# ML-KFHE: Multi-label ensemble classification algorithm exploiting sensor fusion properties of the Kalman filter

Arjun Pakrashi<sup>1,2\*</sup> and Brian Mac Namee<sup>1,2</sup>

<sup>1</sup>School of Computer Science, University College Dublin, Dublin, Ireland.

<sup>2</sup>Insight Centre for Data Analytics, Dublin, Ireland.

\*Corresponding author(s). E-mail(s): arjun.pakrashi@ucd.ie; Contributing authors: brian.macnamee@ucd.ie;

#### Abstract

Despite the success of ensemble classification methods in multi-class classification problems, ensemble methods based on approaches other than bagging have not been widely explored for multi-label classification problems. The Kalman Filter-based Heuristic Ensemble (KFHE) is an ensemble method that exploits the sensor fusion properties of the Kalman filter to combine several classifier models, and that has been shown to be very effective. This work proposes a multi-label version of KFHE, ML-KFHE, demonstrating the effectiveness of the KFHE method on multi-label datasets. Two variants are introduced based on the underlying component classifier algorithm, ML-KFHE-HOMER, and ML-KFHE-CC which uses HOMER and Classifier Chain (CC) as the underlying multilabel algorithms respectively. ML-KFHE-HOMER and ML-KFHE-CC sequentially train multiple HOMER and CC multi-label classifiers and aggregate their outputs using the sensor fusion properties of the Kalman filter. Extensive experiments and detailed analysis were performed on thirteen multi-label datasets and eight other algorithms, which included state-of-the-art ensemble methods. The results show, for both versions, the ML-KFHE framework improves the predictive performance significantly with respect to bagged combinations of HOMER (named E-HOMER), also introduced in this paper, and bagged combination of CC, Ensemble Classifier Chains (ECC), thus demonstrating the effectiveness of ML-KFHE. Also, the ML-KFHE-HOMER variant was

found to perform consistently and significantly better than the compared multi-label methods including existing approaches based on ensembles.

 ${\bf Keywords:}$  multi-label classification, ensembles, Kalman filter, classifier fusion

### 1 Introduction

A multi-class classification task assigns an object to at most one class. Realworld classification problems exist, however, where an object can be assigned to more than one class simultaneously [1]. In other words, an object can be *labelled* with more than one class at the same time. For example, an image of a landscape may contain mountains, sea and sky and therefore it can be a member of each of the corresponding classes [2]. Similarly, music can be tagged with more than one genre. Such problems are known as *multi-label* classification problems.

Formally, multi-label classification problems can be defined as follows. Let  $\boldsymbol{x}_i$  be a datapoint from a d-dimensional input space  $\mathcal{X}$  of real and/or categorical attributes. Also, let the set of all possible labels for a specific multi-label classification problem be  $\mathcal{L} = \{\lambda_1, \lambda_2, \ldots, \lambda_q\}$ , from which a subset of labels,  $\mathcal{L}_i \subseteq \mathcal{L}$ , is applicable to the datapoint  $\boldsymbol{x}_i$ . Here labels in  $\mathcal{L}_i$  are called the *relevant* labels, and  $(\mathcal{L} - \mathcal{L}_i)$  are called the *irrelevant* labels for  $\boldsymbol{x}_i$ . Then a typical multi-label dataset is defined as  $\mathcal{D} = \{(\boldsymbol{x}_i, \boldsymbol{y}_i) || 1 \leq i \leq n\}$ , where n is the number of datapoints in the dataset;  $\boldsymbol{x}_i = \{x_{i1}, x_{i2}, \ldots, x_{id}\}$  is a d-dimensional vector indicating the  $i^{th}$  datapoint; and  $\boldsymbol{y}_i = \{y_{i1}, y_{i2}, \ldots, y_{iq}\}$  is a binary vector indicating the label assignments  $\mathcal{L}_i$  for the  $i^{th}$  datapoint. Here  $y_{ij} = 1$  if  $\lambda_j \in \mathcal{L}_i$ , that is, the  $j^{th}$  label is relevant to the  $i^{th}$  datapoint, and  $y_{ij} = 0$  if  $\lambda_j \notin \mathcal{L}_i$ . The objective of multi-label classification is to learn a model  $\boldsymbol{h}$ , which predicts the relevance of every label for a new datapoint  $\boldsymbol{d}$ , written as  $\boldsymbol{h}(\boldsymbol{d})$ .

Multi-label classification algorithms can be categorised as either *problem* transformation or algorithm adaptation methods [3]. Problem transformation methods—for example classifier chains [4]—break the multi-label problem down into smaller multi-class classification problems. Algorithm adaptation methods—for example BPMLL [5]—modify multi-class algorithms to directly train on multi-label datasets.

Ensemble classification methods train multiple component classifiers and aggregate them. Generally, ensemble methods perform better than the single component classifiers [6] and ensemble classifiers based on boosting generally perform better than bagging methods [7]. In the multi-label classification literature, several methods have been proposed that combine multiple multi-label models to form an ensemble. These, however, are mostly problem transformation methods based on *bagging* or majority voting based approaches to building ensembles [4, 8–13]. There are very few boosting approaches in

the multi-label classification literature. AdaBoost.MH [14] is probably the most prominent boosting approach in multi-label classification, but does not perform well when compared to other methods [15]. Although based on bagging and majority voting methods, algorithms like Classifier Chains (CC) [4] and RAKEL [12] which still stay a very competitive [16, 17] and has received test of time award in ECML 2017 and 2019, respectively. As boosting based methods in multi-class classification perform so much better than approaches based on bagging, but in the multi-label context they didn't perform as much as they perform in the multi-class domain, probably due to the high degree of imbalance which multi-label classification that are based on ideas similar to boosting.

Kalman Filter-based Heuristic Ensemble (KFHE) [18] is a recently proposed algorithm that frames ensemble training as a state estimation problem which is solved using a Kalman filter [19, 20]. Although Kalman filters are most commonly used to solve problems associated with time series data, this is not the case for KFHE. Rather, this work exploits the data fusion property of the Kalman filter to combine individual multi-class component classifier models to construct an ensemble. This can be interpreted as effectively being in between boosting and bagging, and therefore is expected to benefit from exploiting the advantages of both types of ensembles. Given the nature of this method, as it effectively falls between boosting and bagging, this can be especially helpful as multi-label datasets are inherently highly imbalanced and can be helpful to balance the drawbacks and benefits of both boosting and bagging.

By utilising the sensor fusion properties of the Kalman filter in the KFHE framework, this article proposes a multi-label classification algorithm ML-KFHE, and demonstrates its effectiveness using two variants of multi-label classification methods, ML-KFHE-HOMER and ML-KFHE-CC. Here, ML-KFHE-HOMER and ML-KFHE-CC which are ensembles of HOMER [21] and classifier chains (CC) [4] algorithms, respectively, using the ML-KFHE framework.

To demonstrate the effectiveness of the ML-KFHE ensemble combination method, ensemble versions of CC and HOMER were compared. Ensemble version of CC called *Ensemble Classifier Chains* (ECC) was proposed in [4]. Although bagged HOMER was not found in the literature, a simple bagged ensemble version of HOMER, named E-HOMER, was also proposed in this work to be able to directly compare the effectiveness of the ML-KFHE-HOMER. The source code of the algorithms are made online: https: //github.com/phoxis/kfhe-homer

The contributions of this paper are:

• An ensemble multi-label classification method, ML-KFHE, which exploits the sensor fusion properties of a Kalman filter in KFHE algorithm to

combine component classifiers. Two variants are presented, ML-KFHE-HOMER and ML-KFHE-CC, which incorporates HOMER and CC, respectively, as component classifiers in the KFHE framework.

- An extensive experiment demonstrating the effectiveness of using the ML-KFHE method to combine multi-label classifiers over using bagging based multi-label ensemble algorithms, as well as demonstrating the overall effectiveness of the proposed methods.
- Introduction of E-HOMER, a simple bagged ensemble version of HOMER to compare the ML-KFHE ensembling and bagged HOMER.

The remainder of the paper is structured as follows. Section 2 describes existing multi-label ensemble classification algorithms, as well as the relevant aspects of CC and HOMER. Section 3 introduces the proposed ML-KFHE-HOMER and ML-KFHE-CC methods, by first introducing KFHE and then ML-KFHE. The design of the evaluation experiments performed is described in Section 4, and the results of these experiments are discussed in detail in Section 5. Finally, Section 6 concludes the article and discusses directions for future work.

## 2 Related Work

This section discusses the current state-of-the-art ensemble methods used for multi-label classification and then mentions the relevant aspects of the HOMER and CC algorithms, which are later used in this work.

#### 2.1 Ensemble Methods for Multi-label Classification

In multi-label classification the same principles are used as in multi-class ensemble algorithms—although most of the multi-label ensembles are bagging based (in which multiple independent classification models vote on the relevance of each label). Some important multi-label ensemble classification algorithms will be described briefly next.

Ensemble Binary Relevance (EBR) and ECC [4] are simple bagged version of the Binary Relevance and Classifier Chains classifier algorithms. Label Powerset models can also be bagged to form *Ensemble Label Powerset* (ELP) models [15]. *Ensemble of Pruned Sets* (EPS) [22] is also a simple bagged version of Pruned Set and prevents overfitting properties of the Pruned Set algorithm. RAkEL-o [12] is also ensemble classifier algorithms, but instead of bagging it divides the label space into smaller overlapping subsets, then learn LP models for each of these subsets. RF-PCT [10], is a Random Forest [23] of Probabilistic Clustering Trees (PCT) [24] and therefore is also a bagging type ensemble model. *Triple Random Ensemble for Multi-label Classification* (*TREMLC*) [8] randomly performs several samples of datapoints and features at the same time and trains multiple LP models on each such sample combining the prediction from these trained models. *Clustering Based for Multi-label Classification* (*CBMLC*) [25] first performs a clustering on the feature space of the dataset to generate k clusters, each of which ideally holds similar datapoints. Next, it trains any multi-label algorithm for each of the k clusters. In summary, these bagged classifier algorithms range from combining existing component classifiers by training on bootstrap samples or sub-samples of the datapoints (EBR, ELP, EPS), then training different sub-samples of the datapoints along with different random subsets of the labels (RAkEL-o) or different random order of the labels (ECC). The other types performs bagging like subsampling but also considers subsampling or pre-processing the attributes as well (RF-PCT, TREMLC, CBMLC).

