

Unlabeled Consensus Modeler: Leveraging Voting Ensemble's Consensus on Unlabeled Data

Dong Cao

Computer Science Department at Earlham College

Richmond, Indiana, USA

dcaohuu18@earlham.edu

ABSTRACT

Voting is an important Ensemble Learning technique. However, there has not been much discussion about leveraging the base classifiers' consensus on unlabeled data to better inform the final prediction. My proposed method identifies the data points where the ensemble reaches consensus and where conflict arises in the unlabeled space. A meta weighted KNN model is trained upon this half-labeled set with the labels of the consensus and the conflict points marked as "Unknown," which is treated as a new, additional class. The predictions of the meta model are expected to better inform the decision of the ensemble in the case of conflict. This research project aims to implement my proposed method and evaluate it on a range of benchmark datasets.

KEYWORDS

Machine Learning, Voting Ensemble, Semi-supervised Learning

1 INTRODUCTION

Ensemble learning has received great attention from the Machine Learning community, as the aggregated output of multiple learners is often better than that of any single one of them. For classification problems specifically, several methods have been proposed to combine multiple classifiers' predictions: Voting, Bagging, Boosting, Stacking, etc. Among these methods, Voting is very important not only because it is a simple, intuitive, and effective method in and of itself, but also because it plays the role of determining the collective prediction in other ensemble frameworks, like Bagging and Boosting. However, there has not been much discussion about leveraging the consensus of the base classifiers on unlabeled data in order to better inform the final prediction. My proposed method identifies the data points where the ensemble reaches consensus and where conflict arises in the unlabeled space. A meta weighted KNN model is trained upon this half-labeled set with the labels of the consensus and the conflict points marked as "Unknown," which is treated as a new, additional class. The predictions of the meta model are expected to better inform the decision of the ensemble in the case of conflict. I named this method Unlabeled Consensus Modeler (UCM). The motivation is that the points agreed upon by the base classifiers represent patterns that we can confidently extract from the training set. However, the training set can also contain noise that leads the base classifiers astray and makes them disagree with each other. In the unlabeled space, the agreed patterns may be clearer around a point of dispute and we can use this information to

strengthen the prediction for that point. This research project is an attempt to implement UCM and compare it with Simple Majority Voting in terms of their performance on a variety of benchmark datasets under some conditions.

In the next section, I will review the related and background knowledge in two domains: Voting Ensemble and Semi-supervised Learning, as well as the differences between UCM and the work introduced. The third section delves deeper into the theoretical motivation and the specific problem that my method targets. It also presents a formal design of UCM and how its components are implemented. Finally, the evaluation framework and results are discussed in the fourth section.

2 RELATED WORK

2.1 Voting Methods

Majority Voting: Majority Voting or Plural Voting is probably the most well-known voting system due to its straightforwardness. A class is assigned to a sample if it receives a majority of votes from the base classifiers. However, there are different definitions of "majority." Depending on the situation, it can mean unanimity, simple majority (i.e. more than 50% of the base classifiers agree on a label). Yet, the most common approach has been that whichever class receiving the most votes "wins" and becomes the ultimate prediction. Simple Majority Voting makes a lot of assumptions about the relative accuracy of the classifiers and each classifiers' performance with respect to a particular class [10]. In reality, most of these assumptions are not accurate. However, the simplicity and efficiency of this method still makes it one of the favorite options of Machine Learning practitioners.

Weighted Voting: Another popular voting scheme is Weighted Voting. Weighted Voting drops one of the assumptions made by Majority Voting [10]. Instead of considering the predictions of the base classifiers equally likely to be accurate, Weighted Voting attaches different weights to the classifiers' predictions based on their performance in the training set. One of the most well-known examples of this approach is the voting system of AdaBoost, which trains a number of weak learners, weights them differently based on their error rates, and aggregates their predictions by taking into account these weights [5]. Other techniques make use of fuzzy sets [3], particle swarm optimization [8], genetic algorithm [12], or instance-wise weight assignment based on a classifier's relative performance with respect to others' [4].

Support Function: With a support function, instead of being confined to a single output class, a base classifier can provide their predictions in terms of the likelihood of a sample belonging to each

of the available classes. The term “Support Function” is used by Woźniak et al. in their survey of classifiers combination [16]. One of the widely used types of likelihood is the *a posteriori* probability. Kittler et al. have proposed a variety of rules by which a class’s *a posteriori* probabilities from different base classifiers can be combined [9]. Later work has focused on comparing the effectiveness of these rules under various conditions [1]. The value of the support function can also be the rankings of the classes. In this case, the method is known as Borda Count, which outperforms Majority Voting in some experiments [11].

