
Is Combining Classifiers Better than Selecting the Best One?

Sašo Džeroski
Bernard Ženko

SASO.DZEROSKI@IJS.SI
BERNARD.ZENKO@IJS.SI

Department of Intelligent Systems, Jožef Stefan Institute, Jamova 39, 1000 Ljubljana, Slovenia

Abstract

We empirically evaluate several state-of-the-art methods for constructing ensembles of heterogeneous classifiers with stacking and show that they perform (at best) comparably to selecting the best classifier from the ensemble by cross validation. We then propose a new method for stacking, that uses multi-response model trees at the meta-level, and show that it clearly outperforms existing stacking approaches and selecting the best classifier by cross validation.

1. Introduction

An ensemble of classifiers is a set of classifiers whose individual predictions are combined in some way (typically by voting) to classify new examples. One of the most active areas of research in supervised learning has been to study methods for constructing good ensembles of classifiers (Dietterich, 1997). The attraction that this topic exerts on machine learning researchers is based on the premise that ensembles are often much more accurate than the individual classifiers that make them up.

Most of the research on classifier ensembles is concerned with generating ensembles by using a single learning algorithm (Dietterich, 2000), such as decision tree learning or neural network training. Different classifiers are generated by manipulating the training set (as done in boosting or bagging), manipulating the input features, manipulating the output targets or injecting randomness in the learning algorithm. The generated classifiers are then typically combined by voting or weighted voting.

Another approach is to generate classifiers by applying different learning algorithms (with heterogeneous model representations) to a single dataset (see, e.g., (Merz, 1999)). More complicated methods for combining classifiers are typically used in this setting. Stacking (Wolpert, 1992) is often used to learn a combin-

ing method in addition to the ensemble of classifiers. Voting is then used as a baseline method for combining classifiers against which the learned combiners are compared. Typically, much better performance is achieved by stacking as compared to voting.

The work presented in this paper is set in the stacking framework. We argue that selecting the best of the classifiers in an ensemble generated by applying different learning algorithms should be considered as a baseline to which the stacking performance should be compared. Our empirical evaluation of several recent stacking approaches shows that they perform comparably to the best of the individual classifiers as selected by cross validation, but not better. We then propose a new stacking method, based on classification by using model trees, and show that this method does perform better than other combining approaches, as well as better than selecting the best individual classifier.

Section 2 first summarizes the stacking framework, then surveys some recent results and finally introduces our stacking approach based on classification via model trees. The setup for the experimental comparison of several stacking methods, voting and selecting the best classifier is described in Section 3. Section 4 presents and discusses the experimental results and Section 5 concludes.

2. Stacking with model trees

We first give a brief introduction to the stacking framework, introduced by (Wolpert, 1992). We then summarize the results of several recent studies in stacking (Merz, 1999; Ting & Witten, 1999; Todorovski & Džeroski, 2000; Seewald & Fürnkranz, 2001; Todorovski & Džeroski, 2002). Motivated by these, we introduce a stacking approach based on classification via model trees (Frank et al., 1998).

2.1 The stacking framework

Stacking is concerned with combining multiple classifiers generated by using different learning algorithms

L_1, \dots, L_N on a single dataset S , which consists of examples $s_i = (x_i, y_i)$, i.e., pairs of feature vectors (x_i) and their classifications (y_i). In the first phase, a set of base-level classifiers C_1, C_2, \dots, C_N is generated, where $C_i = L_i(S)$. In the second phase, a meta-level classifier is learned that combines the outputs of the base-level classifiers.

To generate a training set for learning the meta-level classifier, a leave-one-out or a cross validation procedure is applied. For leave-one-out, we apply each of the base-level learning algorithms to almost the entire dataset, leaving one example for testing:

$$\forall i = 1, \dots, n : \forall k = 1, \dots, N : C_k^i = L_k(S - s_i).$$

We then use the learned classifiers to generate predictions for s_i : $\hat{y}_i^k = C_k^i(x_i)$. The meta-level dataset consists of examples of the form $((\hat{y}_i^1, \dots, \hat{y}_i^n), y_i)$, where the features are the predictions of the base-level classifiers and the class is the correct class of the example at hand. When performing, say, 10-fold cross validation, instead of leaving out one example at a time, subsets of size one-tenth of the original dataset are left out and the predictions of the learned classifiers obtained on these.

In contrast to stacking, no learning takes place at the meta-level when combining classifiers by a voting scheme (such as plurality, probabilistic or weighted voting). The voting scheme remains the same for all different training sets and sets of learning algorithms (or base-level classifiers). The simplest voting scheme is the plurality vote. According to this voting scheme, each base-level classifier casts a vote for its prediction. The example is classified in the class that collects the most votes.

2.2 Recent advances

The most important issues in stacking are probably the choice of the features and the algorithm for learning at the meta-level. Below we review some recent research on stacking that addresses the above issues.