In boosting AdaBoost.MH and AdaBoost.MR [14] are the most frequently used multi-label ensemble classification models in existence. These two extensions of AdaBoost [26] for multi-label classification. AdaBoost.MH minimises the Hamming loss, and considers the labels independently and applies AdaBoost. Whereas, AdaBoost.MR considers pairwise label ranking by minimising a function similar to rankloss which penalises incorrect label ordering. Then it gives more weight to the datapoints which has a higher rank loss.

Some other less known algorithms are  $AdaBoost.MH^{kr}$  [27] which uses a relatively more complex model per boosting iteration by replacing the usual weak hypotheses with a sub-committee of weak hypotheses. *RFBoost* [28] improves AdaBoost.MH by training the weak learner on only a few top ranked features instead of using all the features. RFBoost was experimented using two different feature ranking methods and shown to perform similar or slightly better than AdaBoost.MH.

#### 2.2 HOMER

*Hierarchy of Multi-label Classifiers (HOMER)* [21] also divides the multi-label dataset into smaller subsets of labels, but in a hierarchical manner. This method divides the dataset based on the labels, but establishes a hierarchical relationship between the partitions. The root of the hierarchy has all labels  $\mathcal{L}$ and the entire dataset associated with it. Every leaf node has one associated label  $\lambda_k$ . Any internal node v have only a subset of labels  $\mathcal{L}_v \subset \mathcal{L}$  associated with itself, which is a union of the label subsets associated to its children. Therefore,  $\mathcal{L}_v = \cup \mathcal{L}_c, c \in children(v)$ , where children(v) indicates the children of the node v in the hierarchy. Each node only keeps the datapoints that have at least one of the associated labels of the node in their relevant set. Therefore, the dataset at the node v is  $\mathcal{D}_v = \{(x_j, y_j) || \forall (x_j, y_j) \in$  $parent(v), \exists_k \lambda_k \in \mathcal{L}_v \land y_{jk} = 1$ . Each internal node v is also associated with a meta label  $\mu_v$ , which is associated with all the datapoints in the corresponding node, and each internal node's datapoints are also associated with meta labels assigned by its children. The meta labels are generated as  $D'_v = \{(\boldsymbol{x}_j, \boldsymbol{y}'_j) \| \boldsymbol{y}'_j = \bigcup_{c \in children(v)} \mu_c \}.$ 

The root node and all the internal nodes v will have a trained model  $h_v$  associated to them, which is trained using the dataset  $D'_v$  with the target being the meta labels associated with each datapoint in the child nodes. Utility of

the meta labels is to indicate which branch or branches in the hierarchy have to be followed to predict the labels. A new datapoint t starts from the root of the hierarchy and travels one or more paths from root to leaf. All the labels represented by the leaves which are encountered by the datapoint t are taken as the prediction.

#### 2.3 Classifier Chains

Classifier Chains [4] take a similar approach to binary relevance but explicitly take the associations between labels into account. Again a one-vs-all classifier is built for each label, but these classifiers are chained together such that the outputs of classifiers early in the chain (the relevance of specific labels) are used as inputs into subsequent classifiers. Let  $\tau : \{1 \dots q\} \rightarrow \{1 \dots q\}$  be a permutation function which gives a new ordering or a chain of the labels  $\lambda_{\tau(1)} \succ \lambda_{\tau(2)} \succ \dots \succ \lambda_{\tau(q)}$ . For a label  $\lambda_{\tau(l)}$  a dataset is formed  $\mathcal{D}_{\tau(l)} =$  $\{([\mathbf{x}_i, pre_{\tau(j)}^i], y_{i\tau(l)}) || 1 \leq i \leq n\}$ , where  $pre_{\tau(l)}^i = [y_{i\tau(1)}, y_{i\tau(2)}, \dots, y_{i\tau(l-1)}]$ , is the concatenation of labels from the first label in the ordering up to the previous label for the  $i^{th}$  data point. Therefore, each dataset,  $\mathcal{D}_{\tau(l)}$ , includes the original input space, as well as the label space from  $\lambda_{\tau(1)}$  up to  $\lambda_{\tau(l-1)}$ , with a target label of  $\lambda_{\tau(l)}$ . This explicitly imposes the dependency of the labels earlier in the chain on the labels later in the chain. Next, for each  $\mathcal{D}_{\tau(l)}$ a binary classifier  $h_{\tau(l)}$  is learned.

For prediction, the chain order generated using the permutation function  $\tau$  is followed. Prediction for t starts from first predicting the  $\lambda_{\tau(1)}$ ,  $\hat{y}_{\tau(1)} = h_{\tau(1)}$ , then this label prediction is concatenated with t and that is fed into next learned model in the chain. To predict  $\lambda_{\tau(l)}$  all the labels  $\lambda_{\tau(1)}$  to  $\lambda_{\tau(l-1)}$  have to be predicted in the chain order first allowing their predictions to be concatenated. Therefore, predicting all the labels is done as follows.  $\hat{y} = \{\hat{y}_{\tau(l)} \| th(h_{\tau(l)}([t, \hat{y}_{\tau(1)}, \dots, \hat{y}_{\tau(l-1)}])), 1 \leq l \leq q\}.$ 

## 3 Proposed Method: ML-KFHE

In this section first the relevant parts of KFHE will be mentioned then ML-KFHE will be described.

#### 3.1 KFHE Algorithm

The discrete Kalman filter is a mathematical framework to estimate an unobservable state of a linear stochastic discrete time controlled process through noisy measurements [29].

The Kalman filter-based Heuristic Ensemble (KFHE) [18] is a multiclass ensemble algorithm which, unlike existing boosting or bagging methods, considers the ensemble to be trained as a hypothesis to be estimated within a hypothesis space. This approach considers the trained classifiers in an ensemble to be noisy measurements which it combines using a Kalman filter. In effect, KFHE behaves like a combination of both boosting and bagging. The remainder of this section describes how a Kalman filter can be used for static state estimation, before describing details of the KFHE approach.

Let there be a state, y, of a linear stochastic system has to be estimated, where y cannot be observed directly. The state of the system can be estimated in two ways. Firstly, given an estimate of the state  $\hat{y}_{t-1}$  with a related variance  $p_{t-1}$  at time step (t-1), a linear model is used to make an *a priori* state estimate  $\hat{y}_t^-$ . The variance related to  $\hat{y}_t^-$  is also updated to  $p_t^-$ . This variance can be imagined as the uncertainty of state. This is known as the *time update* step. Secondly, an external sensor can be used to get an estimate through a measurement,  $z_t$ , of the state with a related variance  $r_t$ , which can also be seen as the uncertainty of the measurement. Given these two noisy state estimates, the *a priori* estimate,  $\hat{y}_t^-$ , its related variance  $p_t^-$ , and the measurement  $z_t$ , its related variance  $r_t$ , the Kalman filter combines them optimally to get an a posteriori state estimate,  $\hat{y}_t$ , which potentially has a lower uncertainty than the previous two. This is known as the *measurement update* step. The Kalman filter iterates through the time update and the measurement update steps. At iteration t, the a priori estimate is used in the measurement update step to get an *a posteriori* estimate, which is fed back to the time update in the next iteration as the *a priori* estimate.

If the state to be estimated is assumed to be static, then the time update step is considered to be non-existent. This kind of scenario can occur in cases when, say, the voltage level of a DC battery or the altitude of a cruising aircraft is being estimated. In both of the cases, the DC voltage and the altitude of the aircraft is supposed to be constant, but unknown. In such cases, the measurement of the static state from a noisy sensor is repeatedly combined using the measurement update step.

The basic idea of KFHE is to view the ideal hypothesis for a specific classification problem as a static state to be estimated in a hypothesis space [30]. As in the above description, when estimating the static state, after T iterations the estimate of the Kalman filter is essentially the combination or an *ensemble* of a sensor output. Similarly, the component classifiers are combined in KFHE using the above principle. The equations used for the algorithm is as follows

$$\hat{y}_{t} = \hat{y}_{t-1} + k_t (z_t - \hat{y}_{t-1})$$
(1)

$$k_t = p_{t-1}/(p_{t-1} + r_t) \tag{2}$$

$$p_t = (1 - k_t)p_{t-1} \tag{3}$$

Here  $z_t$ , the measurement, can be an external source or sensor (voltage or altitude sensor),  $r_t$  is the related measurement variance indicating the uncertainty of the estimate. The  $k_t$  is the Kalman gain, which optimally combines the *a priori* estimate and the measurement. A complete and detailed explanation of Kalman filters can be found in [29, 31].



Fig. 1 The high level interactions between kf-m and kf-w [18]

**Table 1** Intermediate representation of a state for KFHE and ML-KFHE. A trained model is represented using the prediction scores for the classes  $(c_1, c_2, \text{ and } c_3)$  of a given set of datapoints. This representation is used with  $\hat{y}_t$ ,  $z_t$  and  $h_t(\mathcal{D})$ .

	$c_1$	$c_2$	$c_3$
$oldsymbol{x}_1$	0.10	0.89	0.01
$oldsymbol{x}_2$	0.08	0.27	0.65
•	•	•	•
•	•	•	•
$oldsymbol{x}_n$	0.77	0.20	0.03

KFHE has two components which interact with each other. The Kalman filter which estimates the ideal hypothesis (as described above) is called the model Kalman filter, abbreviated as kf-m, estimates the final model or hypothesis by combining multiple noisy measurements. The other component named the weight Kalman filter, kf-w, computes the weight using which the kf-m performs the sampling of training datapoints for the component classifier.

The measurement in this case is defined as,

$$\boldsymbol{z}_t^{(y)} = (h_t(\mathcal{D}) + \boldsymbol{\hat{y}}_{t-1})/2 \tag{4}$$

Where  $h_t = \mathcal{H}(\mathcal{D}, \hat{\boldsymbol{w}}_{t-1})$  is a classifier model trained using algorithm  $\mathcal{H}$  (decision tree, SVM, etc.) using a dataset defined by a set of datapoint weights updated in the previous iteration,  $\hat{\boldsymbol{w}}_{t-1}$ . A datapoint is weighted more if it was misclassified previously, and less if correctly classified which is similar to the approach taken in boosting. Although, unlike AdaBoost [32], the weights for the datapoints are determined by another Kalman filter, the weight Kalman filter or kf-w.