While the aforementioned methods tackle the problem differently, they all agree that it is critical for the ensemble to be diverse. This concept of diversity can be interpreted and measured in a variety of ways. Usually, it is important that the classifiers don’t make the same mistake together. Some diversification strategies are using different underlying algorithms for the base classifiers or different feature sets [9], and bootstrapping.

2.2 Semi-supervised Learning

Because UCM utilizes information in both the labeled and unlabeled spaces, it can be linked to semi-supervised learning. However, since semi-supervised learning is a broad field, I will only focus on the areas that are relevant to making use of the ensemble’s consensus on unlabeled data. Before I go into the details of each area, let us quickly touch upon the rudiments of semi-supervised learning. The big problem that semi-supervised learning tries to solve is that labeled data for training is often insufficient and difficult to acquire while unlabeled data is abundant. Semi-supervised learning aims to fully exploit the few labeled samples available to extract patterns from the pool of ample unlabeled data [14].

Inductive versus Transductive: The landscape of semi-supervised learning methods comprises of two major approaches: Inductive and Transductive. The goal of an inductive framework is to build a mechanism that can independently predict unlabeled samples one by one. This goal is shared with most of the supervised algorithms but the training process of an inductive algorithm takes in both labeled and unlabeled data. On the other hand, a transductive method seeks to optimize the predictions for each space of data. This space contains samples that are either labeled or unlabeled and a transductive algorithm attempts to use the distribution of the entire space to provide a set of predictions for all data points. In other words, the input for a transductive algorithm is the whole data space, not a single data point [14]. In this aspect, UCM is similar to the transductive approach when the meta model needs to be fed the complete unlabeled set. However, while most transductive algorithms use graph theory to model the similarity among the data points [14], UCM looks to draw a connection between the patterns in the training set and those in the unlabeled set via inspecting the consensus of the base classifiers.

Tri-Training: Tri-Training is an inductive method that uses three classifiers, all of which are trained upon the same complete dataset. When it comes to leveraging unlabeled data for “refinement,” a classifier is given a sample to train with the label agreed upon by

the other two [17]. A variation of Tri-Training is Multi-Train when more than three classifiers are used and a sample is accepted for the refinement of one classifier if a majority of the rest of the classifiers return the same label [6]. It is not difficult to point out the similarities between my idea and that of Tri-Training, when my meta KNN model learns from the data labeled based on the base classifiers’ consensus. However, there are fundamental differences between UCM and Tri-Training:

- As mentioned, Tri-Training falls under the inductive approach while UCM is generally transductive.
- My meta KNN model does not learn from the labeled data in the training set.
- There is no co-training. In other words, my meta KNN model does not affect the base classifiers in any way. Hence, there is no refinement of the base classifiers using unlabeled data.
- My meta KNN model also takes into account the uncertainty of the nearby conflict points.

3 UNLABELED CONSENSUS MODELER

3.1 Theoretical Motivation

UCM is expected to take advantage of the collective decisions of a voting system and summarize how strong these collective patterns are in the unlabeled set to help with the classification of the conflict points. In Figure 1, the graph on the left presents the training set, which has two classes “A” and “B.” On the right is the graph of the unlabeled set, which is pseudo-labeled based on the majority decisions of two classifiers cl_1 and cl_2 . The circled question mark indicates a conflict point where cl_1 and cl_2 disagree. It is also shown in the training set for the sake of convenient comparison.

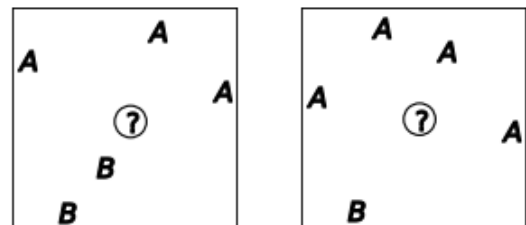


Figure 1: Training (left) and unlabeled (right) sets

In the unlabeled set, it is apparent that the conflict point is more likely to have the label “A.” However, the additional B in the training set causes confusion and dissent between cl_1 and cl_2 . UCM settles this dispute by adding another voice based on the consensus in the unlabeled space. The feasibility of UCM rest on two assumptions. First, the pseudo-labels inferred from consensus are reliable. This assumption can be satisfied with a diverse ensemble. If base classifiers with different learning “lenses” all agree on the label of a point then this label is credible. The second assumption of UCM is that the distribution of the unlabeled set is more trustworthy and contains less noise than the training set, especially around the point of dispute.