It is common knowledge that ensembles of diverse base-level classifiers (with weakly correlated predictions) yield good performance. (Merz, 1999) proposes a stacking method called SCANN that uses correspondence analysis to detect correlations between the predictions of base-level classifiers. The original meta-level feature space (the class-value predictions) is transformed to remove the dependencies, and a nearest neighbor method is used as the meta-level classifier on this new feature space.

(Ting & Witten, 1999) use base-level classifiers whose predictions are probability distributions over the set

of class values, rather than single class values. The meta-level attributes are thus the probabilities of each of the class values returned by each of the base-level classifiers. The authors argue that this allows to use not only the predictions, but also the confidence of the base-level classifiers. Multi-response linear regression (MLR) is recommended for meta-level learning, while several learning algorithms are shown not to be suitable for this task.

(Seewald & Fürnkranz, 2001) propose a method for combining classifiers called grading that learns a meta-level classifier for each base-level classifier. The meta-level classifier predicts whether the base-level classifier is to be trusted (i.e., whether its prediction will be correct). The base-level attributes are used also as meta-level attributes, while the meta-level class values are + (correct) and - (incorrect). Only the base-level classifiers that are predicted to be correct are taken and their predictions combined by summing up the probability distributions predicted.

(Todorovski & Džeroski, 2000) introduce a new meta-level learning method for combining classifiers with stacking: meta decision trees (MDTs) have base-level classifiers in the leaves, instead of class-value predictions. Properties of the probability distributions predicted by the base-level classifiers (such as entropy and maximum probability) are used as meta-level attributes, rather than the distributions themselves. These properties reflect the confidence of the base-level classifiers and give rise to very small MDTs, which can (at least in principle) be inspected and interpreted.

(Todorovski & Džeroski, 2002) report that stacking with MDTs clearly outperforms voting and stacking with decision trees, as well as boosting and bagging of decision trees. On the other hand, MDTs perform only slightly better than SCANN and selecting the best classifier with cross validation (SelectBest). (Ženko et al., 2001) report that MDTs perform slightly worse as compared to stacking with MLR. Overall, SCANN, MDTs, stacking with MLR and SelectBest seem to perform at about the same level.

It would seem natural to expect that ensembles of classifiers induced by stacking would perform better than the best individual base-level classifier: otherwise the extra work of learning a meta-level classifier doesn't seem justified. The experimental results mentioned above, however, do not show clear evidence of this. This has motivated us to investigate the performance of state-of-the-art stacking methods in comparison to SelectBest and seek new stacking methods that would be clearly superior to SelectBest.

2.3 Stacking with multi-response model trees

We assume that each base-level classifier predicts a probability distribution over the possible class values. Thus, the prediction of the base-level classifier C when applied to example x is a probability distribution:

$$\mathbf{p}^C(x) = (p^C(c_1|x), p^C(c_2|x), \dots, p^C(c_m|x)),$$

where $\{c_1, c_2, \dots, c_m\}$ is the set of possible class values and $p^C(c_i|x)$ denotes the probability that example x belongs to class c_i as estimated (and predicted) by classifier C . The class c_j with the highest class probability $p^C(c_j|x)$ is predicted by classifier C . The meta-level attributes are thus the probabilities predicted for each possible class by each of the base-level classifiers, i.e., $p^{C_j}(c_i|x)$ for $i = 1, \dots, m$ and $j = 1, \dots, N$.

The experimental evidence mentioned above indicates that although SCANN, MDTs, stacking with MLR and SelectBest seem to perform at about the same level, stacking with MLR has a slight advantage over the other methods. It would thus seem as a suitable starting point in the search for better method for meta-level learning to be used in stacking. Stacking with MLR uses linear regression to perform classification. A natural direction to look into is the use of model trees (which perform piece-wise linear regression) instead of MLR: model trees have namely been shown to perform better than MLR for classification via regression (Frank et al., 1998).

MLR is an adaptation of linear regression. For a classification problem with m class values $\{c_1, c_2, \dots, c_m\}$, m regression problems are formulated: for problem j , a linear equation LR_j is constructed to predict a binary variable which has value one if the class value is c_j and zero otherwise. Given a new example x to classify, $LR_j(x)$ is calculated for all j , and the class k is predicted with maximum $LR_k(x)$.

In our approach, we use model tree induction instead of linear regression and keep everything else the same. Instead of m linear equations LR_j , we induce m model trees MT_j . M5' (Wang & Witten, 1997), a re-implementation of M5 (Quinlan, 1992) included in the data mining suite Weka (Witten & Frank, 1999), is used to induce the trees. Given a new example x to classify, $MT_j(x)$ is calculated for all j , and the class k is predicted with maximum $MT_k(x)$. We call our approach stacking with multi-response model trees, analogously to stacking with MLR.

3. Experimental setup

In the experiments, we investigate the following issues:

- The (relative) performance of existing state-of-the-

Table 1. The datasets used and their properties.