Note that, the ensemble model  $h_t$  cannot directly be used with the equations in Eq. (1) and (4), therefore an intermediate proxy representation is used for the states in kf-m. The intermediate representation of a trained model is the label-wise prediction scores of a given dataset by the model of the corresponding state, as shown in Table 1. Therefore, the intermediate representation of a model (individual or ensemble) would be the prediction  $\hat{y}_t$  as shown in Table 1. For example, the first datapoint has the highest prediction score assigned to class-label  $c_2$ , and thus the first datapoint is considered as a member of class  $c_2$  (among two other potential classes  $c_1$  and  $c_3$ ). This

representation of a model is used as the state in the Kalman filter framework. In the final estimated state the class assignment is done by taking the class with the highest score.

The kf-w estimates  $\hat{w}_t$ , which is a vector of weights to be used by the measurement step of kf-m. kf-w is otherwise identical to kf-m.

The training step stores the component classifiers  $h_t$  and the Kalman gains  $k_t^{(y)}$ . When a new datapoint is encountered during the prediction step the Eq. (1) is repeatedly used using the component classifiers  $h_t$  and the Kalman gains  $k_t^{(y)}$  found during the training stage.

An overall interaction of the kf-m and kf-w is shown in Figure 1. The superscript (y) indicates that the variables are related to kf-m, and the superscript (w) indicates these variables are related to kf-w.  $\hat{y}$  is the state estimate by kf-m, and  $\hat{w}_t$  is estimated by kf-w.

The setting of the measurements and the related errors are the heuristic components of the method, which are set by making assumptions. A detailed explanation of the concept, derivation and explanation of KFHE can be found in [18], and a few related work utilising the Kalman filter based framework to combine classifier models can be found in [33, 34].

#### 3.2 ML-KFHE

In this work, ML-KFHE, proposes a multi-label classification algorithm by combining multiple multi-label classifier models exploiting the sensor fusion properties of the Kalman filter. Depending on which underlying multi-label classifiers are ensembled, two variants are proposed in this work. ML-KFHE-HOMER, which ensembles multiple HOMER models and ML-KFHE-CC which ensembles multiple CC models. HOMER is a multi-label classification method that has been shown to have competitive performance with other leading approaches [16, 17]. HOMER was selected for this task because it has a lower training time which makes it suitable for this purpose to train more ensemble components in a shorter time. Also, CC was chosen as the existing ECC [4] already attains very good classification performance [15] as well as takes label associations into consideration.

As explained in Section 3.1, there are two components of KFHE: kf-m that estimates the hypothesis, and kf-w that computes the weights of the training datapoints during each measurement. To make KFHE work in a multi-label setting, the measurements of the kf-m and kf-w steps were adapted in this work.

For ML-KFHE, the measurement at each step is the average of a trained multi-label classifier and the previous estimate of the ensemble as shown in Eq. (4). The related measurement uncertainty  $r_t^{(y)}$  is the Hamming loss (*hloss*) [35] of the trained model. Each multi-label model at every step is trained on different weights,  $\hat{w}_t$ , assigned to different datapoints, where the weights are determined by the *kf-w* component. The *kf-w* estimates one single vector of weights  $\hat{w}_t$  using which a sampling with replacement of the training dataset

is done. Although the measurement  $\boldsymbol{z}_{t}^{(w)}$  for kf-w is taken as per-datapoint weighted Hamming loss, which can be defined as follows

$$\boldsymbol{z}_{t}^{(w)} = [z_{ti}^{(w)} \| z_{ti}^{(w)} = \hat{w}_{ti} \times exp(hloss(\boldsymbol{x}_{i}, \boldsymbol{l}_{i})) \ 1 \le i \le n]$$
(5)

Eq. (5) uses Hamming loss and the exponential function (KFHE-e variant) to highlight misclassified datapoints in the measurement which will later be used by the measurement update step to get the weights  $\hat{\boldsymbol{w}}_t$  to be used in kf-m in the next iteration. In this case the related uncertainty is calculated as in KFHE.

The model  $h_t$  in this case is a trained multi-label classifier model  $\mathcal{H}(\mathcal{D}, \hat{w}_{t-1}, \mathcal{M})$ . Here  $\mathcal{H}$  is CC or HOMER (the underlying multi-label classifier) algorithm and  $\mathcal{M}$  is the hyperparameters of  $\mathcal{H}$ . To weight the datapoints for training, the multi-label lassifiers are trained using samples using the distribution  $\hat{w}_{t-1}$ , the last updated weights. Based on the underlying multi-label classification algorithm two variants are presented in this work.

- ML-KFHE-HOMER: This variant uses HOMER as the underlying classifier. To train each component HOMER classifier the following three hyperparameter are modified.  $\mathcal{M} = \{\mathcal{C}, k, \phi\}$ . Here  $\mathcal{C}$  is the clustering algorithm used by HOMER is randomly selected from  $\{random, k\text{-means}, balanced k\text{-means}\}$ , k is the number of clusters which is randomly selected too. Also, the kernel  $\phi$  of the underlying SVM used by HOMER, is also selected randomly.
- ML-KFHE-CC: This variant uses CC as the underlying classifier and the hyperparameters adjusted randomly to increase diversely of the models in this case are  $\mathcal{M} = \{\mathcal{O}, \phi\}$ . Here,  $\mathcal{O}$  is the chain order for a component CC classifier, which is selected randomly. Also, like before, the kernel type  $\phi$  of underlying SVM used by CC is selected randomly.

Next, the measurement is done using Eq. (4).

The above method is applied to increase diversity of the models. The reason to increase diversity, for example, the HOMER models in the case of ML-KFHE-HOMER by randomly selecting the clustering algorithm, cluster size, and the SVM kernel type is as follows. Given a set of different HOMER models trained using different hyperparameters, many of them may lead to a poor measurement. The ML-KFHE framework combines the measurements based on the measurement errors. If the measurement uncertainty  $r_t^{(y)}$  is higher than the uncertainty of the ensemble found up to the tth iteration  $p_t^{(y)}$ , then the measurement is weighted less and the Kalman gain is lower than 0.5, and when the measurement error is lower the measurement is incorporated more, as a result of the Kalman gain being greater than 0.5. Therefore, based on this property, the HOMER models which have a poor performance will have a much less impact on the entire ensemble, whereas a more accurate HOMER model will have more impact on the entire ensemble. This also applies on the ML-KFHE-CC version, where the diversity is induced by selecting the random chain order and the randomly selected underlying SVM kernel.

The values of  $\hat{\boldsymbol{y}}_0$ ,  $p_0^{(y)}$  and  $\hat{\boldsymbol{w}}_0$ ,  $p_0^{(w)}$  have to be initialised.  $\hat{\boldsymbol{y}}_0$  is initialised using a single  $\mathcal{H}$  classifier model,  $h_0$ . The value of  $p_0^{(y)}$  is set to 1 indicating maximum uncertainty. Equal weight is given to every point in  $\hat{\boldsymbol{w}}_0$ , and  $p_0^{(w)}$  is also initialised with 1.

Algorithm 1 shows the ML-KFHE training algorithm . The superscripts (y) and (w) indicate that the corresponding variables are related to kf-m and kf-w respectively. On Lines 7-17 the different hyperparameters of HOMER are selected randomly. Next, the component classifier model is trained on Line 18, and the measurement is done on Line 19. Line 21 computes the Kalman gain,  $k_t^{(y)}$ , for kf-m and Line 22 computes the proxy representation of the ensemble  $\hat{y}_t$ , based on the ML-KFHE ensemble predictions on the training dataset. The kf-w steps are similar and are performed on Lines 25-29f. The process runs until a maximum number of ensemble iterations T.

The prediction algorithm is the same as for KFHE and is shown in Algorithm 2. Here the trained models and the Kalman gain values learned during the training along with a new query datapoint is given. Using the models in Line 5 the Kalman gain is repeatedly used to combine the measurements on Line 4. After T iterations the predicted labels for the new datapoint d, the estimate  $\hat{y}_T^{(y)}$  are returned. To find the label assignments, these scores are thresholded at 0.5.

Algorithm 1 ML-KFHE training 1: procedure TRAIN( $\mathcal{D} = \{(\boldsymbol{x}_i, \boldsymbol{l}_i) || 1 \le i \le n\}, T$ )  $p_0^{(w)} = 1, \, \hat{w}_0 = [1/n, \dots, 1/n]$ 2:  $\hat{h}_t = \mathcal{H}(\mathcal{D}, \hat{\boldsymbol{w}}_0, \mathcal{C}, k, \phi), \ \hat{\boldsymbol{y}}_0 = \hat{h}_0(\mathcal{D})$ 3: t = 14: for t < T do 5: ⊳ kf-m 6: if  $\mathcal{H}$  is HOMER then 7:  $\mathcal{M} = \{\mathcal{C}, k, \phi\}$ 8: Choose  $\mathcal{C} \in \{k\text{-means, balanced } k\text{-means, random}\}$  randomly 9: Choose  $k \in \{2, \ldots, \lceil \sqrt{\|\mathcal{L}\|} \rceil\}$  randomly 10: Choose  $\phi \in \{linear, radial\}$  randomly 11: else if  $\mathcal{H}$  is CC then 12:Select  $\mathcal{M} = \{\mathcal{O}, \phi\}$ 13: $\mathcal{O}$  is a random label ordering 14:Choose  $\phi \in \{linear, radial\}$  randomly 15:end if 16 $h_t = \mathcal{H}(\mathcal{D}, \hat{w}_{t-1}, \mathcal{M})$ 17: $z_t^{(y)} = (h_t(\mathcal{D}) + \hat{y}_{t-1})/2$ ▷ Measurement 18:  $r_t^{(y)} = hloss(\mathcal{D}, \boldsymbol{z}_t^{(y)})$ 19: 
$$\begin{split} \dot{k}_{t}^{(y)} &= p_{t-1}^{(y)} / (p_{t-1}^{(y)} + r_{t}^{(y)}) \\ \dot{\boldsymbol{y}}_{t} &= \hat{\boldsymbol{y}}_{t-1} + k_{t}^{(y)} (\boldsymbol{z}_{t}^{(y)} - \hat{\boldsymbol{y}}_{t-1}) \end{split}$$
⊳ Kalman gain 20: $\triangleright$  Measurement update 21: $p_{t}^{(y)} = (1 - k_{t}^{(y)})p_{t-1}^{(y)}$ 22. ⊳ kf-w 23: $\begin{aligned} & \boldsymbol{z}_{t}^{(w)} = [z_{ti}^{(w)} \| \overline{z_{ti}^{(w)}} = \hat{w}_{ti} \times exp(hloss(\boldsymbol{x}_{i}, \boldsymbol{l}_{i})) \ 1 \leq i \leq n] \\ & r_{t}^{(w)} = r_{t}^{(m)} \\ & k_{t}^{(w)} = p_{t-1}^{(w)} / (p_{t-1}^{(w)} + r_{t}^{(w)}) \\ & \triangleright \mathbf{F}_{t}^{(w)} \end{aligned}$ 24 25: ⊳ Kalman gain 26: $\hat{w}_t = \hat{w}_{t-1} + k_t^{(w)} (z_t^{(w)} - \hat{w}_{t-1})$  $\triangleright$  Measurement update 27: $p_t^{(w)} = (1 - k_t^{(w)}) p_{t-1}^{(w)}$ 28:29: t = t + 130: 31: end for **return** ({ $h_t, k_t^{(y)} || \forall_{1 \le t \le T}$ }) 32: 33: end procedure