It is also possible that in the unlabeled space, there are other conflict

points around the point in question. Figure 2 signifies the other conflict points with unringed question marks. These points are assigned the class "Unknown."

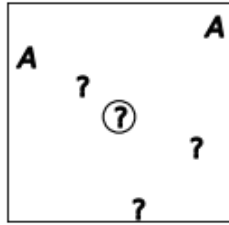


Figure 2: Accounting for uncertainty in the unlabeled set

In Figure 2, although there are two A's near the point in question, UCM also takes notice of the three unknown samples around it and is less positive that the point in question also has the label "A." The analysis of the theoretical motivation of UCM makes it clearer that the distinctions between UCM and Tri-Training reflect the different objectives that the two methods are pursuing. While Tri-Training, a representative of semi-supervised learning, tackles the lack of labeled data, UCM aims at detecting fake patterns that exist in the training set but not in the unlabeled set, thereby reducing the chance of overfitting. Even though the approach of UCM is semi-supervised, it is intended to serve supervised frameworks. There is no need for refinement using unlabeled data since the training data should be sufficient for the base classifiers to perform decently on their own and UCM will only play the role of assisting them with making the final prediction where discord occurs. Another reason for no retraining of the base classifiers is that many semi-supervised techniques suffer from degradation due to their biased conjecture about the unlabeled data [18]. By keeping the opinions of the base classifiers intact and only putting another voice on top of their opinions when needed, UCM is anticipated to be less prone to the issue of degradation.

3.2 Design and Implementation

Figure 3 is the architectural diagram of UCM. I will now dissect each of its components, most of which are implemented with Pandas [15], Numpy [7], and Scikit-learn [13].

Data Cleaning: Constant columns, identity columns, and duplicated rows are dropped. Missing values are consistently signified across all datasets.

Basic preprocessing: Missing values are imputed using the median for numerical variables and the mode for categorical features. If a categorical feature has already been ordinally encoded (for example, in the case of the lymphography dataset), it is kept as such and is treated as a numeric feature. Otherwise, one-hot encoding is applied to satisfy many algorithms' requirement that categorical data be represented as numeric values. Although it may make sense for some categorical features to be transformed according to an ordinal scale, the vast number of evaluation datasets and the fact that they

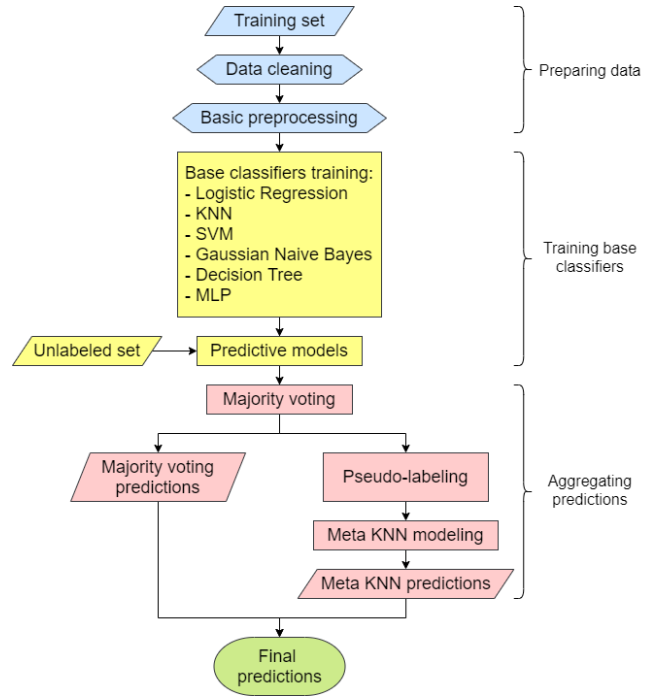


Figure 3: UCM framework

spread across a range of specialized domains make it difficult to determine the ordinality of each categorical feature. Moreover, that all the categorical features are one-hot encoded also makes the findings more reproducible. Most importantly, the objective of this research is not to achieve the best performance on the benchmark datasets. Rather, it is to carry out a comparative experiment of two frameworks and the choice of the preprocessing method does not interfere with this objective significantly. In addition to the categorical encoding, all numerical features are standardized and clamped to the same scale.

