

MANAGEMENT OF BIG DATA IN SCIENTIFIC COMPUTING

Marcel ILIE¹, Augustin SEMENESCU²

Rezumat. Dezvoltările recente în sistemul de procesare al calculatoarelor au condus la generarea unor cantități mare de date. Procesarea acestor large baze de date reprezintă adevărate provocări pentru comunitatea științifică. O problemă particulară este interpretarea datelor, în particular separarea datelor utile de cele inutile. Mai mult decât atât clasificarea datelor în subgrupuri de date aduce provocări suplimentare. În general este acceptat ca analiza acestor date este foarte dificilă. De aceea modele eficiente și acurate de data-mining ar înlesni analiza acestor mari baze de date. În acest studiu propunem un model computațional bazat pe algoritmul k-means, identificat ca algoritmul de fuzzy-clustering. Folosirea algoritmului de fuzzy-clustering reduce timpul computațional cu 72% comparat cu metodele computaționale normale.

Abstract. The recent developments in computer processing has led to the generation of significant amount of data. However, the post-processing of this large amount of data poses significant challenges for the scientific community. A particular issues is the data interpretation, particularly regarding the segregation of valuable data from arbitrary data. Moreover, data classification in subgroups poses further challenges. It has been largely accepted that the analysis of these data sets may be cumbersome. Therefore, efficient and accurate data mining models would enable the analysis if large data-sets. In this research we propose a computational model based on the k-means algorithm, identified as fuzzy clustering algorithm. The use of the fuzzy clustering algorithm reduces the computational time by 72%, when compared with regular computational approaches.

Keywords: data management, numerical modeling, adaptive mesh refinement, k-means algorithm

1. Introduction

The rapid progress in the developments of computing technologies and software such as high-performance computing (HPC) have generated large amounts of data that require further post-processing, interpretation and dissemination [1-5]. Post-processing and visualization of this large amount of data pose significant challenges, particularly when the data needs to be assembled/coupled from various instantaneous time-frames [5, 6]. Some of these cases are encountered in the medical field, biology research, engineering applications, etc. Out of all these

¹PhD, Assistant Professor, Dept. of Mechanical Engineering, Georgia Southern University, Statesboro, GA 30458, USA, e-mail: milie@georgiasouthern.edu

²PhD, Professor, Dept. of Material Sciences, University Politehnica Bucharest, Bucharest, Romania, augustin.semenescu@upb.ro