$\triangleright$ Measurement
$\triangleright$ Measurement update

• E-HOMER: A simple bagged version of HOMER, E-HOMER, is also introduced in this section mainly with the intension to compare with ML-KFHE-HOMER to evaluate the effectiveness of ML-KFHE-HOMER. The hyperparameters for each HOMER model in the ensemble (the cluster type, the number of clusters, and the type of underlying SVM kernel) are all selected randomly, as in ML-KFHE-HOMER. The difference between E-HOMER and ML-KFHE-HOMER is that the combination of ML-KFHE-HOMER uses the ML-KFHE framework to ensemble the HOMER component classifier models and E-HOMER ensembles the component HOMER classifiers using simple bagging. Algorithm 3 describes the E-HOMER training process, and Algorithm 4 describes the prediction process.

Algorithm 3 E-HOMER training

1: procedure TRAIN( $\mathcal{D} = \{(\boldsymbol{x}_i, \boldsymbol{l}_i) || 1 \le i \le n\}, T$ )  $h_t = \mathcal{H}(\mathcal{D}, \hat{\boldsymbol{w}}_0, \mathcal{C}, k, \phi), \ \hat{\boldsymbol{y}}_0 = h_0(\mathcal{D})$ 2: t = 13: for t < T do 4: Randomly select  $\mathcal{C}$ , k and  $\phi$ , where 5: $C \in \{k\text{-means, balanced } k\text{-means, random}\},\$ 6:  $k \in \{2, \ldots, \lceil \sqrt{\|\mathcal{L}\|} \},\$ 7:  $\phi \in \{linear, radial\}$ 8.  $\boldsymbol{b} = bootstrap \quad sample(2 \times n)$ 9:  $h_t = \mathcal{H}(\mathcal{D}, \boldsymbol{b}, \mathcal{C}, k, \phi)$ 10: t = t + 111: end for 12: 13:return  $(\{h_t \| \forall_{1 \le t \le T}\})$ 14: end procedure

#### Algorithm 4 E-HOMER prediction

```
1: procedure PREDICT(\boldsymbol{d}, \{h_t \| \forall_{1 \leq t \leq T}\}, T)

2: \hat{\boldsymbol{y}}_0^{(y)} = h_0(\boldsymbol{x}), t = 1

3: for t \leq T do

4: \hat{\boldsymbol{y}}_t = \hat{\boldsymbol{y}}_{t-1} + h_t(\boldsymbol{d})

5: t = t + 1

6: end for

7: return (\frac{\hat{\boldsymbol{y}}_T^{(y)}}{T})

8: end procedure
```

## 4 Experiment

To evaluate the effectiveness of ML-KFHE, experiments were performed on thirteen well-known multi-label benchmark datasets<sup>1</sup> listed in Table 2. In Table 2, different properties of the multi-label datasets are summarised. *Instances, Inputs* and *Labels* are the number of datapoints, the dimension of the datapoints, and the number of labels, respectively. *Labelsets* indicates the number of unique combinations of labels. *Cardinality* measures the average number of labels assigned to each datapoint and *MeanIR* [36] indicates the degree of imbalance of the labels, where higher values indicate higher imbalance.

Dataset	Instances	Inputs	Labels	Labelsets	Cardinality	MeanIR
flags	194	26	7	24	3.392	2.255
yeast	2417	103	14	77	4.237	7.197
scene	2407	294	6	3	1.074	1.254
emotions	593	72	6	4	1.869	1.478
medical	978	1449	45	33	1.245	89.501
enron	1702	1001	53	573	3.378	73.953
birds	322	260	20	55	1.503	13.004
genbase	662	1186	27	10	1.252	37.315
cal500	502	68	174	502	26.044	20.578
llog	1460	1004	75	189	1.180	39.267
foodtruck	407	21	12	116	2.290	7.094
Water_quality	1060	16	14	852	5.073	1.767
PlantPseAAC	978	440	12	32	1.079	6.690

Table 2 Multi-label datasets used in this work

A brief description of the datasets are as follows:

• *flags* [37]: Different properties of flags are used to predict seven different colours of flags. Each colour is considered as a label.

 $<sup>^{1}</sup> Dataset \qquad sources: \\ http://www.uco.es/kdis/mllresources/$ 

- *yeast* [38]: Is a widely used datasest from the biological domain where genes are associated one or more of 14 different biological functions (labels).
- *scene* [2]: Several pictures of the scene is processed picture of scenary. Each scene can have one or more of the following labels: "beach", "sunset", "field", "fall-foliage", "mountain" and "urban".
- *emotions* [39]: Each datapoint represents a piece of music. Each instance can be labelled with six emotions: "sad-lonely", "angry-aggressive", "amazed-surprised", "relaxing-calm", "quiet-still" and "happy-pleased".
- *medical* [40]: Each datapoint represents a document which includes a brief free text summary of patient symptom history and their prognosis, with 45 ICD-9-CM (International Classification of Diseases, Ninth Revision, Clinical Modification) codes<sup>2</sup>.
- enron [41]: This dataset is a version of Enron email corpus <sup>3</sup> found in [41], where each leaf categories in the original Enron corpus hierarchy is considered as a label, making a total of 53 labels.
- *birds* [42]: Each of the datapoint is a bird song recording and other natural sound recordings which can be labelled with 19 bird species or none of them. For this thesis, an additional label is added indicating "not a bird song" is added to avoid empty prediction problem [43].
- genbase [44]: Similar to the yeast dataset, this dataset is from the microbiology domain concerned with gene functions, where each gene can be labelled with 27 labels.
- cal500 [45]: Each datapoint is a western song, which were hand annotated with genre, emotions, etc. Then each of the song/datapoint is labelled with 174 possible "musically relevant" labels in [45] and named Computer Audition Lab 500 dataset.
- *llog* [41]: This dataset is a version of data found in the Language Log Forum <sup>4</sup> used in [41]. This assigns 75 possible topics (labels) to each datapoint describing the document.
- foodtruck [46]: This dataset was created by using the answers provided by the 407 survey on 21 objective questions about food truck preferences of the participants and their users' profile. Food truck preference questions were ingredients, place to sit, menu, hygiene preferences, along with personal users' information questions were age group, average income etc. The responses were recorded as categorical attributes.
- *Water\_quality* [47]: This dataset is used to predict the quality of water of Slovenian rivers, using 16 attributes such as, the temperature, pH, hardness, nitrous oxide or carbon di-oxide.
- *PlantPseAAC* [48]: This dataset contains 978 sequences for Plant species. Gene ontology, amino acids, pseudo-amino acids and diptide components are provided. There are 12 labels indicating subcellular locations (cell membrace, cell wall, chloroplast, cytoplasm, endoplasmic reticulum, extracellular, golgi

 $<sup>^{2}</sup> https://www.cdc.gov/nchs/icd/icd9cm.htm$ 

<sup>&</sup>lt;sup>3</sup>http://www.cs.cmu.edu/ enron/

 $<sup>^{4}</sup>$  http://languagelog.ldc.upenn.edu/nll/

apparatus, mitochondrion, nucleus, peroxisome, plastid, and vacuole), which needs to be predicted.

Label-based macro-averaged F-Score [35] was used to measure the performance of models in these evaluations. This was chosen over Hamming loss, which has been used in several previous studies (e.g. [49–51]) because in the highly imbalanced multi-label datasets used in this study (see the high MeanIR scores for several datasets in Table 2) with Hamming loss performance on the majority classes may overwhelm the performance of the minority class.

For all evaluations 2 times 5 fold cross-validation experiments were performed. While generating the folds for cross-validation method a multi-label specific stratification method, iterative stratification [52], was used.

Performance of ML-KFHE-HOMER and ML-KFHE-CC was compared with various algorithms. The hyperparameter configurations of each will be described next.

ECC, RAkEL, state-of-the-art ensemble based multi-label classifiers, AdaBoost.MH, RF-PCT and E-HOMER (a bagging-based ensemble using HOMER, Section 3.2). ECC and RAkEL was specifically compared as they were the top performer in the extensive experiment in [15]. The hyperparameter setting for RAkEL (named RAkEL2) was used as in [15] as it is shown to perform very well. RAkEL2 is RAKEL with label subset size of 3 and number of such random label overlapping subsets is 2q, where the q is the number of labels in the corresponding dataset. Support vector machine (SVM) is used as the underlying classifiers for ECC, RAkEL and all algorithms.

The compared hyperparameter tuned individual classifier (non-ensemble) algorithms compared are as follows. For, CC, HOMER-K (using k-means clustering) and HOMER-B (using balanced clustering) models were included to understand how much the ensembles led to improved performance over single base models when they are hyperparameter tuned for performance for each dataeset. During all the experiments, the sample size of training data (bag fraction) is 2n, or twice the size of original number of datapoints for all the ensemble algorithms. The cluster size for HOMER was selected using the best values found in the benchmark experiments in [17]. The HOMER and CC models used support vector machines (SVM) as their underlying learner (as all the others), as they have proved to perform very well [15–17].

Now the configuration of the proposed methods will be explained. At each iteration of E-HOMER and ML-KFHE-HOMER, for the underlying HOMER algorithm, the type of clustering, C, was selected randomly from  $\{balanced \ k-means, \ k-means, \ random\}$ , the number of clusters k was selected randomly from the range  $k \in \{2, \ldots, \lceil \sqrt{\|\mathcal{L}\|} \rceil\}$ , and the SVM kernel  $\phi$  from  $\{linear, radial\}$ , at each ensemble iteration. For ML-KFHE-CC, the kernel types  $\phi$  for each of the base SVM models was selected at each ensemble iteration randomly, from  $\{linear, radial\}$ . The chain ordering of each component CC classifier was also selected randomly at each ensemble iteration.

