

Classification of Imbalanced Data in a Streaming Environment: a Literature Review

R. Anderson, K. Chen and E. Jin

September 12, 2014

Abstract

We review the literature available around classification of imbalanced data in a streaming environment. First, we introduce the particular problems faced within streaming environments and with imbalanced classes. Through examining current research in the field, we explore the facets of classification in this environment: sampling techniques; approaches to detecting concept drift; base learners for building models to classify unseen data; and methods of utilising ensembles to improve classification accuracy. Finally, we examine evaluation of these methods, and the pertinent measures used.

- 1 Introduction**
- 2 Sampling methods**
- 3 Concept drift detection**
- 4 Base learners**

Supervised learning requires models to be trained from labelled data. When model-building, we need an approach that can use explanatory variables and classes in a training set of data to create a set of rules for classifying unseen instances. This approach is our base learner: our tool that will build our classifier to apply to unseen instances of our data [?].

In a streaming environment, however, we can only make one pass through our data. Concept drift requires a learner that can be run at regular intervals, so as to understand the current distribution that our instances fall into. This requires a trade-off between time, accuracy and memory. A successful learner will be sufficiently accurate while running quickly and with manageable memory overhead. These requirements are much more important than in a static environment, where we have a finite amount of data and time to process it.

Decision trees have particular features which make them a good fit for classifying streams with concept drift. The C4.5 decision tree, and its predecessor, the ID3 tree, are very popular within the research community [?]. Their appeal lies in their simplicity – at every node, a binary decision is made related to one of the explanatory variables, resulting in a divide-and-conquer approach which results in fast classification and low overhead for the model. This leads to some clear drawbacks – only rectangular regions can be used to discriminate between classes in the feature space for instance. Splits are made through information gain measures, seeking to maximise the change in impurity of the nodes with each split. These trees will naturally overfit, if fully grown, but can be pruned to achieve a desired level of fit to the data. Their adaptability and speed make them a natural fit in classification of streams.

Domingos and Hulten [?] highlight a major concern with C4.5, ID3 and CART – they assume that all training examples can be held in memory, and are limited in the examples they can consider. Other trees such as SPRINT and SLIQ have attempted to resolve this issue by using a window to scan large datasets sequentially, but this is very slow when building complex trees, and still requires all data to reside in disk-space. They propose the Very Fast Decision Tree (VFDT), which requires that data items are only read once, in a small and constant time. They utilise Hoeffding Bounds [?] to guarantee that the statistically best attribute to divide a node on can very commonly be found using a small number of examples of instances as could be with infinite. They propose the Hoeffding tree, which regularly checks whether splitting a node on the best attribute will lead to an improvement in a selected measure (information gain or Gini) beyond a set threshold. This allows it to continue to adapt, even with infinite data, without sacrificing significant quality nor having burdensome memory requirements. Parameters need to be set by the user: the minimum number of examples before reviewing whether to split a node; the threshold to allow a split; and whether the algorithm can rescan prior examples. Even if these parameters are well selected by the user, there is a danger that they will lead to poor performance in environments featuring significant concept drift.

Recognising the issues posed by concept drift, Hulten et al. [?] proposed the Concept-Adapting Very Fast Decision Tree (CVFDT), an adaptation on the VFDT to give it the adaptability to handle changing data. Through adopting a sliding window approach, the CVFDT increments counts of recently seen data while reducing counts of older data. This should have no effect when there is no concept drift, but where there is drift, prior splits will no longer pass the Hoeffding tests that deemed them worthwhile. At this point, an alternate subtree is grown by the algorithm from that node. If it accurately classifies new data better than the existing subtree, the existing subtree is replaced. Overall, this requires additional memory to keep alternate subtrees in memory, and summary statistics of split quality at each node. However, these memory requirements are still constant with the data items received. CVFDTs cannot utilise old sub-trees which are discarded, which could provide potential performance improvements where concepts drifts can revert back to prior probability distributions.

Decision trees have trouble accurately classifying rare classes. Consider a

binary class problem: if a very large majority belongs to one class, often a classifier achieves best performance by classifying all instances as that class. Lyon et al [?] propose an improvement on the VFDT by introducing a new criterion for splitting: Hellinger distance. This measure seeks the attributes which are most disjoint in the data, and seeks to create tree splits based on these features. Instead of prioritising information gain, this tree prioritises class discrimination. This will often lead to a smaller portion correctly classified, but should amplify the number of the rare class successfully classified. Calculating Hellinger distance can be costly, as it scales to the number of classes and needs to discretize continuous classes into bins, leading to further multiplicative scaling. Lyon et al propose the HD-VFDT which uses the splitting criterion above, and also the Gaussian Hellinger Very Fast Decision Tree, which assumes normality of the distribution to simplify calculation of the Hellinger Distance. However, instances in a data stream are not independent, as time affects their distribution, and so this could lead to flawed results in some circumstances. Experimental analysis of this approach showed statistically significant improvements using these methods over simple Hoeffding Trees where the minority class makes up 1% or less of the total data. They do not evaluate the approach thoroughly with regard to time nor memory use. Specifically, it would be interesting to measure the difference in accuracy, speed and time required between the two variants of the Hellinger tree proposed.

Liechtenwalter and Chawla [?] have found potential in using Hellinger trees in environments with concept drift. By weighting the Hellinger distance with Information Gain to measure distance between datasets, a learner can decide whether model learnt on previous data is valid on current data. Pazzolo et al [?] showed, using the HDDT proposed by Cieslak and Chawla [?] and C4.5 tree, that trees with this weighting generally outperformed trees without the weighting in class-imbalanced environments with drift. These papers do not address the multi-class problem, which would be a natural extension of this research.

5 Ensemble learners

6 Evaluating approaches

7 Conclusion