the $m$-dimensional compressed data values on screen;

3.  an excel file with five sheets, containing the compressed data points, the covariance matrix, the absolute error, the relative error and the cumulative percentage of the total variance explained (fig. 7) (which are formulae respectively (21), (25), (26) and (27) from [8]).
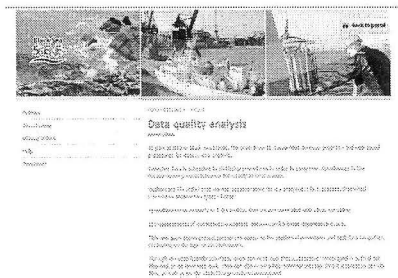
If we choose to visualize using multi-dimensional scaling, then the partial differences would appear in the input and in the output of the system. The parameters which determine the course of calculation of the multi-dimensional scaling are defined in a panel shown in fig. 8.

1.  1) whether to use standardized data or not (the first option is selected in this example), as in the visualization with factor analysis.

2.  2) the output dimension to be 2 or to be 3 ( $m=2$ is selected in this example);

3.  3) the type of distance measure (Manhattan distance in this example).

Choosing to proceed, the system generates output as follows:

1.  1) scatter plot of the first two coordinates of the compressed data (as in fig. 6 but here this is the only scatter plot because the compressed data is two-dimensional);

2.  2) the compressed $m$-dimensional data values on screen;

3.  3) an excel file with four sheets (fig.9), containing the compressed $m$-dimensional data points, ( $\vec{f}_j$ ), the distances in the original space ( $\delta_{k,j}$ ), the distances in the compressed space ( $\delta_{k,j}$ ) and the difference between the corresponding pairs of distances ( $\delta_{k,j} - d_{k,j}$ ).

Such a procedure might be employed to identify outliers or certain internal connections (e.g. sometimes data would lie on the same straight line, which would indicate high level of multicolinearity in the data).

## 2.2  Imputing Missing Values

A test file is created containing 500 numbers of 10-dimensional test measurements, where some of the values are missing. The structure of the file is shown in fig. 10. The names of the coordinates (i.e. the environmental variables) are given on the second row, columns from 2 to $p+1=11$ of the file. The third row (same columns) contains the variable dimensions. The fourth row, columns from 2 to $p+1=11$ contains the values of $\vec{x}_1$, and the other vectors are written consecutively on the next rows. The last one is on row $n+3=500+3=503$. The first column of the file contains the consecutive number of the vectors.

The missing values section of the system is given in fig. 11. After choosing the *.xls format and the file itself (fig. 12) the user might choose:

1.  the dimension of the compressed space ($m=5$ in the example)

2.  whether to use standardized data or not (the first option in this example), i.e. whether to divide by $\sqrt{k_{i,i}}$ or not in step 3) of the algorithm in [8];

3.  the drifting moments (either moving or fixed moments on $X$ on each iteration, i.e. whether to execute step 9) in the algorithm of [8] or not; it is executed for this example);

4.  whether to have additional fixing of the recovered data or not (whether to execute step 17) in the algorithm of [8] or not; it is executed for this example).

Choosing to proceed, the system generates output as follows:

1.  graphics of the imputed data (see fig. 13 for the first coordinate, where the imputed points are in blue; such a plot is created for each coordinate);

2.  all the imputed values on the screen;

3.  an excel file with a single sheet containing the imputed set $D_{imp}$ with structure the same as the input data, however the imputed values are placed instead of the missing ones in red (fig. 14).

## 2.3 Real Data Analysis

A file with real environmental data is created, related to chemical characteristics of sea water, containing 110 numbers of 13-dimensional measurements. Some of the values in the dataset are missing. The structure of the file is shown in fig. 15. The names of the coordinates (i.e. the chemical parameters) are given on the second row, columns from 2 to $p+1=14$ of the file. The third row (same columns) contains the variable dimensions. The fourth row, columns from 2 to $p+1=14$ contains the values

of $\vec{x}_1$, and the other vectors are written consecutively on the next rows. The last one is on row $n+3=110+3=113$. The first column of the file contains the consecutive number of the vectors.