For ECC and E-HOMER the bootstrap sample was selected to be twice the size of the training dataset, to keep it consistent with the ML-KFHE variants.

A total of 100 component classifiers were trained for all the ogher ensemble algorithms. Therefore, the experimental environment were kept identical for all ensemble methods for a fair comparison.

ML-KFHE and E-HOMER are implemented in R scripting language<sup>5</sup>, AdaBoost.MH and RF-PCT was used from the MULAN [53] and CLUS <sup>6</sup> libraries respectively, and for the other methods the *utiml* library [54] is used.

## 5 Results

 $<sup>^5\</sup>mathrm{A}$  version of ML-KFHE and E-HOMER is available at: https://github.com/phoxis/kfhe-homer  $^6\mathrm{https://dtai-static.cs.kuleuven.be/clus/}$ 

	aBoost.MH	$796 \pm 0.08 (10)$	$222 \pm 0.00 (10)$	$00 \pm 0.00 (10)$	$573 \pm 0.02 \ (10)$	$33 \pm 0.03 (8)$	$190 \pm 0.03 (10)$	$105 \pm 0.08 (10)$	$593 \pm 0.04 (9)$	$502 \pm 0.01 (9)$	$113 \pm 0.06 (3)$	$700 \pm 0.00 (10)$	$330 \pm 0.02 (10)$	$83 \pm 0.03 (10)$	9.15
	F-PCT Ads	$6473 \pm 0.05$ (5 ) 0.57	$3671 \pm 0.01 (8) 0.12$	$7161 \pm 0.01$ (8) 0.00	$6617 \pm 0.02 (8) 0.03$	$3356 \pm 0.05 (10) 0.49$	$1760 \pm 0.04 (9) 0.14$	$2176 \pm 0.04 (9)$ 0.11	$2333 \pm 0.08 (10) 0.25$	$1231 \pm 0.03$ (5 ) 0.03	$2460 \pm 0.03 (1)$ 0.24	$1782 \pm 0.03$ (5 ) 0.07	$5388 \pm 0.01$ (5 ) 0.08	$0117 \pm 0.00 (9) 0.00$	2.08
	OMER-K R	$6484 \pm 0.03 (4) = 0.03$	$3524 \pm 0.02 (9)$ 0.	$2011 \pm 0.03 (9) = 0.03$	$6809 \pm 0.04$ (5) 0.	$3674 \pm 0.04 (9) = 0.04$	$2123 \pm 0.02$ (5) 0.	$3203 \pm 0.03$ (7) 0.	$8810 \pm 0.03$ (7) 0.	$0642 \pm 0.02 (7) = 0.02$	$1683 \pm 0.02 (10) 0.01$	$2335 \pm 0.03 (2)$ 0.	$5743 \pm 0.02 (4)$ 0.	$0969 \pm 0.02 (7) = 0.02$	6.54
	RAKEL2 H	$0.5891 \pm 0.05$ (9) 0.	$0.4288 \pm 0.01$ (6) 0.	$0.7990 \pm 0.02$ (4) 0.	$0.6915 \pm 0.03$ (3) 0.	$0.6126 \pm 0.03$ (5) 0.	$0.1877 \pm 0.01$ (8) 0.	$0.3161 \pm 0.05$ (8) 0.	$0.9217 \pm 0.03$ (5) 0.	$0.0537 \pm 0.00$ (8) 0.	$0.2404 \pm 0.02$ (4) 0.	$0.1099 \pm 0.01$ (9) 0.	$0.4239 \pm 0.01$ (9) 0.	$0.0995 \pm 0.01 (6) 0.01 (6)$	6.46
	CC	$0.5957 \pm 0.05$ (8) (	$0.4565 \pm 0.01$ (3) (	$0.7818 \pm 0.02$ (5) (	$0.6595 \pm 0.04$ (9) (	$0.6114 \pm 0.02$ (6) (	$0.1985 \pm 0.02$ (6) (	$0.3297 \pm 0.04$ (5) (	$0.9245 \pm 0.02$ (4) (	$0.0859 \pm 0.01$ (6) (	$0.2454 \pm 0.02$ (2) (	$0.1280 \pm 0.02$ (8) (	$0.4415 \pm 0.01$ (8) (	$0.0968 \pm 0.02$ (8) (	6.00
	HOMER-B	$0.6681 \pm 0.04 (3)$	$0.3812 \pm 0.02$ (7)	$0.7777 \pm 0.02$ (7)	$0.6964 \pm 0.03$ (2)	$0.5505 \pm 0.03$ (7)	$0.1888 \pm 0.01$ (7)	$0.3256 \pm 0.05 \ (6)$	$0.7534 \pm 0.06$ (8)	$0.0331 \pm 0.01 \ (10)$	$0.2253 \pm 0.01 (8)$	$0.1692 \pm 0.02 (6)$	$0.5800 \pm 0.01$ (3)	$0.1708 \pm 0.03$ (3)	5.92
	ECC	$0.6387 \pm 0.04 (6)$	$0.4403 \pm 0.02 (4)$	$0.7806 \pm 0.02 (6)$	$0.6731 \pm 0.03 \ (6.5)$	$0.6274 \pm 0.02 (2)$	$0.2435 \pm 0.03 (4)$	$0.3463 \pm 0.06 (4)$	$0.9293 \pm 0.02 (2)$	$0.1310 \pm 0.01 (4)$	$0.2236 \pm 0.01 (9)$	$0.1601 \pm 0.01$ (7)	$0.4527 \pm 0.01$ (7)	$0.1694 \pm 0.03 (4)$	5.04
	AL-KFHE-CC	$0.6271 \pm 0.04$ (7)	$0.4661 \pm 0.01 (2)$	$0.8005 \pm 0.01$ (3 )	$0.6731 \pm 0.03 \ (6.5)$	$0.6235 \pm 0.02 (4)$	$0.2458 \pm 0.02$ (3 )	$0.3586 \pm 0.04 (3)$	$0.9273 \pm 0.03$ (3)	$0.1341 \pm 0.01$ (3)	$0.2267 \pm 0.01$ (7 )	$0.2090 \pm 0.04 (3)$	$0.4978 \pm 0.01 (6)$	$0.1652 \pm 0.03$ (5 ) (	4.27
	E-HOMER N	$0.6710 \pm 0.05$ (2) (	$0.4400 \pm 0.04$ (5) (	$0.8008 \pm 0.02$ (2) (	$0.6891 \pm 0.03$ (4) (	$0.6254 \pm 0.03$ (3) (	$0.2587 \pm 0.02$ (2) (	$0.3834 \pm 0.04$ (2) (	$0.8911 \pm 0.03$ (6) (	$0.1455 \pm 0.01$ (2) (	$0.2308 \pm 0.01$ (6) (	$0.1976 \pm 0.01$ (4) (	$0.5836 \pm 0.01$ (1) (	$0.1749 \pm 0.02$ (2) (	3.15
algorithms.	<b>IL-KFHE-HOMER</b>	$.6862 \pm 0.02 (1)$	$.4899 \pm 0.01 (1)$	$.8090 \pm 0.02 (1)$	$.7046 \pm 0.03 (1)$	$.6387 \pm 0.02 (1)$	$.2612 \pm 0.02 (1)$	$.3928 \pm 0.05 (1)$	$.9402 \pm 0.03 (1)$	$.1458 \pm 0.01 (1)$	$.2315 \pm 0.01 (5)$	$.2384 \pm 0.02 (1)$	$.5822 \pm 0.01 (2)$	$.1760 \pm 0.02 (1)$	1.38
rank of each	N.	flags 0.	yeast 0.	scene 0.	emotions 0.	medical 0.	enron 0.	birds 0.	genbase 0.	cal500 0.	llog 0.	foodtruck 0.	Water_quality 0.	PlantPseAAC 0.	Avg. rank

standard deviations. The rank of each score for a dataset across the algorithms compared is shown in parenthesis. The last row shows the average Table 3 Experiment results. Values in cells are mean label-based macro-averaged F-Scores from the cross-validation experiments, and their

18

ML-KFHE

The results are presented and analysed in this section. First results from the experiments are presented in Section 5.1. Next, in Section 5.2, a detailed statistical analysis is performed to understand the overall differences between the algorithms, as well as the per-dataset performance differences.

### 5.1 Performance Comparison

Table 3 shows the results of the experiments performed. The columns indicate the algorithms and the rows indicate the datasets. In each cell, the mean and standard deviation label-based macro-averaged F-Score (higher values are better) across the cross-validation performed are shown. The values in the parenthesis indicate the relative ranking (lower values are better) of the algorithm with respect to the corresponding dataset. The last row of Table 3 indicates the overall average ranks of the algorithms compared.

Performance of the compared algorithms in Table 3 is interpreted in three ways. Firstly, overall comparison with all the algorithms. Secondly, to compare the individual classifier models compare with ML-KFHE. Finally, and most importantly, to compare ML-KFHE-HOMER with E-HOMER and compare ML-KFHE-CC with ECC directly to understand the effectiveness of the ML-KFHE ensembling.

Table 3 shows the overall picture, where ML-KFHE-HOMER attains the best average rank of 1.38. In fact, ML-KFHE-HOMER attained the top rank for all the datasets, except for *llog* and *Water\_quality* where it got fifth and second rank respectively. E-HOMER attained the second best overall average rank of 3.15, whereas ML-KFHE-CC attained the third best overall rank of 4.27. ECC attained the fourth best overall average rank of 5.04. HOMER-B comes next with overall average ranks of 5.92. The classifiers, CC, RAkEL2, HOMER-K comes next with average ranks of 6.00, 6.46 and 6.54 respectively. RF-PCT was not able to perform well and was ranked 7.08. The worst performing multi-label classification models in this case was by AdaBoost.MH with an average rank of 9.15. Similar results using RF-PCT and AdaBoost.MH was also found in [15].

The difference between E-HOMER and ML-KFHE-HOMER is the aggregation method of the component HOMER classifier models, and ML-KFHE-HOMER has performed better than E-HOMER in all the cases. E-HOMER has similar benefits and drawbacks of a bagged method. E-HOMER models could be trained in parallel as the component classifiers do not depend on each other. Unlike ML-KFHE, E-HOMER gives equal weights to all the component classifiers, good or bad, when combining the models, due to this E-HOMER has a poor predictive performance compared to ML-KFHE-HOMER.

