

2013 International Conference on Computational Science

Interactive data mining by using multidimensional scaling

Piotr Pawliczek^{a,b}, Witold Dzwiniel^{*b}

^aUniversity of Texas, Department of Biochemistry and Molecular Biology, Houston, TX 77030, USA

^bAGH University of Science and Technology, Department of Computer Science, Al.Mickiewicza 30, 30-059 Kraków, Poland

Abstract

Blind choice and parameterization of data mining tools often yield vague or completely misleading results. Interactive visualization enables not only extensive exploration of data but also better matching of clustering/classification schemes to the type of data being analyzed. The multidimensional scaling (MDS), which employs particle dynamics to the error function minimization, is a good candidate to be a computational engine for interactive data mining. However, the main disadvantage of MDS is both its memory and time complexity. We developed novel SUBSET algorithm of a lower complexity, which is competitive to the best, currently used, MDS algorithms in terms of efficiency and accuracy. SUBSET employs reduced dissimilarity matrix, which structure allows for efficient usage of both multi-core CPU and SIMD GPU processor architectures. Consequently, SUBSET enables visualization of datasets consisting of an order of 10^5 data items on a standard personal computer or laptop. We compare a few strategies of dissimilarity matrix reduction and we present typical timings obtained by respective MDS algorithms on selected multithread CPU and GPU architectures.

© 2013 The Authors. Published by Elsevier B.V. Open access under [CC BY-NC-ND license](https://creativecommons.org/licenses/by-nc-nd/4.0/).

Selection and peer review under responsibility of the organizers of the 2013 International Conference on Computational Science

Keywords: multidimensional scaling; method of particles; incomplete distance matrix; multicore CPU; GPU.

1. Introduction

Data mining of large datasets consisting of data items $O_i \in \Omega$, ($i=1, \dots, M$), where Ω is an abstract data space, involves application of many machine learning tools such as classifiers, regression and clustering schemes. Many of them are specialized for analysis of a special case of Ω represented by Y space of multidimensional

* Corresponding author. tel.: +48 662130188, fax: +48 12 617 51 72

E-mail address: dzwinel@agh.edu.pl

feature vectors y_i , $\mathbf{Y} = \{y_i = (y_{i1}, \dots, y_{iN})\}_{i=1, \dots, M}$, where $N = \dim \mathbf{Y}$ (see Fig. 1a). However, in general, the data items O_i can have more sophisticated structures, which cannot be directly represented by the feature vectors. This problem can be partially overcome provided that it is possible to define a dissimilarity measure $\delta(O_i, O_j) \rightarrow \mathbb{R}^1$ between data items O_i and O_j . We assume that $\delta(\cdot, \cdot)$ is symmetric and $\delta(O_i, O_i) = 0$. Particularly, $\delta(\cdot)$ can be a distance obeying also the condition of triangle inequality. In general, the Ω space topology is represented by dissimilarity matrix $\Delta = \{\delta_{ij}\}_{M \times M}$. The vector representation of Ω can be derived by employing multidimensional scaling (MDS) procedure [1-3].