DATASET	EXS	CLS	(D/C)	ATT	MAJ	ENT
AUSTRALIAN	690	2	(8/6)	14	0.56	0.99
BALANCE	625	3	(0/4)	4	0.46	1.32
BREAST-W	699	2	(9/0)	9	0.66	0.92
BRIDGES-TD	102	2	(4/3)	7	0.85	0.61
CAR	1728	4	(6/0)	6	0.70	1.21
CHESS	3196	2	(36/0)	36	0.52	0.99
DIABETES	768	2	(0/8)	8	0.65	0.93
ECHO	131	2	(1/5)	6	0.67	0.91
GERMAN	1000	2	(13/7)	20	0.70	0.88
GLASS	214	6	(0/9)	9	0.36	2.18
HEART	270	2	(6/7)	13	0.56	0.99
HEPATITIS	155	2	(13/6)	19	0.79	0.74
HYPO	3163	2	(18/7)	25	0.95	0.29
IMAGE	2310	7	(0/19)	19	0.14	2.78
IONOSPHERE	351	2	(0/34)	34	0.64	0.94
IRIS	150	3	(0/4)	4	0.33	1.58
SOYA	683	19	(35/0)	35	0.13	3.79
TIC-TAC-TOE	958	2	(9/0)	9	0.65	0.93
VOTE	435	2	(16/0)	16	0.61	0.96
WAVEFORM	5000	3	(0/21)	21	0.34	1.58
WINE	178	3	(0/13)	13	0.40	1.56

art stacking methods, especially in comparison to SelectBest.

- The performance of stacking with multi-response model trees relative to the above methods.
- The influence of the number of base-level classifiers on the (relative) performance of the above methods.

We look into the last topic because the recent studies mentioned above use different numbers of base-level classifiers, ranging from three to eight. The Weka data mining suite (Witten & Frank, 1999) was used for all experiments, within which all the base-level and meta-level learning algorithms used in the experiments have been implemented. Ten-fold cross validation is used to construct the meta-level datasets.

3.1 Datasets

In order to evaluate the performance of the different combining algorithms, we perform experiments on a collection of twenty-one datasets from the *UCI Repository of machine learning databases* (Blake & Merz, 1998). These datasets have been widely used in other comparative studies. The datasets and their properties (number of examples, classes, (discrete/continuous) attributes, probability of the majority class, entropy of the class probability distribution) are listed in Table 1.

3.2 Base-level algorithms

We performed two batches of experiments: one with three and one with seven base-level learners. The set

of three contains the following algorithms:

- J4.8: a Java re-implementation of the decision tree learning algorithm C4.5 (Quinlan, 1993),
- IB k : the k -nearest neighbor algorithm of (Aha et al., 1991), and
- NB: the naive Bayes algorithm of (John & Langley, 1995).

The second set of algorithms contains, in addition to the above three, also the following four algorithms:

- K*: an instance-based algorithm which uses an entropic distance measure (John & Leonard, 1995),
- KDE: a simple kernel density estimation algorithm,
- DT: the decision table majority algorithm of (Kohavi, 1995),
- MLR: the multi-response linear regression algorithm, as used by (Ting & Witten, 1999) and described in Section 2.3.

All algorithms are used with their default parameter settings, with the exceptions described below. IB k in the set of three learners used inverse distance weighting and k was selected with cross validation from the range of 1 to 77. (IB k in the set of seven learners uses the default parameter values, i.e., no weighting and $k = 1$.) The NB algorithm in both sets uses the kernel density estimator rather than assume normal distributions for numeric attributes.

3.3 Meta-level algorithms

At the meta-level, we evaluate the performance of six different schemes for combining classifiers (listed below), each applied with the two different sets of base-level algorithms described above.

- VOTE: The simple plurality vote scheme (see Section 2.1),
- SELB: The SelectBest scheme selects the best of the base-level classifiers by cross validation.
- GRAD: Grading as introduced by (Seewald & Fürnkranz, 2001) and briefly described in Section 2.2.
- SMDT: Stacking with meta decision-trees as introduced by (Todorovski & Džeroski, 2000) and briefly described in Section 2.2.
- SMLR: Stacking with multiple-response regression as used by (Ting & Witten, 1999) and described in Sections 2.2 and 2.3.
- SMM5: Stacking with multiple-response model trees, as proposed by this paper and described in Section 2.3.

3.4 Evaluating and comparing algorithms

In all the experiments presented here, classification errors are estimated using ten-fold stratified cross validation. Cross validation is repeated ten times using different random generator seeds resulting in ten different sets of folds. The same folds (random generator seeds) are used in all experiments. The classification error of a classification algorithm C for a given dataset as estimated by averaging over the ten runs of ten-fold cross validation is denoted $\text{error}(C)$.

