

ACTIVE SEEDS SELECTION WITH A K-NEAREST NEIGHBORS GRAPH

Vũ Việt Vũ, Nicolas Labroche, Violaine Antoine, Lê Bá Dũng

TÓM TẮT:

Active learning allows semi-supervised clustering algorithms to solicit domain experts to retrieve a few set of class labels (or seeds) to improve their efficiency or the relevance of their results. However, some recent studies show that, even in the case of a good answer from the domain expert, semi-supervised clustering can see their performances drop with badly chosen seeds. Until now, only few works address the problem of determining the best queries for a clustering algorithm in an active learning context, and most of these studies are limited because of their hypothesis on the size and the shape of expected clusters. In this paper, we propose a new active seed selection algorithm that makes no hypothesis on the underlying data distribution. Experiments conducted on real data sets show the efficiency of this new approach compared to existing ones.

Active learning allows semi-supervised clustering algorithms to solicit domain experts to retrieve a few set of class labels (or seeds) to improve their efficiency or the relevance of their results. However, some recent studies show that, even in the case of a good answer from the domain expert, semi-supervised clustering can see their performances drop with badly chosen seeds. Until now, only few works address the problem of determining the best queries for a clustering algorithm in an active learning context, and most of these studies are limited because of their hypothesis on the size and the shape of expected clusters. In this paper, we propose a new active seed selection algorithm that makes no hypothesis on the underlying data distribution. Experiments conducted on real data sets show the efficiency of this new approach compared to existing ones.