# Ensembles of wrappers for automated feature selection in fish age classification

*Sergio Bermejo*

Departament d'Enginyeria Electrònica,

Universitat Politècnica de Catalunya (UPC),

Jordi Girona 1-3 (C4 building), 08034 BARCELONA, SPAIN

PHONE: +34 4016758, E-MAIL: sergio.bermejo@upc.edu

**Abstract.** In feature selection, the most important features must be chosen so as to decrease the number thereof while retaining their discriminatory information. Within this context, a novel feature selection method based on an ensemble of wrappers is proposed and applied for automatically select features in fish age classification. The effectiveness of this procedure using an Atlantic cod database has been tested for different powerful statistical learning classifiers. The subsets based on few features selected, e.g. otolith weight and fish weight, are particularly noticeable given current biological findings and practices in fishery research and the classification results obtained with them outperforms those of previous studies in which a manual feature selection was performed.

## 1. Introduction

One of the most challenging problems in the field of pattern recognition (PR) is feature extraction (Guyon et al., 2006), which aims finding the most compact and discriminative set of properties or "features" presented in data. Although many research in feature extraction has been addressed to automate such a process, it has traditionally been considered a task much more problem- or domain-dependent than others in PR (Duda et al., 2001) since a good knowledge of the domain could be used to obtain such features, at least tentatively.

Fish age classification, a PR task of vital relevance among others for stock assessment and management (Girdler et al., 2010), usually relies on such manual procedures for feature extraction. In this direction, several fish features have been proposed for use in statistical fish age prediction and classification, with special emphasis in recent years to fish otolith features based on Fourier descriptors (Fablet and Le Josse, 2005; Galley et al., 2006) and different morphological parameters (Burke et al., 2008; Bermejo et al., 2007; Robotham et al., 2010; Hua et al., 2012).

However, the generalization error of statistical classifiers –i.e. their ability to mistake new examples taken on the same problem– tends to increase as of the number of features (Raudys and Jain, 1991) and, accordingly, the use of an arbitrary number of them leads to poor performance. One example of such behavior was demonstrated in (Bermejo, 2014) using multi-class support vector machines for fish age classification of an Atlantic cod database. Hence, if automatic feature extraction methods were additionally employed for reducing the complexity of the feature space a better performance could presumably be obtained. Other important benefits of such strategy includes speeding up computation (e.g. decreasing training times) and data understanding or reverse engineering (i.e. to increase knowledge about the problem, which can be of vital significance in natural sciences like fisheries science).

While some authors (e.g. Webb, 2002) consider feature extraction a process only concerning transformation of the original variables, it is generally agreed that feature extraction comprises the following steps: feature construction or generation that performs some kind of preprocessing –e.g. a linear or non-linear transformation– of the original raw variables (Theodoridis and Koutroumbas, 2008) and feature selection (Guyon and Elisseeff, 2003) that chooses a subset of the original or transformed variables.

There are three main approaches to feature selection (Blum and Langley, 1997; Guyon and Elisseeff, 2003, 2006): filter methods, wrappers and embedded methods. While filters can be viewed as a preprocessing step since they select a subset of variables independently of the chosen predictor (e.g. a classifier), wrappers use it as a black box or subroutine to score subsets of variables and embedded methods perform variable selection in its training phase. In this way, wrappers are based on an arguably better estimate of accuracy obtained with the predictor that will employ the feature subset than a separate measure that may have a completely unrelated inductive bias, but, at the expense of a higher computational cost (Blum and Langley, 1997). However, the inherent variance (or instability) of feature subset selection methods (Guyon and Elisseeff, 2006) produces a plethora of very different subsets attained for different conditions, i.e. different parameter tuning, small perturbations of the dataset or presence of redundant features.

In this paper, a novel wrapper that use a form of ensemble learning (Dietterich, 2003), which are based on a strategic combination of several predictors, have been proposed to attain a greater stabilization and thus a better generalization of the feature selection process. Feature subsets obtained with the ensemble of wrappers which employ as base classifiers support vector machines and nearest neighbor classifiers allow achieving a classification performance that outperforms a previous study

79    (Bermejo, 2014). Moreover, these subsets that have very few features, e.g. only otolith

80    weight and fish weight, are of relevance in accordance with recent findings in fisheries

81    research.

82

83    **2. Materials and methods**

84

85    *2.1. Atlantic cod database*

86    This dataset contains morphological and biological features for codfish age

87    classification. Traditional methods for determining the age of fish usually focus on

88    analyzing hard parts of the body, such as otoliths, which are small particles in the inner

89    ear composed of a gelatinous matrix and calcium carbonate, since the macroscopic

90    growth patterns of otholiths are correlated with the fish' age.

91    The fish database consists of one hundred forty-five Atlantic cod of known age

92    (varying from two to six years) from the Plateau stock that were hatched the same year

93    and later kept and reared in pen cages. This dataset was created from originally fish of

94    known-age sampled at different years in captivity since a number of samples were

95    recaptured once a year.  Otoliths were taken from this stock and weighed and also four

96    morphological features were recorded following an image analysis method defined in

97    (Bermejo et al., 2007).  Additionally, fish length, weight and sex were available for each

98    sample.