For pair-wise comparisons of classification algorithms, we calculate the relative improvement and the paired t -test, as described below. In order to evaluate the accuracy improvement achieved in a given domain by using classifier C_1 as compared to using classifier C_2 , we calculate the relative improvement: $1 - \text{error}(C_1)/\text{error}(C_2)$. In Table 4, we compare the performance of SMM5 to other approaches: C_1 in this table thus refers to ensembles combined with SMM5. The average relative improvement across all domains is calculated using the geometric mean of error reduction in individual domains: $1 - \text{geometric_mean}(\text{error}(C_1)/\text{error}(C_2))$. Note that this may be different from $\text{geometric_mean}(\text{error}(C_2)/\text{error}(C_1)) - 1$.

The classification errors of C_1 and C_2 averaged over the ten runs of 10-fold cross validation are compared for each dataset ($\text{error}(C_1)$ and $\text{error}(C_2)$ refer to these averages). The statistical significance of the difference in performance is tested using the paired t -test (exactly the same folds are used for C_1 and C_2) with significance level of 95%: $+/-$ to the right of a figure in the tables with results means that the classifier C_1 is significantly better/worse than C_2 .

We also study how the improvement of performance of SMM5 over SMLR and SELB is related to the diversity of the base-level classifiers. We use the measure of error correlation (higher diversity means lower error correlation) proposed by (Gama, 1999). For two classifiers C_i and C_j , this measure $\hat{\phi}(C_i, C_j)$ is defined as the conditional probability that both classifiers make the same error, given that one of them makes an error: $p(C_i(x) = C_j(x) | C_i(x) \neq c(x) \vee C_j(x) \neq c(x))$, where $C_i(x)$ and $C_j(x)$ are the predictions of classifiers C_i and C_j for a given example x and $c(x)$ is the true class of x . The error correlation for a set of multiple classifiers \mathcal{C} is defined as the average of the pairwise error correlations:

$$\hat{\phi}(\mathcal{C}) = \frac{1}{|\mathcal{C}|(|\mathcal{C}| - 1)} \sum_{C_i \in \mathcal{C}} \sum_{C_j \neq C_i} \hat{\phi}(C_i, C_j).$$

Table 2. The relative performance of 3-classifier ensembles with different combining methods. The entry in row X and column Y gives the relative improvement of X over Y in % and the number of wins/loses.

	VOTE	SELB	GRAD	SMDT	SMLR	SMM5	TOTAL
VOTE		-21.53 7+/10-	-4.12 6+/5-	-22.45 6+/11-	-27.43 5+/11-	-47.06 2+/10-	26+/47-
SELB	17.72 10+/7-		14.33 11+/3-	-0.76 0+/2-	-4.85 2+/5-	-21.00 1+/9-	24+/26-
GRAD	3.96 5+/6-	-16.72 3+/11-		-17.60 1+/12-	-22.39 2+/14-	-41.24 1+/13-	12+/56-
SMDT	18.34 11+/6-	0.75 2+/0-	14.97 12+/1-		-4.07 4+/5-	-20.10 2+/8-	31+/20-
SMLR	21.53 11+/5-	4.63 5+/2-	18.29 14+/2-	3.91 5+/4-		-15.40 1+/7-	36+/20-
SMM5	32.00 10+/2-	17.36 9+/1-	29.20 13+/1-	16.73 8+/2-	13.35 7+/1-		47+/7-

Table 3. The relative performance of 7-classifier ensembles with different combining methods. The entry in row X and column Y gives the relative improvement of X over Y in % and the number of wins/loses.

	VOTE	SELB	GRAD	SMDT	SMLR	SMM5	TOTAL
VOTE		-19.21 5+/12-	-6.73 2+/7-	-18.04 4+/9-	-24.40 2+/10-	-42.04 0+/10-	13+/48-
SELB	16.10 12+/5-		10.46 11+/4-	0.97 3+/3-	-4.37 5+/7-	-19.17 2+/7-	33+/26-
GRAD	6.30 7+/2-	-11.68 4+/11-		-10.60 5+/7-	-16.56 2+/12-	-33.09 0+/12-	18+/44-
SMDT	15.29 9+/4-	-0.97 3+/3-	9.59 7+/5-		-5.39 5+/6-	-20.33 0+/11-	24+/29-
SMLR	19.62 10+/2-	4.19 7+/5-	14.21 12+/2-	5.11 6+/5-		-14.18 1+/5-	36+/19-
SMM5	29.60 10+/0-	16.08 7+/2-	24.86 12+/0-	16.89 11+/0-	12.42 5+/1-		45+/3-

