Bruno Baruque and Emilio Corchado

Fusion Methods for Unsupervised Learning Ensembles

# Studies in Computational Intelligence, Volume 322

**Editor-in-Chief**
Prof. Janusz Kacprzyk
Systems Research Institute
Polish Academy of Sciences
ul. Newelska 6
01-447 Warsaw
Poland
*E-mail:* kacprzyk@ibspan.waw.pl

---

Bruno Baruque and Emilio Corchado

# Fusion Methods for Unsupervised Learning Ensembles

Dr. Bruno Baruque
Departamento de Ingeniería Civil
Escuela Politécnica Superior
Universidad de Burgos
Avda. Cantabria, s/n
09006 Burgos, Spain
E-mail: bbaruque@ubu.es

Dr. Emilio Corchado
Departamento de Informática y Automática
Facultad de Ciencias
Universidad de Salamanca
Plaza de la Merced, s/n
37008 Salamanca
Spain
E-mail: escorchado@usal.es

# Abstract

The application of a "committee of experts" or ensemble learning to artificial neural networks that apply unsupervised learning techniques is widely considered to enhance the effectiveness of such networks greatly. This book examines in one of its chapters the potential of the ensemble meta-algorithm by describing and testing a technique based on the combination of ensembles and statistical PCA that is able to determine the presence of outliers in high-dimensional data sets and to minimize outlier effects in the final results. After that, it presents its central contribution, which consists on an algorithm for the ensemble fusion of topology-preserving maps, referred to as Weighted Voting Superposition (WeVoS), which has been devised to improve data exploration by 2-D visualization over multi-dimensional data sets. This generic algorithm is applied in combination with several other models taken from the family of topology preserving maps, such as the SOM, ViSOM, SIM and Max-SIM. A range of quality measures for topology preserving maps that are proposed in the literature are used to validate and compare WeVoS with other algorithms. The experimental results demonstrate that, in the majority of cases, the WeVoS algorithm outperforms earlier map-fusion methods and the simpler versions of the algorithm with which it is compared. All the algorithms are tested in different artificial data sets and in several of the most common machine-learning data sets in order to corroborate their theoretical properties. Moreover, a real-life case-study taken from the food industry demonstrates the practical benefits of their applications to more complex problems.

# Contents

# List of Figures

# List of Tables

# List of Algorithms