Similarly, in the case ML-KFHE-CC, it performed better than ECC. As in both the cases the difference between E-HOMER, ML-KFHE-HOMER and ECC, ML-KFHE-CC is the combination method of the component classifiers, and all the other processes are kept identical, this demonstrates



Fig. 2 Critical difference plot of Friedman rank sum test with Finner *p*-value correction. The scale indicates the average ranks. The methods which are not connected with the horizontal lines are significantly different with a significance level of 0.05.

the effectiveness of ML-KFHE method. Therefore, it can be concluded that ML-KFHE is a better way of ensembling classifiers.

On the other hand, interestingly, E-HOMER has performed better in almost all cases compared to a single HOMER model, as well as having performed better than ECC and ML-KFHE-CC, which demonstrates the effectiveness of ensembling the HOMER method in general. From this experiment it can also be concluded that the KFHE-ML method's overall performance is dependent on the underlying multi-label classifier used, but KFHE-ML method would almost always be able to improve the classification performance when compared to a bagged combination of the underlying classifiers.

#### 5.2 Statistical Significance Testing and Further Analysis

To further analyse the overall difference of the methods over the different datasets and the differences between per-dataset performances statistical significance tests are performed.

#### 5.2.1 Multiple Classifier Comparison

**Table 4** Upper diagonal: win/lose/tie. Lower diagonal: Results of the Friedman rank test with Finner *p*-value correction. \*  $\alpha = 0.1$ , \*\*  $\alpha = 0.05$  and \*\*\*  $\alpha = 0.01$ 

	ML-KFHE-HOMER	E-HOMER	KFHE-CC	ECC	HOMER-B	CC	RAkEL2	HOMER-K	RF-PCT	AdaBoost.MH
ML-KFHE-HOMER		12/1/0	13/0/0	13/0/0	13/0/0	12/1/0	12/1/0	13/0/0	12/1/0	12/1/0
E-HOMER	0.2166		10/3/0	10/3/0	12/1/0	10/3/0	10/3/0	12/1/0	12/1/0	12/1/0
ML-KFHE-CC	0.0384 **	0.4321		8/5/0	9/4/0	12/1/0	11/2/0	9/4/0	10/3/0	12/1/0
ECC	0.0095 ***	0.1985	0.6315		7/6/0	10/3/0	10/3/0	9/4/0	9/4/0	12/1/0
HOMER-B	0.0007 ***	0.0485 **	0.2563	0.5123		5/8/0	7/6/0	9/4/0	10/3/0	11/2/0
CC	0.0007 ***	0.0432 **	0.2359	0.4812	0.9517		9/4/0	7/6/0	7/6/0	13/0/0
RAkEL2	0.0002 ***	0.0184 **	0.1319	0.3007	0.6933	0.7139		7/6/0	8/5/0	12/1/0
HOMER-K	0.0002 ***	0.0163 **	0.1195	0.2779	0.6664	0.6933	0.9517		9/4/0	11/2/0
RF-PCT	0.0000 ***	0.0043 ***	0.0485 **	0.1453	0.4321	0.4512	0.6664	0.6933		11/2/0
AdaBoost.MH	0.0000 ***	0.0000 ***	0.0003 ***	0.0023 ***	0.0208 **	0.0236 **	0.0518 *	0.0583 *	0.1453	

A Friedman rank test was performed with the Finner *p*-value correction [55]. The results of this evaluation is summarised in Figure 2, where the scale indicates the average ranks and if the methods are not connected with a horizontal line then they are significantly different over different datasets with a significance level of 0.05. This shows that ML-KFHE-HOMER was significantly better than all the methods except E-HOMER in which case the null hypothesis of Friedman rank test could not be rejected with a significance level of 0.05. Overall, ML-KFHE-HOMER attained better ranks in all the datasets.

An overall pairwise table of the p-values of the Friedman test is shown in Table 4. The lower diagonal of Table 4 a value in a cell is the p-values of the Friedman rank test with the Finner p-value correction for of the corresponding pair of algorithms in the rows and columns. Also, in the upper diagonal of the Table 4 each cell has the win/lost/tie count of the algorithm in the corresponding row, over the algorithms in the corresponding column.

#### 5.2.2 Per-dataset Isolated Pairwise Comparison

	ML-KFHE-HOMER vs.									
	ML-KFHE-CC	E-HOMER	ECC	HOMER-B	CC	RAkEL2	HOMER-K	RF-PCT	AdaBoost.MH	
flags	⊻		Ľ		Ľ	Ľ	Ľ	×	⊻	
yeast	1/2	Ľ	Ľ	Ľ	Ľ	Ľ	Ľ	Ľ	1/2	
scene	¥	ĸ ′	Ľ	Ľ	×		Ĺ	Ľ	Ľ	
emotions	1/2		×		×		×	Ľ	1/2	
medical	×	×		Ľ	×	Ľ	Ĺ	Ľ	Ľ	
enron	×		ĸ	Ľ	Ľ	Ľ	Ľ	Ľ	1/2	
birds	×		×	×	×	Ľ	Ľ	Ľ	1/2	
genbase		Ľ		Ľ		×	Ĺ	Ľ	Ľ	
cal500	×		$\checkmark$	Ľ	Ľ	Ľ	Ľ	$\checkmark$	1/2	
llog					× _		Ľ			
foodtruck	×	Ľ	Ľ	Ľ	Ľ	L		Ľ	Ľ	
Water_quality	¥		Ľ		Ľ	L		Ľ	Ľ	
PlantPseAAC	ĸ				Ľ	Ľ	Ľ	Ľ	Ľ	

In the previous section the overall performance of the algorithms over different datasets were analysed. Now, how different the performance (labelbased macro-averaged F-Score) on each individual datasets are, will be analysed. To understand if ML-KFHE-HOMER and ML-KFHE-CC did attain significantly different (better or worse) results than the other methods for each dataset, a two-tailed paired Wilcoxon's signed rank sum test [55] was performed over the folds of each cross-validation experiment. Two tests were done. First, ML-KFHE-HOMER was set as the control method and compared to the other methods per dataset. Next, ML-KFHE-CC was set as the control method and compared with the other methods per dataset.

**Table 6** Per-dataset comparison with ML-KFHE-CC as the control method of two tailed Wilcoxon's signed rank sum test. The symbols " $\not \subseteq$ ", " $\not \subseteq$ " and " $\not \subseteq$ " indicate that the method in the column was *significantly worse* than ML-KFHE-CC at a significance level of 0.01, 0.05 and 0.1, respectively on the specific dataset. The " $\bigtriangledown$ ", " $\rightthreetimes$ " and " $\twoheadleftarrow$ ", "symbols indicate that the method in the column was *significantly better* than ML-KFHE-CC with a significance level of 0.01, 0.05 and 0.1, respectively on the specific dataset. The " $\backsim$ ", " $\circlearrowright$ " and " $\twoheadleftarrow$ ", "symbols indicate that the method in the column was *significantly better* than ML-KFHE-CC with a significance level of 0.01, 0.05 and 0.1, respectively on the specific dataset. The "." indicates that the null hypothesis for the mentioned two tailed Wilcoxon's signed rank sum test could not be rejected with any significance level.

	ML-KFHE-CC vs.									
	ML-KFHE-HOMER	E-HOMER	ECC	HOMER-B	CC	RAkEL2	HOMER-K	RF-PCT	AdaBoost.MH	
flags	R	5		۲,		×	ĸ		×	
yeast	R	×	×	Ľ	×	Ľ	Ľ	Ľ	¥	
scene	ĸ		×	Ľ	×		Ľ	Ľ	1/2	
emotions	~			r.,					¥	
medical	۲,			Ľ			Ľ	¥	¥	
enron	×	*		Ľ	Ľ	L	Ľ	¥	¥	
birds	5			×	×′	×	×	Ľ	¥	
genbase		×		Ľ			×	¥	¥	
cal500	×	尺		Ľ	L	Ľ	Ľ	*	Ľ	
llog					ĸ	r,	Ľ	ĸ		
foodtruck	۲,		Ľ	×	Ľ	L	5	×	¥	
Water_quality	R	1	Ľ	R.	Ľ	Ľ	R.	1	¥	
PlantPseAAC	×				Ľ	Ľ	Ľ	Ľ	¥	

The result from the first experiment where KFHE-ML-HOMER is the control method, is shown in Table 5 and the results of the second experiment with ML-KFHE-CC is the control method, is shown in Table 6. This means, when for example the *yeast* dataset, there are a set of scores from each fold of the 2 times 5 fold crossvalidation experiment (total 10 label-based macro-averaged F-Scores) for ML-KFHE-HOMER and ECC each. The arrows indicate the result of comparing these two sets of scores using the Wilcoxon's signed rank sum test indicating that ECC was significantly worse than ML-KFHE-HOMER.

method in the column was significantly worse than the control method in the corresponding table at a significance level of 0.01, 0.05 and 0.1, respectively on the specific dataset. The  $\mathbb{R}, \mathbb{K}, \mathbb{K}, \mathbb{K}$  symbols indicate that the method in the column was *significantly better* than control method in the corresponding table with a significance level of 0.01, 0.05 and 0.1, respectively on the specific dataset. The "." symbol in both Tables 5 and 6 indicates that the null hypothesis for the mentioned two tailed Wilcoxon's signed rank sum test could not be rejected with any significance level. For example, in Table 5 (with ML-KFHE-HOMER as the control method), ECC was significantly worse than ML-KFHE-HOMER with a significance level of 0.01 in the case of the following datasets: flags, yeast, scene, foodtruck, Water quality. In the case of emotions, birds and cal500 ECC was significantly worse with a level of 0.05 and in the case of *enron* it was significantly worse with a level of 0.1. Similarly, in Table 6 (with ML-KFHE-CC as the control method), ECC performed significantly worse with a level of 0.01 on *foodtruck* and *Water quality* datasets, and significantly worse with a level of 0.05 in the case of *yeast* and *scene* dataset.