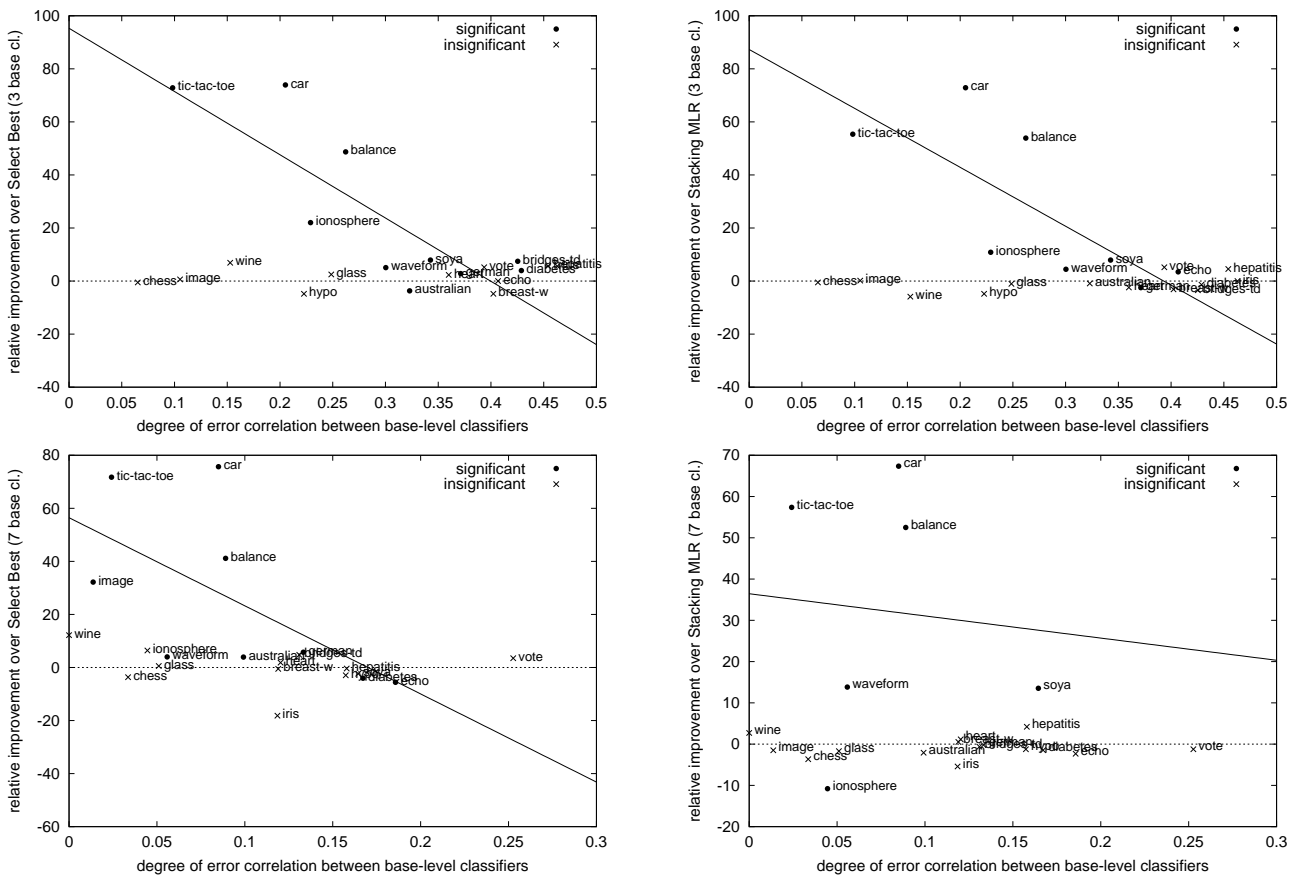


Figure 1. The relation between error correlation and relative improvement of SMM5 over SMLR and SELB.

4. Experimental results

The error rates of the 3-classifier and 7-classifier ensembles induced as described above on the twenty-one dataset and combined with the different combining methods are given in Table 5. However, for the purpose of comparing the performance of different combining methods, Tables 2 and 3 are of much more interest: they give the average relative improvement of X over Y for each pair of combining methods X and Y , as well as the number of significant wins/losses. Table 4 presents a more detailed comparison (per dataset) of SMM5 to the other combining methods. Below we highlight some of our more interesting findings.

4.1 State-of-the-art stacking methods

Inspecting Tables 2 and 3, we find that we can partition the five combining algorithms (we do not consider SMM5 at this stage of the analysis) into three groups. VOTE and GRAD are at the lower end of the performance scale, SELB and SMDT are in the middle, while SMLR performs best. While SMLR clearly outperforms VOTE and GRAD, the advantage over SELB is slim (3 and 2 more wins than losses, about 4% relative improvement) and the advantage over SMDT even slimmer (1 more win than loss in both cases, 4 and 5% of relative improvement).

4.2 Stacking with multi-response model trees

Returning to Tables 2 and 3, this time paying attention to the relative performance of SMM5 to the other combining methods, we find that SMM5 is in a league of its own. It clearly outperforms all the other combining methods, with a wins – loss difference of at least 4 and a relative improvement of at least 10%. The difference is smallest when compared to SMLR. Note that most of the wins of SMM5 over SMLR are in multi-class domains (e.g., 4 of the 5 wins in the 7 base-level classifier case).

