

OSCAR LUACES, FRANCISCO HERRERA, JOSÉ A. GÁMEZ, LUIS MARTÍNEZ, EDURNE BARRENECHEA,
JOSÉ RIQUELME, ALICIA TRONCOSO, BRUNO BARUQUE, MIKEL GALAR, HÉCTOR QUINTIÁN,
EMILIO CORCHADO (EDS.)

Actas de la XVII Conferencia de la Asociación Española para la Inteligencia Artificial


Ediciones Universidad
Salamanca

AQUILAFUENTE

224

©

Ediciones Universidad de Salamanca y
de cada autor

Motivo de cubierta:
Diseñadora María Alonso Miguel

1.º edición: septiembre, 2016
ISBN: 978-84-9012-632-5 (PDF)

Ediciones Universidad de Salamanca
www.eusal.es
eusal@usal.es

Realizado en España – Made in Spain

*Todos los derechos reservados.
Ni la totalidad ni parte de este libro
pueden reproducirse ni transmitirse sin permiso escrito de
Ediciones Universidad de Salamanca*

Obra sometida a proceso de
evaluación mediante sistema de revisión por pares a ciegas
a tenor de las normas del congreso

Ediciones Universidad de Salamanca es miembro de la UNE
Unión de Editoriales Universitarias Españolas
www.une.es

CEP

Parte IV.- VIII Simposio Teoría y Aplicaciones de Minería de Datos (TAMIDA 2016)

Aplicaciones:

Integración de Datos Medioambientales procedentes de Open Data mediante Preprocesado con R	
PAVEL H. LLAMOCCA, VICTORIA LÓPEZ	717
Búsqueda Paralela Exhaustiva Aplicada a Cadenas de ARNi	
JESÚS GARCÍA-RAMÍREZ, IVO H. PINEDA T. ,MARÍA J. SOMODEVILLA, MARIO ROSSAINZ, CONCEPCIÓN PÉREZ DE CELIS	727
Realimentación automática en evaluación por pares de respuestas abiertas mediante factorización de matrices	
JORGE DÍEZ, OSCAR LUACES, ANTONIO BAHAMONDE	735
Estudio y caracterización de variedades de ciruelas utilizando análisis de imagen y técnicas de Deep Learning	
FRANCISCO J. RODRÍGUEZ, ANTONIO GARCÍA, PEDRO J. PARDO, FRANCISCO CHÁVEZ, RAFAEL M. LUQUE-BAENA	745
Aplicación de técnicas de text mining para analizar las interacciones de los estudiantes en el proceso de aprendizaje de gestión de proyectos	
RUBÉN OLARTE-VALENTÍN, ANA GONZÁLEZ-MARCOS, FERNANDO ALBA-ELÍA, JOAQUÍN OR-DIERES-MERÉ	755
A Forecasting Methodology for Workload Forecasting in Cloud Systems	
FRANCISCO J. BALDÁN, SERGIO RAMÍREZ-GALLEGO, CHRISTOPH BERGMEIR, FRANCISCO HERRERA, JOSÉ M. BENÍTEZ	765

Bigdata:

Minería de patrones en Big Data	
FRANCISCO PADILLO, JOSÉ MARÍA LUNA, SEBASTIÁN VENTURA	769
Búsquedas exhaustivas de subgrupos con MapReduce en Big Data	
FRANCISCO PADILLO, JOSÉ MARÍA LUNA, SEBASTIÁN VENTURA	779
Enfoque MapReduce para el democratizado de métodos de selección de instancias	
ALEJANDRO GONZÁLEZ-ROGEL, ÁLVAR ARNAIZ-GONZÁLEZ, CARLOS LÓPEZ-NOZAL, JOSÉ F. DIEZ-PASTOR	789

Clasificación:

Anuran sound classification using MPEG-7 frame descriptors	
JAVIER ROMERO, AMALIA LUQUE, ALEJANDRO CARRASCO	801
Whatever you know, just tell me something: Crowd learning with free supervision	
JERÓNIMO HERNÁNDEZ-GONZÁLEZ, IÑAKI INZA, JOSÉ A. LOZANO	811
MLSMOTE: Approaching Imbalanced Multilabel Learning Through Synthetic Instance Generation	
FRANCISCO CHARTE OJEDA, ANTONIO J. RIVERA RIVAS, MARÍA J. DEL JESUS DÍAZ, FRANCISCO HERRERA TRIGUERO	821
Random Balance: Ensembles of Variable Priors Classifiers for Imbalanced Data	
JOSÉ F DIEZ-PASTOR, JUAN J. RODRIGUEZ, CÉSAR GARCIA-OSORIO, LUDMILA KUNCHEVA	823
Automatic device detection in web interaction	
PERONA, A. YERA, O. ARBELAITZ, J. MUGUERZA, N. RAGKOUSIS, M. ARRUE, J.E. PEREZ	825
Estrategia efectiva para el aprendizaje activo multi-etiqueta	
OSCAR REYES, SEBASTIÁN VENTURA	835