Some interesting patterns can be observed from Table 5 and 6. It is clear that the variant ML-KFHE-HOMER was significantly better in the case of almost all the datasets compared to ML-KFHE-CC, which clearly indicates that ML-KFHE-HOMER is the better variant. Considering ML-KFHE-HOMER vs E-HOMER and HOMER-K and HOMER-B in Table 5 it can be seen that ML-KFHE-HOMER was able to significantly improve on many datasets and never got a worse rank (except in the case of *llog*) as shown in Table 3. Therefore, ML-KFHE-HOMER is always better than the component classifiers as well as a better choice than bagging aggregation approach (E-HOMER). For the other methods, ML-KFHE-HOMER was able to significantly improve the label-based macro-averaged F-Scores in almost all the datasets. Interestingly, CC was significantly better than ML-KFHE-HOMER in the case of *llog* with a significance level of 0.1.

Considering ML-KFHE-CC as the control method in Table 6 it can be seen that ECC has performed significantly worse in the case of *foodtruck* and *Water\_quality* with a significance level of 0.01 and significantly worse in the case of *yeast* and *scene* with a significance level of 0.05. Although ECC performed better in some cases than ML-KFHE-CC (table 3), but for those datasets the null hypotheses could not be rejected in any of the significance levels in this experiment. Also, it is clear that ML-KFHE-CC has performed significantly better in the case of almost all the datasets compared to CC (except *llog*). This shows that ML-KFHE-CC is most of the cases a better ensemble technique compared to ECC, but not as good as ML-KFHE-HOMER.

It must be emphasised that the Wilcoxon's signed rank sum test *cannot* be used to perform multiple classifier comparison without introducing Type I error (rejecting the null hypothesis when it cannot be rejected), as it does not control the Family Wise Error Rate (FWER) [55] in the above analysis. Therefore, each pair from this experiment should *only* be interpreted in isolation from any other algorithms. Multiple classifier comparison is done in Section 5.2.1.

Overall, it can be concluded that the ML-KFHE-HOMER improves the label-based macro-averaged F-Scores significantly in almost all the datasets when compares to any of the algorithms. E-HOMER (also introduced in this work) is an effective technique as well, but ML-KFHE-HOMER almost always performs better than E-HOMER, thus demonstrating the effectiveness of the KFHE framework for ensembling compared to a bagged method for ensemblig. ML-KFHE-CC also performed well, but not as good as ML-KFHE-CC. When compared to ECC, ML-KFHE-CC was able to improve performance of the model on several datasets, but the difference between ECC and ML-KFHE-CC was not as large as what it is between E-HOMER and ML-KFHE-HOMER. Also, the training time growth of HOMER is much faster than CC, which makes HOMER scalable. These leads to the conclution that ML-KFHE-HOMER is a much superior variant.

## 6 Conclusions and Future Work

In multi-label literature there are several ensemble methods which are mostly based on bagging or majority voting methods [15], which perform well. As

in multi-class classification, boosting methods generally performs much better than a single classifier model. But boosting or boosting-like methods are rarely explored in the multi-label literature.

This work introduces a multi-label classification method, ML-KFHE, that exploits the sensor fusion properties of the Kalman filter, used in the Kalman Filter-based Heuristic Ensemble (KFHE). Given the nature of the algorithm, effectively, this falls in the middle of boosting and bagging. ML-KFHE views the ensemble classifier model to be trained as a state to be estimated and does so using a Kalman filter that combines multiple noisy measurements, where each measurement is a trained classifier and the noise is its related classification error.

In ML-KFHE, the sensor fusion properties of the Kalman filter is used to aggregate multiple and diversely trained HOMER or CC models. The method ensembles multiple HOMER or CC models trained on weighted samples of a training dataset and using different hyperparameter settings. The KFHE framework combines these models based on the classification error of the HOMER or CC models. Summary of the findings are as follows

- ML-KFHE was able to perform consistently better than its component classifiers.
- The aggregation method of ML-KFHE using the Kalman filter is more effective than existing and common methods of bagging-like combination as in ECC or E-HOMER, therefore showing the effectiveness of the KFHE framework.
- The ML-KFHE-HOMER variant performed better than ECC, RAkEL and the other multi-label ensemble methods evaluated.

ML-KFHE might converge too fast if the several component classifiers in a sequence have high bias but low variance. This can result in ML-KFHE to perform suboptimally. Presently the random hyperparameter selection of the component classifiers introduce the diversity to stop this happening. Also, in the later iterations of ML-KFHE when the uncertainty of the ensemble is reducing (Kalman gain reduces), if a new component classifier model is found (due to a randomly selected good hyperparameter) to be more accurate, due to the lower uncertainty of the ensemble, the new more accurate measurement may not be incorporated into the model due to a lower Kalman gain value. To stop the method converging too fast, process noise or a slowdown mechanism can be introduced, which may improve performance in some cases where the Kalman gain becomes 1 (one of the measurements was perfect). Also, it would also be interesting to formally compare the training and prediction runtime of the different methods.

Also, In the future a *per-label* version of ML-KFHE could be explored, where instead of the combination of multiple labels using one Kalman gain, per-label Kalman gains will be maintained. The present algorithm does not have a time update step, which can also be introduced and studied.

## Data Availability Statement

The datasets generated during and/or analysed during the current study are available in the following repositories: http://mulan.sourceforge.net/datasets-mlc.html, http://www.uco.es/kdis/mllresources/ .

## Declarations

**Conflict of interests:** The authors have no financial or non-financial conflicts of interests to disclose.

## Acknowledgements

This publication has emanated from research conducted with the financial support of Science Foundation Ireland under Grant number [16/RC/3835]. For the purpose of Open Access, the author has applied a CC BY public copyright licence to any Author Accepted Manuscript version arising from this submission. The authors also would like to acknowledge Dr. Derek Greene for valuable inputs to improve the quality of the draft.

## References

- Herrera, F., Charte, F., Rivera, A.J., del Jesús, M.J.: Multilabel Classification - Problem Analysis, Metrics and Techniques. Springer, ??? (2016). https://doi.org/10.1007/978-3-319-41111-8
- [2] Boutell, M.R., Luo, J., Shen, X., Brown, C.M.: Learning multi-label scene classification. Pattern Recognition 37(9), 1757–1771 (2004). https://doi. org/10.1016/j.patcog.2004.03.009
- [3] Tsoumakas, G., Katakis, I.: Multi-label classification: An overview. Int J Data Warehousing and Mining 2007, 1–13 (2007). https://doi.org/10. 4018/jdwm.2007070101
- [4] Read, J., Pfahringer, B., Holmes, G., Frank, E.: Classifier chains for multilabel classification. Machine Learning 85(3), 333–359 (2011). https://doi. org/10.1007/s10994-011-5256-5
- [5] Zhang, M.-L., Zhou, Z.-H.: Multilabel neural networks with applications to functional genomics and text categorization. IEEE Trans. on Knowl. and Data Eng. 18(10), 1338–1351 (2006). https://doi.org/10.1109/ TKDE.2006.162
- [6] Kelleher, J.D., Mac Namee, В., D'arcy, A.: Fundamentals of Machine Learning for Predictive Data Analytics: Worked MIT Algorithms, Examples, and Case Studies.

Press, ??? (2015). https://mitpress.mit.edu/9780262044691/ fundamentals-of-machine-learning-for-predictive-data-analytics/

- [7] Narassiguin, A., Bibimoune, M., Elghazel, H., Aussem, A.: An extensive empirical comparison of ensemble learning methods for binary classification. Pattern Analysis and Applications 19(4), 1093–1128 (2016). https://doi.org/10.1007/s10044-016-0553-z
- [8] Nasierding, G., Kouzani, A.Z., Tsoumakas, G.: A triple-random ensemble classification method for mining multi-label data. In: 2010 IEEE International Conference on Data Mining Workshops, pp. 49–56 (2010). https://doi.org/10.1109/ICDMW.2010.139
- [9] Tenenboim-Chekina, L., Rokach, L., Shapira, B.: Identification of label dependencies for multi-label classification. In: Working Notes of the Second International Workshop on Learning from Multi-Label Data, pp. 53–60 (2010)
- [10] Kocev, D., Vens, C., Struyf, J., Džeroski, S.: Ensembles of Multi-Objective Decision Trees, pp. 624–631. Springer, Berlin, Heidelberg (2007). https: //doi.org/10.1007/978-3-540-74958-5 61
- [11] Read, J., Pfahringer, B., Holmes, G.: Multi-label classification using ensembles of pruned sets. In: 2008 Eighth IEEE International Conference on Data Mining, pp. 995–1000 (2008). https://doi.org/10.1109/ICDM. 2008.7
- [12] Tsoumakas, G., Katakis, I., Vlahavas, I.: Random k-labelsets for multilabel classification. IEEE Transactions on Knowledge and Data Engineering 23(7), 1079–1089 (2011). https://doi.org/10.1109/TKDE. 2010.164
- [13] Rokach, L., Schclar, A., Itach, E.: Ensemble methods for multi-label classification. Expert Systems with Applications 41(16), 7507–7523 (2014). https://doi.org/10.1016/j.eswa.2014.06.015
- [14] Schapire, R.E., Singer, Y.: BoosTexter: A boosting-based system for text categorization. Machine Learning 39(2), 135–168 (2000). https://doi.org/ 10.1023/A:1007649029923
- [15] Moyano, J.M., Gibaja, E.L., Cios, K.J., Ventura, S.: Review of ensembles of multi-label classifiers: Models, experimental study and prospects. Information Fusion 44, 33–45 (2018). https://doi.org/10.1016/j.inffus. 2017.12.001
- [16] Madjarov, G., Kocev, D., Gjorgjevikj, D., Džeroski, S.: An extensive experimental comparison of methods for multi-label learning. Pattern

Recognition **45**(9), 3084–3104 (2012). https://doi.org/10.1016/j.patcog. 2012.03.004