Base classifiers training: I follow the strategy of using different underlying algorithms to diversify the ensemble. There are six base classifiers, each of which corresponds to one algorithm in the diagram. The classifiers are constructed from Scikit-learn's standard implementations of the listed algorithms. The `random_state` argument, if available, is set to 1. The `max_iter` of Logistic Regression is increased to 1000 to ensure convergence. The same argument is raised to 5000 for the Multilayer Perceptron (MLP) algorithm. For MLP, there is only one hidden layer and the number of hidden nodes is $\lfloor \frac{n+1}{2} \rfloor$ where n is the number of features, i.e. the number of input nodes. This is based on the suggestion that the number of hidden neurons should be "somewhere between the input layer size and the output layer size." [2] All other hyperparameters are set to the package's default values, including the ReLU activation function. For the KNN algorithm, the `weights` hyperparameter is set to 'distance' instead of 'uniform' while the `n_neighbors` is fixed at 3. This is because we want the meta KNN model to have the same configuration and some validation sets may have no more

than 5 samples. Apart from the above exceptions, all the algorithms' default hyperparameter values are retained.

Majority voting: After the unlabeled points are predicted by the base classifiers, their predictions go through majority voting. The predictions of the majority voting system are in terms of probabilities. For instance, if a sample is classified as "A" by five out of six learners and as "B" by only one learner then the prediction will be $\frac{5}{6}$ "A" and $\frac{1}{6}$ "B." UCM can work with any voting schemes whose output can be interpreted probabilistically and if there is a way to determine consensus. Thus, it is well-suited with a support function. Nevertheless, a more complicated voting technique is unnecessary since UCM is not a voting system on its own but is built upon an existing voting framework. Simple Majority Voting, therefore, is good enough for assessing UCM and its contribution, if any, to the improvement of the voting ensemble.

Pseudo-labeling: A certainty threshold needs to be set. For example, if a data point is agreed upon by at least five out of six (or approximately 83%) classifiers then it is labeled as the majority's decision. Otherwise, it is indicated as "Unknown." For this project, we use a certainty threshold of 51%.

Meta KNN modeling: A distance-weighted KNN is then applied to the pseudo-labeled set to predict each of the unknown points. For each of these points, the meta model also considers the other unknown samples around it. The model's output is the probabilities of the point belonging to one of the original classes or the class "Unknown," which consolidates the amount of uncertainty into the predictions and serves as a regulating factor. This is why it is crucial for the predictions to be probabilistic.

KNN is chosen to be the algorithm of the meta model because it is an intuitive way of thinking about the dissimilarity in distribution between the training set and the unlabeled set. Other algorithms that make a strong use of the data distribution and that can produce probabilistic predictions like SVM may also be good candidates. However, for each unknown point to be classified, it needs to be removed from the pseudo-labeled set before the learner is fit. KNN, as a lazy learning algorithm, nicely meets this "leave-one-out" requirement, although the relaxation of this requirement may be acceptable for some eager learning algorithms.

The meta KNN model is built based on Scikit-learn's implementation of the algorithm and is configured in a similar way to the base KNN classifier. Whereas this is not a hard requirement, it is to make sure that any performance discrepancy between Majority Voting and UCM, if any, can be more confidently attributed to the additional information acquired from the unlabeled space, rather than the meta KNN model's configuration happening to be more (or less) suited to the benchmark datasets. This allows a more rigorous examination of UCM's core assumption regarding the potentially useful information about the unlabeled set's distribution. However, this decision also brings up a risk of the meta KNN model biasedly favoring the base KNN classifier.

Producing the final predictions: For each sample and class, the two probabilities from majority voting and meta KNN modeling are

added and whichever class receives the highest probability score becomes the ultimate label for that sample.

4 EXPERIMENT AND EVALUATION

4.1 Evaluation

The performance of UCM is compared with that of mere Majority Voting to see if the additional technique of modeling the consensus brings any benefit. The accuracy rate and weighted F_1 score are the metrics due to their popularity and applicability and are estimated using 10-fold cross-validation to reduce bias, especially with small datasets. The experiment involves 23 public benchmark datasets from the UCI repository. These datasets form a subset of the datasets that Kuncheva et al. and Dogan et al. employ in their studies [10][4]. Information about the datasets can be found in Table 1.