Preprocesamiento:

Impact of discretization with multivariate sequential patterns to do the classification of the survival prediction in Intensive Care Burn Unit	
ISIDORO J. CASANOVA, MANUEL CAMPOS, JOSÉ M. JUAREZ, ANTONIO FERNANDEZ-FERNANDEZ-ARROYO, JOSÉ A. LORENTE	847
Una herramienta para analizar conjuntos de datos multi-etiqueta	
JOSE M. MOYANO, EVA L. GIBAJA, SEBASTIÁN VENTURA	857

MLSMOTE: Approaching Imbalanced Multilabel Learning Through Synthetic Instance Generation

Francisco Charte Ojeda¹, Antonio J. Rivera Rivas², María J. del Jesus Díaz²,
Francisco Herrera Triguero¹

¹Dept. of Computer Science and Artificial Intelligence, E.T.S.I.I.T., University of Granada, Granada 18071, Spain.

²Dept. of Computer Science, E.P.S., University of Jaén, Jaén 23071, Spain.
francisco@fcharte.com, arivera@ujaen.es, mjjesus@ujaen.es,
herrera@decsai.ugr.es

Abstract. This is a summary of our article published in Knowledge-Based Systems [1] to be part of the MultiConference CAEPIA'16 Key-Works.

Keywords: Multilabel classification, Imbalanced learning, Resampling, Instance generation, SMOTE

1 Summary

The learning of classification models from imbalanced data is usually a challenging task, since most algorithms produce classifiers that tend to be biased towards the majority class. The higher is the imbalance level, the greater the likelihood of inducing this bias. Commonly, the classifier would achieve a good performance by simply predicting the majority class for all data patterns. As a consequence the minority class, which frequently is the most interesting to the researchers, suffers from miss-classification errors.

Resampling methods are among the most popular approaches to face imbalanced learning, since they provide a classifier independent mechanism to do so. The goal of these methods is to balance the distribution of classes, either by generating new samples associated to the minority class or by removing those linked to the majority class. Random oversampling and random undersampling are basic ways of performing this work. One the of the most successful ways to achieve this balanced distribution consists in generating synthetic instances, as proposed in the SMOTE [2] (*Synthetic Minority Over-sampling Technique*) algorithm.

Imbalance is a quite common problem in multilabel learning [3]. This is a non-standard classification task in which each data instance is associated to several class labels at once. Therefore, a multilabel classifier has to be able of predicting several outputs for each processed pattern. In this kind of datasets there are usually several minority labels and several majority ones. In addition

the imbalance levels, the ratio between the most common labels and the rarest ones, are usually huge. As a result, balancing the labels distribution in this kind of datasets is a more complex job than in standard classification. In late years several resampling algorithms for multilabel data have been proposed, including random undersampling and oversampling, as well as heuristic undersampling. Some ensemble-based solutions have been presented as well.

The method we introduced in [1] is a multilabel version of the well-known SMOTE [2] algorithm. It aims to produce synthetic instances linked to minority labels, instead of mere clones as previous oversampling proposals did. The main characteristics of MLSMOTE are the following:

- It takes into account the presence of several minority labels, producing new data samples associated to each one of them. Both SMOTE and previous multilabel oversampling techniques consider one class only, so new samples are exclusively associated to the rarest label.
- New data samples produced by MLSMOTE are given a synthetic set of input attributes, instead of being clones of existing instances. This way new patterns are not located into the same position than those taken as reference. The new attributes values are obtained by means of interpolation techniques, as in SMOTE.
- A synthetic labelset is computed for each new data sample. For doing so, the labels in the reference instance and their neighbors are taken into account. Three different strategies to produce these synthetic labelsets were tested. All previous resampling multilabel methods cloned the labelset of the reference instance.

An extensive experimentation was conducted to assess the performance of MLSMOTE. In it, a dozen multilabel datasets were processed with five imbalance aware algorithms, including oversampling (LP-ROS, ML-ROS, SmoteUG), undersampling (BR-IRUS) and ensemble (EML) based approaches. Each configuration was then given to five different multilabel classifiers (BR, RAKEL, CLR, HOMER and IBLR-ML), whose results were evaluated using common multilabel performance metrics. MLSMOTE was the best performer in all cases, achieving statistically significant differences in some of them. These result proved that synthetic instance generation through MLSMOTE, which includes the creation of synthetic labelsets, can be a successful approach when it comes to tackle imbalanced multilabel learning.

References

1. F. Charte, A. J. Rivera, M. J. del Jesus, and F. Herrera. MLSMOTE: Approaching imbalanced multilabel learning through synthetic instance generation. *Knowledge-Based Systems*, 89:385–397, 2015.
2. N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer. SMOTE: Synthetic minority over-sampling technique. *J. Artificial Intelligence Res.*, 16:321–357, 2002.
3. F. Herrera, F. Charte, A. J. Rivera, and M. J. del Jesus. *Multilabel Classification. Problem analysis, metrics and techniques*. Springer, 2016.