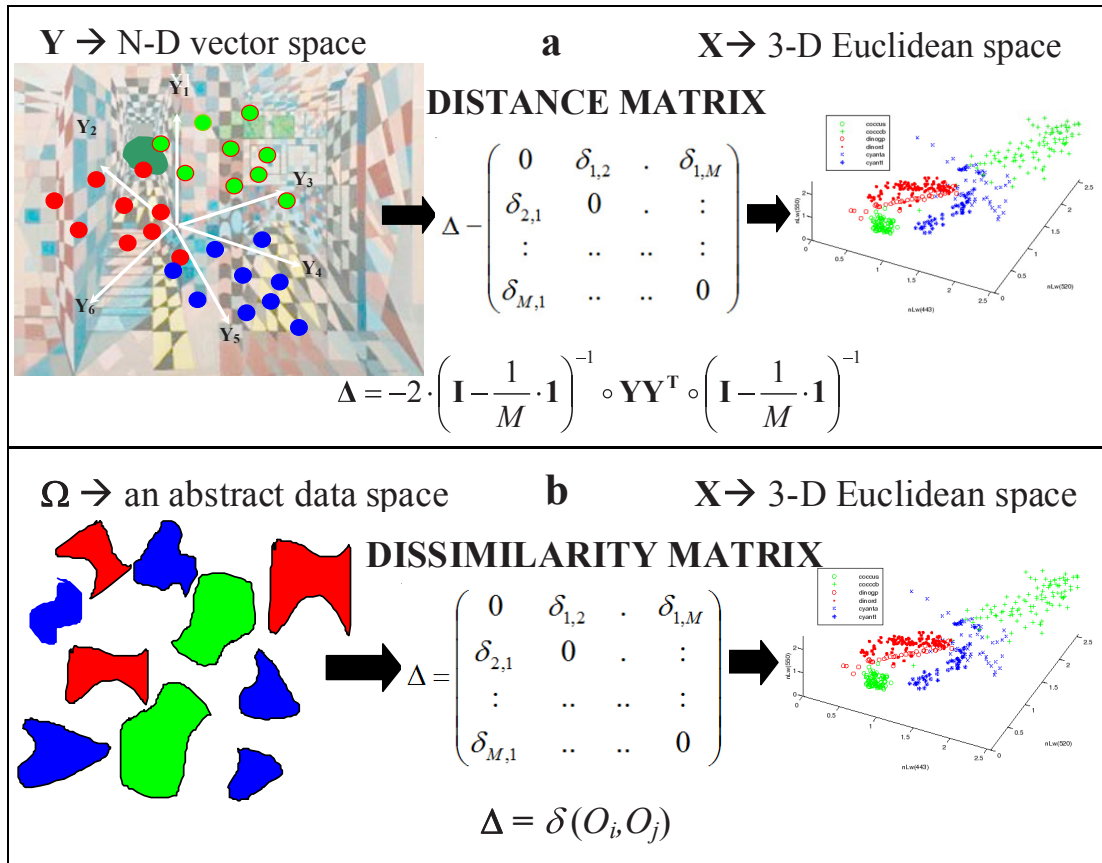


Fig.1. Multidimensional scaling applied for visualization in 3-D Euclidean space \mathbf{X} data from a) a multidimensional feature space \mathbf{Y} , b) an abstract Ω space for which only dissimilarity matrix Δ is known (e.g., dissimilarity measure between shapes).

Multidimensional scaling (MDS) (see e.g. [1,2,3]) is a bijection $B: \Omega \rightarrow \mathbf{X}$ of a “source” space of abstract items $\Omega = \{O_i; i=1, \dots, M\}$ onto a “target” vector space $\mathbb{R}^n \ni \mathbf{X} = \{x_i = (x_{i1}, \dots, x_{in})\}_{i=1, \dots, M}$, where $n = \dim \mathbf{X}$, which reproduce topological structure of Ω in \mathbf{X} in respect to a given error criterion (see Fig. 1b).

Let us define the Euclidean matrix $\mathbf{d} = \{d_{ij}\}_{M \times M}$ in the target vector space \mathbf{X} , where d_{ij} is the Euclidean distance between vectors x_i and x_j which correspond to O_i and O_j , respectively. We assume that to preserve topological structure of Ω in \mathbf{X} an overall error $F(\|\Delta - \mathbf{d}\|)$ should be minimized, where $F(\cdot)$ - called the “stress” function [1-3] - is an increasing function $F: \mathbb{R}^1 \rightarrow \mathbb{R}^1$. The value of $F(\cdot)$ represents a discrepancy measure between dissimilarities Δ from Ω and corresponding distances \mathbf{d} from \mathbf{X} . The resulting matrix $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_M)$, which minimizes the “stress” function $F(\cdot)$, is the final outcome of multidimensional scaling. This way

Download English Version:

<https://daneshyari.com/en/article/490477>

Download Persian Version:

<https://daneshyari.com/article/490477>

[Daneshyari.com](https://daneshyari.com)