Table 1: Basic characteristics of the evaluation datasets

ID	Dataset	Attributes	Instances	Classes
1	abalone	8	4177	29
2	anneal	38	798	6
3	arrhythmia	279	452	16
4	audiology	69	226	24
5	breast-cancer	9	286	2
6	breast-cancer-w	9	699	2
7	car	6	1728	4
8	crx	15	690	2
9	dermatology	34	366	6
10	ecoli	7	336	4
11	glass	10	214	7
12	ionosphere	34	351	2
13	iris	5	150	3
14	kr-vs-kp	36	3196	2
15	labor-neg	16	57	2
16	letter	16	20000	26
17	liver disorders	6	345	2
18	lymphography	20	148	4
19	nursery	8	12960	5
20	page-blocks	10	5473	5
21	segment	21	2310	7
22	sonar	208	60	2
23	spambase	57	4601	2

As described earlier, there is a possibility of the meta KNN model being "partial" to the base classifier of the same algorithm and mainly acting as an additional base KNN classifier. To verify this hypothesis, the performance of a Weighted Voting ensemble with the base KNN classifier's weight doubled is also measured. This Weighted Voting ensemble is simply a Majority Voting ensemble with seven base classifiers, two of which use the KNN algorithm.

Table 2: Majority Voting's, UCM's, and Weighted Voting's performance (%)

ID	Dataset	MV Ac	UCM Ac	WV Ac	MV F_1	UCM F_1	WV F_1
1	abalone	26.56	26.41	26.03	23.66	23.41	23.48
2	anneal	92.6	91.35	92.22	92.64	91.11	92.27
3	arrhythmia	66.74	65.41	66.08	58.04	55.62	57.03
4	audiology	82.23	78.72	82.25	78.47	73.88	78.7
5	breast-cancer	71.55	70.48	70.18	67.95	67.31	67.96
6	breast-cancer-w	96.85	96.85	97	96.87	96.87	97.01
7	car	88.72	89.35	87.5	89.18	89.59	87.67
8	crx	84.18	84.18	83.46	82.74	83.05	82.76
9	dermatology	97.24	96.97	97.24	97.23	96.94	97.21
10	ecoli	86.54	86.25	86.55	85.62	85.28	85.66
11	glass	66.28	65.76	70.48	61.5	60.56	66.45
12	ionosphere	92.86	91.14	90.29	92.66	90.8	89.88
13	iris	95.33	95.33	95.33	95.29	95.29	95.29
14	kr-vs-kp	96.87	96.4	95.93	96.86	96.37	95.9
15	labor-neg	94.33	94.33	96.33	94.19	94.19	96.19
16	letter	92.36	91.39	94.76	92.41	91.43	94.78
17	liver disorders	72.94	72.06	73.24	72.81	71.23	72.69
18	lymphography	80.95	80.29	81.67	79.71	78.79	80.6
19	nursery	90.57	90.9	89.29	89.76	90.1	88.64
20	page-blocks	96.05	96.18	96.03	95.57	95.76	95.71
21	segment	95.93	96.32	96.19	95.87	96.28	96.14
22	sonar	66.71	66.74	67.64	65.39	65.87	66.96
23	spambase	93.91	93.98	93.83	93.88	93.96	93.81
	Average	83.84	83.34	83.89	82.53	81.9	82.73

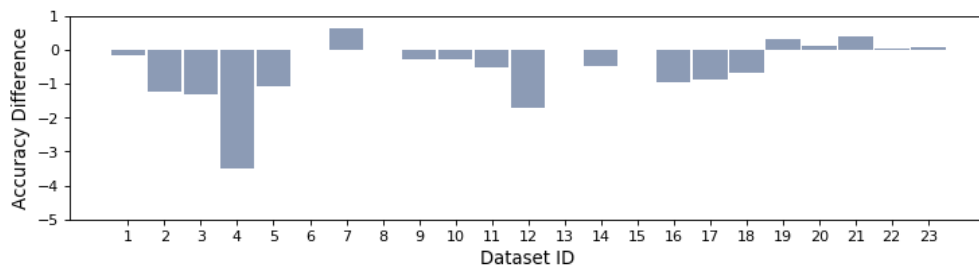


Figure 4: Differences between UCM's and Majority Voting's accuracy: UCM Ac - MV Ac