We set the following input parameters for the missing value procedure:

1.  6-dimensional compressed space;

2.  standardized data;

3.  moving drifting moments on $X$ on each iteration;

4.  additional fixing of the recovered data.

The system generates the following result:

1.  graphics of the imputed data (see fig. 16 for the second coordinate; such a plot is created for each coordinate);

2.  the imputed data set on screen;

3.  an excel file with a single sheet structured the same as the input data, but the imputed values are placed instead of the missing ones in red (fig. 17).

The resulting file is sent for visualization with multi-dimensional scaling at $m=3$ using Mahalonobis distance measure. As an output, the system generated the following:

1. 1) graphics of the compressed data (three two-dimensional plots; fig. 18 gives the plot for the second and third compressed coordinates);

2. 2) the $m$-dimensional compressed data values on the screen;

3. 3) an excel file with four sheets as in section II.1 (see fig. 19).

The resulting file from the missing values procedure is also sent for visualization with factor analysis at $m=2$. As an output, the system generated the following:

1. graphics of the compressed data (see fig. 20);

2. the $m$-dimensional compressed data values on screen;

3. an excel file with five sheets, containing the compressed data points, the covariance matrix, the absolute error, the relative error and the cumulative percentage of the total variance explained (fig. 21) as in section II.1.



Fig. 1



Fig. 2



Fig. 3



Fig. 4

Fig. 5



Fig. 6



Fig. 7



Fig. 8



Fig. 9



Fig. 10

**Fig. 11**



**Fig. 12**



**Fig. 13**



**Fig. 14**



**Fig. 15**



**Fig. 16**



**Fig. 17**

Fig. 18



Fig. 19



Fig. 20



Fig. 21

## 3. Conclusions

This paper presented the operation of an Internet based tool which solves two interconnected problems – visualization of multi-dimensional environmental vectors and imputation of missing values in such data sets. The procedures use modern algorithms and techniques to generate results. The tool is available free of charge to all users.

## 4. Acknowledgments

## References

[1]   Algorithmic Group, "*MDSJ: Java Library for Multidimensional Scaling (Version 2)*", University of Konstanz (2009)

[2]   Bechev, Ch., Nikolova, N.D., Teohareva, M., and Tenekedjiev, K., "Testing the Quality of Environmental Data by an Internet-Based Tool", *Journal of Environmental Protection and Ecology*, Vol. 10, No. 3, (2009), pp. 877-888

[3]   Dempster, A.P., Laird, N.M., and Rubin, D.B., "Maximum Likelihood from Incomplete Data via the EM Algorithm", *Journal of the Royal Statistical Society, Series B*, Vol. 39, (1977), pp. 1-39

[4]   Duda, R., Hart, P, and Stork, D. "*Pattern Classification and Scene Analysis*", A Wiley Interscience Publication, (1973)

[5]   Golub, G.H., and Van Loan, Ch. "*Matrix Computations*", Third Edition, Johns Hopkins (1996)

[6]   Johnson R., and Wicherin D., "*Applied Multivariate Statistical Analysis*", 6th ed., Prentice Hall, (2007)

[7]   Little, R. J.A., "A Test of Missing Completely at Random for Multivariate Data with Missing Values", *Journal of the American Statistical Association*, Vol. 83, No. 404, (1988), pp. 1198-1202

[8]   Nikolova, N.D., Toneva-Zheynova, D., Naydenov, D., and Tenekedjiev, K. "Imputing missing values of environmental multi-dimensional vectors using a modified Roweis algorithm", *International Workshop on Dynamics and Control in Agriculture and Food Processing*, (2012), Plovdiv, Bulgaria, (in print)

[9]   Roweis, S., "EM Algorithms for PCA and SPCA", *Neural Informatics Processing Systems, NIPS'1997*, (1977), pp. 626-632

[10]   Wu, J., "On the Convergence Properties of the EM Algorithm", *The Annals of Statistics*, Vol. 11, No. 1, (1983), pp. 95-103