We next look into the influence of the diversity of the base-level classifiers on the performance improvement of SMM5 over the other combining methods. Figure 1 depicts the relative improvements as a function of the degree of error correlation for SMM5 vs. SELB and SMLR (for 3 and 7 base-level classifiers). The lines are fitted by linear regression to the bulleted points that represent domains where the differences between SMM5 and the other combiners are significant. It is clear that relative improvement increases as error correlation decreases (the lower the error correlation, the higher the diversity): this indicates that SMM5 uses the diversity of the base-level classifiers better than the competing combining methods.

4.3 The influence of the number of base-level classifiers

Studying the differences between Tables 2 and 3, we can note that the relative performance of the different combining methods is not affected too much by the change of the number of base-level classifiers. GRAD and SMDT seem to be affected most. The relative performance of GRAD improves, while that of SMDT worsens, when we go from 3 to 7 base-level classifiers: GRAD becomes better than VOTE, while SMDT becomes ever-so-slightly worse than SELB. SMM5 and SMLR are clearly the best in both cases.

5. Conclusions and further work

We have empirically evaluated several state-of-the-art methods for constructing ensembles of heterogeneous classifiers with stacking and shown that they perform (at best) comparably to selecting the best classifier from the ensemble by cross validation. We have propose a new method for stacking, that uses multi-response model trees at the meta-level. We have shown that it clearly outperforms existing stacking approaches and selecting the best classifier from the ensemble by cross validation.

It is not a surprise that stacking with multi-response model trees performs better than stacking with multi-response linear regression. The results of Frank et al. (Frank et al., 1998), who investigate classification via regression, shows that classification via model trees performs extremely well, i.e., better than multi-response linear regression and better than C5.0 (a successor of C4.5 (Quinlan, 1993)), especially in domains with continuous attributes. This indicates that multi-response model trees are a very suitable choice for learning at the meta-level, as confirmed by our experimental results.

Note that our approach is intended for combining classifiers that are heterogeneous (derived by different learning algorithms, using different model representations) and strong (i.e., each of the base-level classifiers performs relatively well in its own right), rather than homogeneous and weak. It is not intended, for example, for combining many classifiers derived by a single learning algorithm on subsamples of the original dataset. Given this, however, our experimental results indicate that stacking with multi-response model trees is a good choice for learning at the meta-level, regardless of the choice of the base-level classifiers.

While conducting this study and a few other recent studies (Ženko et al., 2001; Todorovski & Džeroski, 2002), we have encountered quite a few contradictions

between claims in the recent literature on stacking and our experimental results. For example, (Merz, 1999) claims that SCANN is clearly better than the oracle selecting the best classifier (which should perform even better than SelectBest). (Ting & Witten, 1999) claim that stacking with MLR clearly outperforms SelectBest. Finally, (Seewald & Fürnkranz, 2001) claim that both grading and stacking with MLR perform better than SelectBest. A comparative study including the datasets in the recent literature and a few other stacking methods (such as SCANN) should resolve these contradictions and provide a clearer picture of the relative performance of different stacking approaches. We believe this is a worthwhile topic to pursue in near-term future work.

We also believe that further research on stacking in the context of base-level classifiers created by different learning algorithms is in order, despite the current focus of the machine learning community on creating ensembles with a single learning algorithm with injected randomness or its application to manipulated training sets, input features and output targets. This should include the pursuit for better sets of meta-level features and better meta-level learning algorithms.

Acknowledgements

Many thanks to Ljupčo Todorovski for the cooperation on combining classifiers with meta-decision trees and the many interesting and stimulating discussions related to this paper. Thanks also to Alexander Seewald for providing his implementation of grading in Weka. We acknowledge the support of the EU funded project METAL: ESPRIT IV No. 26.357.