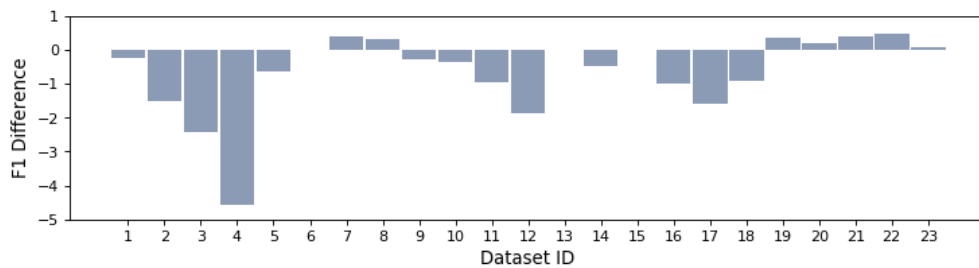


Figure 5: Differences between UCM's and Majority Voting's F_1 scores: UCM F_1 - MV F_1

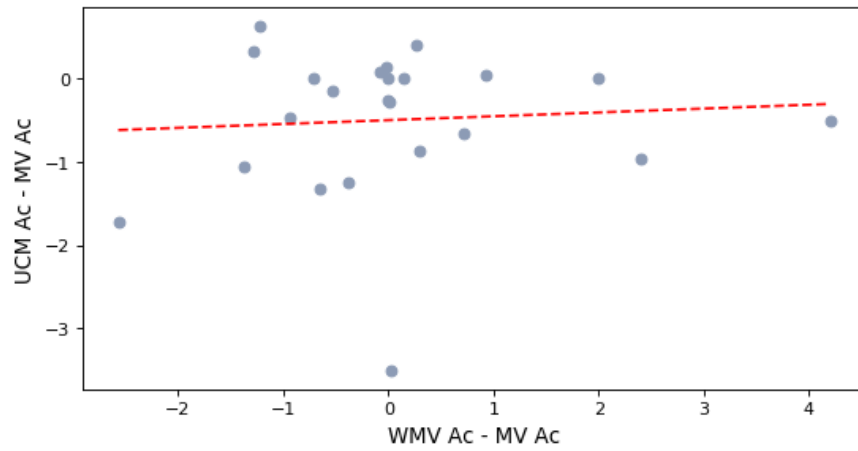


Figure 6: Weighted Voting's and UCM's relative accuracy against Majority Voting

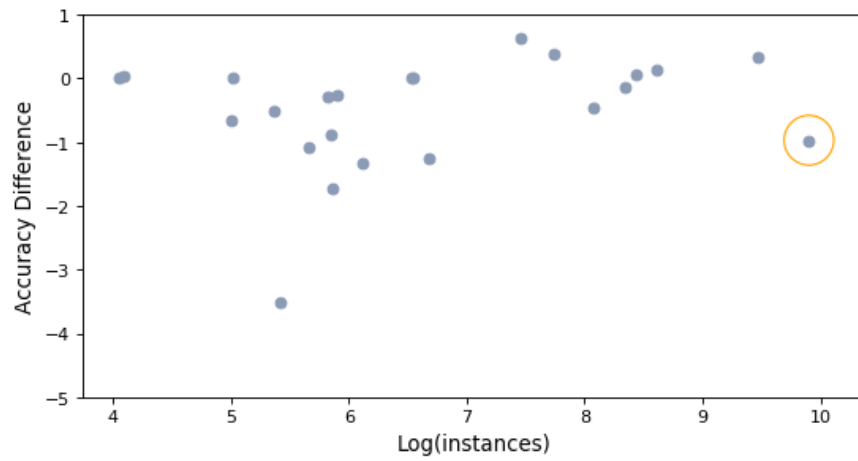


Figure 7: The number of instances and the accuracy difference between Majority Voting and UCM

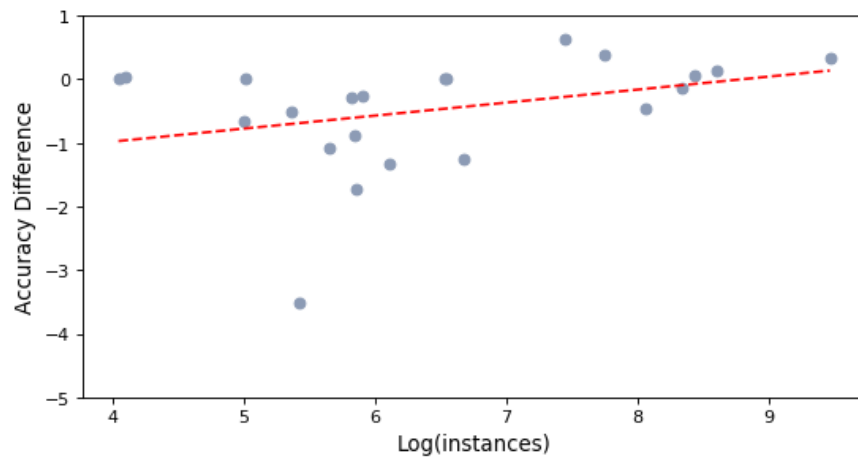


Figure 8: The number of instances and the accuracy difference with the outlier removed