- [17] Pakrashi, A., Greene, D., Mac Namee, B.: Benchmarking multilabel classification algorithms. In: 24th Irish Conference on Artificial Intelligence and Cognitive Science (AICS'16), Dublin, Ireland, 20-21 September 2016 (2016). https://ceur-ws.org/Vol-1751/AICS\_2016\_ paper 33.pdf
- [18] Pakrashi, A., Mac Namee, B.: Kalman filter-based heuristic ensemble (KFHE): A new perspective on multi-class ensemble classification using kalman filters. Information Sciences 485, 456–485 (2019)
- [19] Kalman, R.E.: A new approach to linear filtering and prediction problems. ASME Journal of Basic Engineering (1960). https://doi.org/10.1115/1. 3662552
- [20] Maybeck, P.S. (ed.): Chapter 6 Design and Performance Analysis of Kalman Filters. Mathematics in Science and Engineering, vol. 141, pp. 289–367. Elsevier, ??? (1979). https://doi.org/10.1016/S0076-5392(08) 62171-2
- [21] Tsoumakas, G., Katakis, I., Vlahavas, I.: Effective and efficient multilabel classification in domains with large number of labels. In: Proc. ECML/PKDD 2008 Workshop on Mining Multidimensional Data (MMD'08), vol. 21, pp. 53–59 (2008). https://doi.org/http://www. ecmlpkdd2008.org/files/pdf/workshops/mmd/4.pdf. sn
- [22] Read, J., Pfahringer, B., Holmes, G.: Multi-label classification using ensembles of pruned sets. In: Data Mining, 2008. ICDM'08. Eighth IEEE International Conference On, pp. 995–1000 (2008). https://doi.org/10. 1109/ICDM.2008.74. IEEE
- [23] Rokach, L.: Decision forest: Twenty years of research. Information Fusion 27, 111–125 (2016). https://doi.org/10.1016/j.inffus.2015.06.005
- [24] Blockeel, H., Raedt, L.D., Ramon, J.: Top-down induction of clustering trees. In: Proceedings of the Fifteenth International Conference on Machine Learning. ICML '98, pp. 55–63. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA (1998)
- [25] Nasierding, G., Tsoumakas, G., Kouzani, A.Z.: Clustering based multilabel classification for image annotation and retrieval. In: 2009 IEEE International Conference on Systems, Man and Cybernetics, pp. 4514– 4519 (2009). https://doi.org/10.1109/ICSMC.2009.5346902
- [26] Freund, Y., Schapire, R.E.: A desicion-theoretic generalization of on-line

learning and an application to boosting. In: European Conference on Computational Learning Theory, pp. 23–37 (1995). https://doi.org/10. 1006/jcss.1997.1504. Springer

- [27] Sebastiani, F., Sperduti, A., Valdambrini, N.: An improved boosting algorithm and its application to text categorization. In: Proceedings of the Ninth International Conference on Information and Knowledge Management, pp. 78–85 (2000). https://doi.org/10.1145/354756.354804. Citeseer
- [28] Al-Salemi, B., Noah, S.A.M., Aziz, M.J.A.: RFBoost: An improved multi-label boosting algorithm and its application to text categorisation. Knowledge-Based Systems 103, 104–117 (2016). https://doi.org/10.1016/ j.knosys.2016.03.029
- [29] Faragher, R., et al.: Understanding the basis of the kalman filter via a simple and intuitive derivation. IEEE Signal processing magazine 29(5), 128–132 (2012). https://doi.org/10.1109/MSP.2012.2203621
- [30] Dietterich, T.G.: Ensemble methods in machine learning. In: Multiple Classifier Systems, pp. 1–15. Springer, Berlin, Heidelberg (2000). https: //doi.org/10.1007/3-540-45014-9\_1
- [31] Bishop, G., Welch, G., et al.: An introduction to the kalman filter. Proc of SIGGRAPH, Course 8(27599-23175), 41 (2001)
- [32] Hastie, T., Rosset, S., Zhu, J., Zou, H.: Multi-class adaboost. Statistics and its Interface 2(3), 349–360 (2009). https://doi.org/10.4310/SII.2009. v2.n3.a8
- [33] Pakrashi, A., Mac Namee, B.: KalmanTune: A Kalman filter based tuning method to make boosted ensembles robust to class-label noise. IEEE Access 8, 145887–145897 (2020). https://doi.org/10.1109/ACCESS.2020. 3013908
- [34] Yu, K., Wang, L., Yu, Y.: Ordering-based kalman filter selective ensemble for classification. IEEE Access 8, 9715–9727 (2020). https://doi.org/10. 1109/ACCESS.2020.2964849
- [35] Zhang, M.-L., Zhou, Z.-H.: A review on multi-label learning algorithms. IEEE transactions on knowledge and data engineering 26(8), 1819–1837 (2014). https://doi.org/10.1109/TKDE.2013.39
- [36] Charte, F., Rivera, A., del Jesus, M.J., Herrera, F.: Concurrence among imbalanced labels and its influence on multilabel resampling algorithms. In: Hybrid Artificial Intelligence Systems, pp. 110–121. Springer, Cham (2014). https://doi.org/10.1007/978-3-319-07617-1\_10

- [37] Goncalves, E.C., Plastino, A., Freitas, A.A.: A genetic algorithm for optimizing the label ordering in multi-label classifier chains. In: 2013 IEEE 25th International Conference on Tools with Artificial Intelligence, pp. 469–476 (2013). https://doi.org/10.1109/ICTAI.2013.76. IEEE
- [38] Elisseeff, A., Weston, J.: A kernel method for multi-labelled classification. In: Dietterich, T., Becker, S., Ghahramani, Z. (eds.) Advances in Neural Information Processing Systems, vol. 14. MIT Press, ??? (2001). https://proceedings.neurips.cc/paper\_files/paper/2001/file/ 39dcaf7a053dc372fbc391d4e6b5d693-Paper.pdf
- [39] Trohidis, K., Tsoumakas, G., Kalliris, G., Vlahavas, I.P.: Multi-label classification of music into emotions. In: ISMIR, vol. 8, pp. 325–330 (2008). https://doi.org/10.1186/1687-4722-2011-426793
- [40] Pestian, J.P., Brew, C., Matykiewicz, P., Hovermale, D.J., Johnson, N., Cohen, K.B., Duch, W.: A shared task involving multi-label classification of clinical free text. In: Proceedings of the Workshop on BioNLP 2007: Biological, Translational, and Clinical Language Processing. BioNLP '07, pp. 97–104. Association for Computational Linguistics, USA (2007). https://aclanthology.org/W07-1013/
- [41] Read, J.: Scalable multi-label classification. PhD thesis, University of Waikato (2010). https://hdl.handle.net/10289/4645
- [42] Briggs, F., Huang, Y., Raich, R., Eftaxias, K., Lei, Z., Cukierski, W., Hadley, S.F., Hadley, A., Betts, M., Fern, X.Z., Irvine, J., Neal, L., Thomas, A., Fodor, G., Tsoumakas, G., Ng, H.W., Nguyen, T.N.T., Huttunen, H., Ruusuvuori, P., Manninen, T., Diment, A., Virtanen, T., Marzat, J., Defretin, J., Callender, D., Hurlburt, C., Larrey, K., Milakov, M.: The 9th annual mlsp competition: New methods for acoustic classification of multiple simultaneous bird species in a noisy environment. In: 2013 IEEE International Workshop on Machine Learning for Signal Processing (MLSP), pp. 1–8 (2013). https://doi.org/10.1109/MLSP.2013. 6661934
- [43] Liu, S.M., Chen, J.-H.: An empirical study of empty prediction of multilabel classification. Expert Systems with Applications 42(13), 5567–5579 (2015). https://doi.org/10.1016/j.eswa.2015.01.024
- [44] Diplaris, S., Tsoumakas, G., Mitkas, P.A., Vlahavas, I.: Protein classification with multiple algorithms. In: Bozanis, P., Houstis, E.N. (eds.) Advances in Informatics, pp. 448–456. Springer, Berlin, Heidelberg (2005). https://doi.org/10.1007/11573036\_42
- [45] Turnbull, D., Barrington, L., Torres, D., Lanckriet, G.: Semantic annotation and retrieval of music and sound effects. IEEE Transactions

on Audio, Speech, and Language Processing **16**(2), 467–476 (2008). https://doi.org/10.1109/TASL.2007.913750

- [46] Rivolli, A., Parker, L.C., de Carvalho, A.C.: Food truck recommendation using multi-label classification. In: Progress in Artificial Intelligence: 18th EPIA Conference on Artificial Intelligence, EPIA 2017, Porto, Portugal, September 5-8, 2017, Proceedings 18, pp. 585–596 (2017). https://doi. org/10.1007/978-3-319-65340-2 48. Springer
- [47] Blockeel, H., Džeroski, S., Grbović, J.: Simultaneous prediction of multiple chemical parameters of river water quality with tilde. In: Principles of Data Mining and Knowledge Discovery: Third European Conference, PKDD'99, Prague, Czech Republic, September 15-18, 1999. Proceedings 3, pp. 32–40 (1999). https://doi.org/10.1007/978-3-540-48247-5\_4. Springer
- [48] Xu, J., Liu, J., Yin, J., Sun, C.: A multi-label feature extraction algorithm via maximizing feature variance and feature-label dependence simultaneously. Knowledge-Based Systems 98, 172–184 (2016). https: //doi.org/10.1016/j.knosys.2016.01.032
- [49] Spyromitros, E., Tsoumakas, G., Vlahavas, I.: An empirical study of lazy multilabel classification algorithms. In: Proc. 5th Hellenic Conference on Artificial Intelligence (SETN 2008) (2008). https://doi.org/10.1007/ 978-3-540-87881-0\_40
- [50] Zhang, M.L., Zhou, Z.H.: ML-kNN: A lazy learning approach to multilabel learning. Pattern Recognition 40, 2038–2048 (2007). https://doi. org/10.1016/j.patcog.2006.12.019
- [51] Cheng, W., Hullermeier, E.: Combining instance-based learning and logistic regression for multilabel classification. Machine Learning 76(2-3), 211–225 (2009). https://doi.org/10.1007/s10994-009-5127-5
- [52] Sechidis, K., Tsoumakas, G., Vlahavas, I.: On the stratification of multilabel data. In: Machine Learning and Knowledge Discovery in Databases, pp. 145–158. Springer, Berlin, Heidelberg (2011). https://doi.org/10. 1007/978-3-642-23808-6 10
- [53] Tsoumakas, G., Spyromitros-Xioufis, E., Vilcek, J., Vlahavas, I.: Mulan: A java library for multi-label learning. Journal of Machine Learning Research 12, 2411–2414 (2011)
- [54] Rivolli, A., de Carvalho, A.C.: The utiml package: Multi-label classification in r. The R Journal 10(2), 24–37 (2018). https://doi.org/ 10.32614/RJ-2018-041

[55] García, S., Fernández, A., Luengo, J., Herrera, F.: Advanced nonparametric tests for multiple comparisons in the design of experiments in computational intelligence and data mining: Experimental analysis of power. Information Sciences 180(10), 2044–2064 (2010). https://doi.org/ 10.1016/j.ins.2009.12.010