

TEXT CLASSIFICATION IN ROMANIAN FOR AN INTELLIGENT NEWS PORTAL

Traian REBEDEA¹, Ștefan TRĂUȘAN-MATU²

Rezumat. *Articolul prezintă un modul de clasificare a textului pentru un portal de știri în limba română. Tehnicile de prelucrare statistică a limbajului natural sunt combinate cu scopul de a realiza funcționarea complet autonomă a portalului. Fiecare știre este colectată automat dintr-un mare număr de surse de știri folosind clasificarea web. De aceea, tehnicile de îmbogățire a cunoștințelor programului sunt folosite pentru clasificarea automată a fluxului știrilor. Mai întâi, termenii sunt clasificați folosind un algoritm aglomerativ răsturnat; clasele rezultate corespund subiectelor principale. De aceea, sunt preluate mai multe informații despre fiecare dintre subiectele principale din mai multe surse de știri. Apoi, algoritmi de clasificare a textului sunt aplicați pentru etichetarea automată a fiecărei clase de știri (al căror număr este predeterminat). Au fost folosite peste o mie de știri, atât pentru antrenarea cât și pentru evaluarea clasificatorilor. În articol este prezentată o comparație completă a rezultatelor obținute prin fiecare metodă. Mai mult, sunt prezentate problemele specifice care apar datorită particularităților limbii române și sunt discutate soluțiile găsite.*

Abstract. *The paper presents a text classification module for a news portal for the Romanian language. Statistical natural language processing techniques are combined in order to achieve a completely autonomous functionality of the portal. The news items are automatically collected from a large number of news sources using web syndication. Afterward, machine-learning techniques are used for achieving an automatic classification of the news stream. Firstly, the items are clustered using a bottom-up agglomerative algorithm and the resulting groups correspond to the main news topics. Thus, more information about each of the main topics is acquired from various news sources. Secondly, text classification algorithms are applied to automatically label each group of news items in a predetermined number of classes. More than a thousand news items were employed for both the training and the evaluation of the classifiers. The paper presents a complete comparison of the results obtained for each method. Moreover, specific problems that arose due to the particularities of the Romanian language are presented and the solutions found are discussed.*

Keywords: natural language processing, text clustering, classification, news portal, intelligent agent

¹Drd. Eng., University "Politehnica" Bucharest, Department of Computer Science and Engineering.

²Prof. Dr. Eng., University "Politehnica" Bucharest, Department of Computer Science and Engineering; corresponding member of the Academy of Romanian Scientists (trausan@gmail.com).