4.2 Results

The performance of Majority Voting, UCM, and Weighted Voting on the benchmark datasets is recorded in Table 2. For both metrics, Majority Voting surpasses UCM on 12 out of 23 datasets and outperforms by a larger amount on average. Figures 4 and 5 are visualizations of the discrepancy between the performance of the two methods.

Weighted Voting is the ultimate winner with an average accuracy rate of 83.89%. This outcome challenges the hypothesis that the meta KNN model just assigns more weight to the base KNN classifier. Otherwise, we would have observed less variation between the scores of Weighted Voting and UCM. The relative performance of the two methods with respect to Majority Voting leads us to the same conclusion.

In Figure 6, the accuracy difference between UCM and Majority Voting is plotted against the accuracy difference between Weighted Voting and Majority Voting. If there were a strong connection between the meta KNN model and the base KNN classifier, there should be a noticeable correlation. However, the line of best fit is nearly horizontal and the R^2 score is very low at 0.0052.

While it is hard to reject any relation between the meta model and the base classifier of the same algorithm, our observation about their little relevance may increase the explanatory power of the unlabeled set's distribution with regard to the performance of UCM. An interesting phenomenon is that 5 out of the 6 datasets on which UCM beats Majority Voting have more than 1000 instances and on only 2 datasets with more than 1000 records does Majority Voting perform better. Figure 7 plots the accuracy difference between UCM and Majority Voting against the number of instances in a logarithmic scale due to the wide range of the volumes of data. There seems to be an increasing relationship between the number of records and the performance of UCM against Majority Voting.

If we revisit the two assumptions made by UCM as discussed in section 3.1, more data may better satisfy both of them. That is it can enhance the diversity of the ensemble and the authenticity of the unlabeled set's distribution. This is because we are applying cross-validation and more samples in the complete dataset mean more samples in the "unlabeled" validation set. However, because the satisfaction of the first assumption also benefits Majority Voting and we are concentrating on the performance difference between Majority Voting and UCM, the pattern displayed in Figure 7 should be better explained by the second assumption. This verdict justifies dropping the outlier on the far right of Figure 7. This outlier corresponds to dataset 16 - the letter recognition dataset. It gives information about images of the 26 capital letters in the English alphabet written in different fonts. The nature of this dataset and its data volume indicate that its training set contains little noise while the second assumption asserts that the unlabeled set should be less noisy than the training set in general. When this outlier is removed, a sharper line of best fit is shown in Figure 8 with a moderate R^2 score of 0.1126. This result suggests the validity of the second assumption about the utility of the unlabeled set's distribution, which can be enriched with more data.

5 CONCLUSION AND FUTURE WORK

A novel method, called UCM, has been proposed. Under this method, samples in the unlabeled set are pseudo-labeled based on a Majority Voting ensemble's consensus. If consensus is not reached, the sample is categorized as "Unknown," which serves as another class. A meta KNN model is then fit on the pseudo-labeled set and its predictions are combined with those of the voting ensemble to generate the final predictions for the unknown points.

UCM was applied on 23 benchmark datasets from the UCI repository. Its performance was compared with Majority Voting. The results show that Majority Voting achieved a higher accuracy and F_1 score on more than half of the datasets. Its average score is also higher than that of UCM in both metrics. Nevertheless, this outcome does not necessarily reject the validity of UCM's theoretical motivation concerning the useful information of the unlabeled set's distribution. A closer look at the results reveals a positive trend between the number of samples and the performance discrepancy of the two methods. This may imply that UCM is more suitable when the unlabeled set has a more representative distribution than the training set, which can be determined by the number of instances and the nature of the data.