References

- Aha, D., Kibler, D., & Albert, M. K. (1991). Instance-based learning algorithms. *Machine Learning*, 6, 37–66.
- Blake, C., & Merz, C. (1998). UCI repository of machine learning databases.
- Dietterich, T. G. (1997). Machine-learning research: Four current directions. *AI Magazine*, 18, 97–136.
- Dietterich, T. G. (2000). Ensemble methods in machine learning. *Proceedings of the First International Workshop on Multiple Classifier Systems* (pp. 1–15). Berlin: Springer.
- Frank, E., Wang, Y., Inglis, S., Holmes, G., & Witten, I. H. (1998). Using model trees for classification. *Machine Learning*, 32, 63–76.
- Gama, J. (1999). Discriminant trees. *Proceedings of the Sixteenth International Conference on Machine Learning* (pp. 134–142). San Francisco: Morgan Kaufmann.
- John, G. C., & Leonard, E. T. (1995). K*: An instance-based learner using an entropic distance measure. *Proceedings of the 12th International Conference on Machine Learning* (pp. 108–114). San Francisco: Morgan Kaufmann.
- John, G. H., & Langley, P. (1995). Estimating continuous distributions in Bayesian classifiers. *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence* (pp. 338–345). San Francisco: Morgan Kaufmann.
- Kohavi, R. (1995). The power of decision tables. *Proceedings of the Eighth European Conference on Machine Learning* (pp. 174–189).
- Merz, C. J. (1999). Using correspondence analysis to combine classifiers. *Machine Learning*, 36, 33–58.
- Quinlan, J. R. (1992). Learning with continuous classes. *Proceedings of the Fifth Australian Joint Conference on Artificial Intelligence* (pp. 343–348). Singapore: World Scientific.
- Quinlan, J. R. (1993). *C4.5: Programs for machine learning*. San Francisco: Morgan Kaufmann.
- Seewald, A. K., & Fürnkranz, J. (2001). An evaluation of grading classifiers. *Proc. Fourth International Symposium on Intelligent Data Analysis* (pp. 221–232). Berlin: Springer.
- Ting, K. M., & Witten, I. H. (1999). Issues in stacked generalization. *Journal of Artificial Intelligence Research*, 10, 271–289.
- Todorovski, L., & Džeroski, S. (2000). Combining multiple models with meta decision trees. *Proceedings of the Fourth European Conference on Principles of Data Mining and Knowledge Discovery* (pp. 54–64). Berlin: Springer.
- Todorovski, L., & Džeroski, S. (2002). Combining classifiers with meta decision trees. *Machine Learning*. *In press*.
- Wang, Y., & Witten, I. H. (1997). Induction of model trees for predicting continuous classes. *Proceedings of the Poster Papers of the European Conference on Machine Learning*. Prague: University of Economics, Faculty of Informatics and Statistics.
- Witten, I. H., & Frank, E. (1999). *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*. San Francisco: Morgan Kaufmann.
- Wolpert, D. (1992). Stacked generalization. *Neural Networks*, 5, 241–260.
- Ženko, B., Todorovski, L., & Džeroski, S. (2001). A comparison of stacking with MDTs to bagging, boosting, and other stacking methods. *Proceedings of the First IEEE International Conference on Data Mining* (pp. 669–670). Los Alamitos: IEEE Computer Society.

Table 4. Relative improvement in accuracy (in %) of stacking with multi-response model trees (SMM5) as compared to other combining algorithms and its significance (+/- means significantly better/worse, x means insignificant).

DATASET	3 BASE LEVEL CLASSIFIERS					7 BASE LEVEL CLASSIFIERS				
	VOTE	SELB	GRAD	SMDT	SMLR	VOTE	SELB	GRAD	SMDT	SMLR
AUSTRALIAN	-3.46 x	-3.68 -	-1.75 x	-3.79 -	-0.92 x	-1.97 x	3.91 +	1.40 x	5.29 +	-2.07 x
BALANCE	50.99 +	48.68 +	50.27 +	48.68 +	53.89 +	50.79+	41.13 +	50.16 +	40.91 +	52.51 +
BREAST-W	18.60 +	-4.79 x	23.64 +	-4.79 x	-3.14 x	25.88+	-0.53 x	25.88 +	-0.53 x	0.53 x
BRIDGES-TD	7.45 x	7.45 +	3.25 x	9.15 +	-3.47 x	-1.91 x	4.76 x	-1.91 x	10.11 +	-0.63 x
CAR	76.56 +	73.91 +	75.05 +	69.67 +	72.89 +	79.45+	75.69 +	74.02 +	62.83 +	67.35 +
CHES	58.89 +	-0.52 x	48.39 +	-0.52 x	-0.52 x	61.10+	-3.66 x	48.57 +	-3.66 x	-3.66 x
DIABETES	-0.38 x	3.94 +	0.64 x	2.58 x	-1.37 x	0.22 x	-4.06 -	1.34 x	0.91 x	-1.48 x
ECHO	5.48 x	0.00 x	9.05 +	0.28 x	3.47 +	2.22 x	-5.60 -	2.22 x	1.25 x	-2.33 x
GERMAN	0.87 x	2.80 +	1.73 x	2.46 +	-2.50 -	3.45 x	5.76 +	4.67 +	4.28 +	-0.22 x
GLASS	-5.35 -	2.48 x	-1.67 x	1.62 x	-1.06 x	2.90 x	0.56 x	5.63 x	2.01 x	-1.71 x
HEART	8.44 +	2.31 x	11.51 +	2.31 x	-2.42 x	8.15+	1.83 x	9.32 +	4.68 +	1.15 x
HEPATITIS	14.07 +	5.69 x	18.60 +	5.69 x	4.53 x	1.57 x	-0.40 x	6.37 +	3.47 x	4.21 x
HYPO	42.72 +	-4.80 x	4.76 x	4.38 x	-4.80 x	50.10+	-2.93 x	26.13 +	42.39 +	-1.23 x
IMAGE	3.39 x	0.61 x	14.49 +	-11.97 -	0.15 x	-6.76 x	32.19 +	-3.72 x	16.99 +	-1.50 x
IONOSPHERE	8.73 x	22.03 +	18.73 +	25.81 +	10.85 +	7.36+	6.42 x	8.28 +	10.36 +	-10.80 -
IRIS	-6.35 x	5.63 x	-1.51 x	5.63 x	0.00 x	-4.00 x	-18.18 x	-6.85 x	-18.18 x	-5.41 x
SOYA	1.52 x	7.91 +	9.92 +	5.81 +	7.91 +	5.02 x	-2.35 x	-0.69 x	-0.46 x	13.52 +
TIC-TAC-TOE	97.18 +	72.83 +	95.70 +	72.83 +	55.35 +	92.42+	71.74 +	88.98 +	71.74 +	57.37 +
VOTE	52.75 +	5.20 x	35.68 +	5.20 x	5.20 x	39.34+	3.51 x	26.99 +	3.51 x	-1.23 x
WAVEFORM	13.92 +	5.05 +	19.68 +	4.94 +	4.44 +	18.99+	4.00 +	19.66 +	2.64 +	13.83 +
WINE	-74.19 -	6.90 x	-68.75 -	6.90 x	-5.88 x	-38.46 x	12.20 x	-38.46 x	7.69 x	2.70 x
AVERAGE	32.00	17.36	29.20	16.73	13.35	29.60	16.08	24.86	16.89	12.42
W/L	10+/2-	9+/1-	13+/1-	8+/2-	7+/1-	10+/0-	7+/2-	12+/0-	11+/0-	5+/1-