Future work on UCM includes but is not limited to:

- (1) Testing UCM on medium-sized datasets (between 1000 and 5000 records) and increasing the size of the validation sets to meet UCM's second assumption.
- (2) Weighting the probabilities addition of the meta model by quantifying and estimating the noise in the unlabeled set as well as the diversity of the ensemble. We may also experiment with another algebraic operation other than addition.
- (3) Examining the potency of UCM as a tie breaker and evaluating it against the baseline of random guessing.
- (4) Tuning the certainty threshold based on the number of classes. For example, a class having three out of six votes may be enough for problems with more than 10 classes. Adjusting this hyperparameter has brought about a 3% increase in accuracy for UCM on some datasets.
- (5) Using a support function where each of the base classifiers produces a probability prediction. Such an output will align more with the prediction provided by the meta model.
- (6) Trying other diversification strategies like bootstrapping. Some of these strategies will allow a higher number of base classifiers.
- (7) Trying other algorithms (e.g. SVM) for the meta model in place of KNN.

ACKNOWLEDGMENTS

I would like to thank Dr. David Barbella and Dr. Charlie Peck for their detailed feedback during the development of this research project.

REFERENCES

- [1] Luis A. Alexandre, Aurélio C. Campilho, and Mohamed Kamel. 2001. On combining classifiers using sum and product rules. *Pattern Recognition Letters* 22, 12 (2001), 1283–1289.
- [2] Adam Blum. 1992. *Neural networks in C++ an object-oriented framework for building connectionist systems*. John Wiley & Sons, Inc.

- [3] Robert Burduk. 2012. Recognition task with feature selection and weighted majority voting based on interval-valued fuzzy sets. In *International Conference on Computational Collective Intelligence*. Springer, 204–209.
- [4] Alican Dogan and Derya Birant. 2019. A weighted majority voting ensemble approach for classification. In *2019 4th International Conference on Computer Science and Engineering (UBMK)*. IEEE, 1–6.
- [5] Yoav Freund and Robert E. Schapire. 1997. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences* 55, 1 (1997), 119–139.
- [6] Shenkai Gu and Yaochu Jin. 2017. Multi-train: A semi-supervised heterogeneous ensemble classifier. *Neurocomputing* 249 (2017), 202–211.
- [7] Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. 2020. Array programming with NumPy. *Nature* 585, 7825 (Sept. 2020), 357–362. <https://doi.org/10.1038/s41586-020-2649-2>
- [8] Asma Kausar, M. Ishtiaq, M. Arfan Jaffar, and Anwar M. Mirza. 2010. Optimization of ensemble based decision using PSO. In *Proceedings of the World Congress on Engineering*, Vol. 1. IAENG, 1–6.
- [9] Josef Kittler, Mohamad Hatef, Robert P.W. Duin, and Jiri Matas. 1998. On combining classifiers. *IEEE transactions on pattern analysis and machine intelligence* 20, 3 (1998), 226–239.
- [10] Ludmila I. Kuncheva and Juan J. Rodríguez. 2014. A weighted voting framework for classifiers ensembles. *Knowledge and Information Systems* 38, 2 (2014), 259–275.
- [11] Florin Leon, Sabina-Adriana Floria, and Costin Bădică. 2017. Evaluating the effect of voting methods on ensemble-based classification. In *2017 IEEE International Conference on INnovations in Intelligent SysTems and Applications (INISTA)*. IEEE, 1–6.
- [12] Yasir Mehmood, Muhammad Ishtiaq, Muhammad Tariq, and M. Arfan Jaffar. 2010. Classifier ensemble optimization for gender classification using genetic algorithm. In *2010 International Conference on Information and Emerging Technologies*. IEEE, 1–5.
- [13] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12 (2011), 2825–2830.
- [14] Jesper E. Van Engelen and Holger H. Hoos. 2020. A survey on semi-supervised learning. *Machine Learning* 109, 2 (2020), 373–440.
- [15] Wes McKinney. 2010. Data Structures for Statistical Computing in Python. In *Proceedings of the 9th Python in Science Conference*, Stéfan van der Walt and Jarrod Millman (Eds.), 56 – 61. <https://doi.org/10.25080/Majora-92bf1922-00a>
- [16] Michał Woźniak, Manuel Grana, and Emilio Corchado. 2014. A survey of multiple classifier systems as hybrid systems. *Information Fusion* 16 (2014), 3–17.
- [17] Zhi-Hua Zhou and Ming Li. 2005. Tri-training: Exploiting unlabeled data using three classifiers. *IEEE Transactions on knowledge and Data Engineering* 17, 11 (2005), 1529–1541.
- [18] Xiaojin Jerry Zhu. 2005. Semi-supervised learning literature survey. (2005).

A SOURCE CODE

The source code of this research project can be accessed via this [hyperlink](#). Please see the enclosed README for instructions on how to use the code. The linked repository also includes metadata about the benchmark datasets and the results reported in this paper.