Table 5. Error rates (in %) of the learned ensembles of classifiers.

DATASET	3 BASE LEVEL CLASSIFIERS						7 BASE LEVEL CLASSIFIERS					
	VOTE	SELB	GRAD	SMDT	SMLR	SMM5	VOTE	SELB	GRAD	SMDT	SMLR	SMM5
AUSTRALIAN	13.81	13.78	14.04	13.77	14.16	14.29	13.99	14.84	14.46	15.06	13.97	14.26
BALANCE	8.91	8.51	8.78	8.51	9.47	4.37	10.14	8.48	10.02	8.45	10.51	4.99
BREAST-W	3.46	2.69	3.69	2.69	2.73	2.82	3.65	2.69	3.65	2.69	2.72	2.70
BRIDGES-TD	15.78	15.78	15.10	16.08	14.12	14.61	15.39	16.47	15.39	17.45	15.59	15.69
CAR	6.49	5.83	6.10	5.02	5.61	1.52	6.73	5.69	5.32	3.72	4.24	1.38
CHES	1.46	0.60	1.16	0.60	0.60	0.60	1.59	0.60	1.20	0.60	0.60	0.62
DIABETES	24.01	25.09	24.26	24.74	23.78	24.10	24.10	23.11	24.38	24.27	23.70	24.05
ECHO	29.24	27.63	30.38	27.71	28.63	27.63	30.92	28.63	30.92	30.61	29.54	30.23
GERMAN	25.19	25.69	25.41	25.60	24.36	24.97	24.08	24.67	24.39	24.29	23.20	23.25
GLASS	29.67	32.06	30.75	31.78	30.93	31.26	25.79	25.19	26.54	25.56	24.63	25.05
HEART	17.11	16.04	17.70	16.04	15.30	15.67	17.26	16.15	17.48	16.63	16.04	15.85
HEPATITIS	17.42	15.87	18.39	15.87	15.68	14.97	16.39	16.06	17.23	16.71	16.84	16.13
HYPO	1.32	0.72	0.80	0.79	0.72	0.76	1.56	0.76	1.05	1.35	0.77	0.78
IMAGE	2.94	2.85	3.32	2.53	2.84	2.84	1.92	3.03	1.98	2.47	2.02	2.05
IONOSPHERE	7.18	8.40	8.06	8.83	7.35	6.55	8.52	8.43	8.60	8.80	7.12	7.89
IRIS	4.20	4.73	4.40	4.73	4.47	4.47	5.00	4.40	4.87	4.40	4.93	5.20
SOYA	6.75	7.22	7.38	7.06	7.22	6.65	6.71	6.22	6.33	6.34	7.36	6.37
TIC-TAC-TOE	9.24	0.96	6.08	0.96	0.58	0.26	3.58	0.96	2.46	0.96	0.64	0.27
VOTE	7.10	3.54	5.22	3.54	3.54	3.36	6.25	3.93	5.20	3.93	3.75	3.79
WAVEFORM	15.90	14.42	17.04	14.40	14.33	13.69	16.64	14.04	16.78	13.85	15.65	13.48
WINE	1.74	3.26	1.80	3.26	2.87	3.03	1.46	2.30	1.46	2.19	2.08	2.02
AVERAGE	11.85	11.22	11.90	11.17	10.92	10.40	11.51	10.79	11.41	10.97	10.76	10.29