99    The leave-one-out (LOO) error using a 1-NN rule (Devroye et al., 1996; pp. 407-

100    421) were computed for this set (19.31%) as a way to estimate the Bayes error, i.e. the

101    minimum amount of classification error achievable. In a previous study with this

102    database using SVMs (Bermejo, 2014), the minimum obtained error was 21.79% for

103    otolith weight, fish length, weight and sex acting as features, which is lower than an

104    error rate of 22% obtained for a related dataset, combining five experts' readings, who

105 were given low and intermediate levels of information about fishes and the conditions

106 that they were obtained (Doering-Arjes et al., 2008). According to the above

107 considerations, some improvement in accuracy is still possible with SVMs taking the

108 value of the LOO estimate as an approximate lower bound to the attainable

109 misclassification rate. Table 1 displays the results of the LOO estimate and also

110 includes other relevant information of this dataset. A more comprehensive description

111 of the cod database is presented in (Bermejo, 2014).

112

113 *2.2. Ensemble of wrappers*

114　　Ensemble learning methods, such as bagging, boosting and variants (Bauer and

115 Kohavi, 1999) are based on the formation of a set of predictors $\{\varphi(\boldsymbol{x};\boldsymbol{D}_k)\}$ trained on a

116 sequence of learning sets $\{\boldsymbol{D}_k\}$, which are typically generated from a single dataset $\boldsymbol{D}$

117 using a resampling technique such as bootstrapping (Efron and Tibshirani, 1994). The

118 second core element of any ensemble method is a combination strategy: the most

119 obvious and effective procedure for combining a sequence of $K$ predictors $\{\varphi_k\}$ whose

120 outputs are continuous is averaging (Breiman, 1996a), i.e. $\overline{\varphi} = \sum_k \varphi_k / K$. Ensembles

121 have been built specifically to select features; for example, variants of AdaBoost for

122 feature selection have been proposed using decision stumps (Long and Vega, 2003) and

123 a mutual information measure (Liu et al., 2008), random subspace methods have also

124 been employed in feature ranking for removal of irrelevant variables (e.g. Tuv et al.,

125 2009), and ensembles based on bootstrapping have been combined with recursive

126 feature elimination and feature ranking (Windeatt et al., 2007). Furthermore, several

127 studies have analyzed the use of averaging and voting for the combination of multiple

128 feature selection criteria with the hope that several criteria would reflect different

129 properties in feature subsets (e.g. Somol et al., 2009), although none of them has

130  analyzed the effect of these procedures using a sole criterion to obtain a single feature

131  subset. Our proposal addresses this problem in the context of wrappers.

132  Wrappers (Kohavi, 1995) select features from a pool of feature sets based on a

133  decision rule of the form $\varphi_W = \arg\min_j L_{CV}\left(C_b^j; \boldsymbol{D}\right)$, that is, they select the j[th] feature

134  set for which $L_{CV}\left(C_b^j; \boldsymbol{D}\right)$ is the minimum, where $L_{CV}$ is the cross-validation error based

135  on the dataset $\boldsymbol{D}$ computed in the base classifier $C_b^j = C\left(\boldsymbol{x}^j; \boldsymbol{D}\right)$, whose inputs belong to

136  the j[th] feature set space. If the database is divided into a learning set $\boldsymbol{D}$ for performing

137  cross-validation and a test set $\boldsymbol{T}$ for final assessment of the classifier after feature

138  selection, a sequence of learning sets $\{\boldsymbol{D}_k\}$ and test sets $\{\boldsymbol{T}_k\}$ can be generated for

139  different random splits of the database. Then, and in accordance to the theoretical

140  analysis given in (Breiman, 1996a, 1996b), we propose in this paper a stabilized feature

141  selection rule that can be obtained through averaging over $L_{CV}$ in order to stabilize the

142  metric used in wrappers directly, so the feature selection rule based on an ensemble of

143  wrappers (EW) can be computed as $\overline{\varphi}_{EW} = \arg\min_j \left(\sum_k L_{CV}\left(C_{b_k}^j; \mathbf{D}_k\right)\big/K\right)$. The proposed

144  stabilization of the assessment criterion can be simply seen as an averaging of several *k-*

145  fold cross-validation estimates (based on the output of the wrapper's base classifier)

146  similarly to the way in which the outputs of several classifiers are stabilized through

147  averaging. The reader is referred to Breiman, 1996a, 1996b for further discussion, and

148  definition, of stability.

149  A baseline algorithm for feature selection with wrappers using internal cross-

150  validation (Flach, 2012) is suggested in Algorithm no. 1. The ensemble approach using

151  rule $\overline{\varphi}_{EW}$ is detailed in Algorithm no. 2 as a straightforward variation of the baseline

152  algorithm, in which feature selection is postponed until all the splits obtained in the first

153  version are evaluated. In this way, the second algorithm uses the same amount of

154 computational resources than the first one but a single decision on what features are

155 more relevant is obtained averaging over all these splits.

156 *2.3. Base classifiers*

157 Reducing the instability of the base classifiers would make it possible to evaluate

158 the degree of stability achieved by $\overline{\varphi}_{EW}$ with respect to $\varphi_W$ and could also provide

159 additional insight into how the stabilized decision rules work. Specifically, if the

160 induction algorithm $C_{\mathbf{D}_k}^j$ is completely stable on a sequence of learning sets $\{\mathbf{D}_k\}$, then

161 $C^j = C\left(\mathbf{x}^j;\mathbf{D}_i\right) = C\left(\mathbf{x}^j;\mathbf{D}_k\right) for\ \forall i,k$. Thus, the metric $\sum_k L_{CV}\left(C^j;\mathbf{D}_K\right)/K = \overline{L}_{CV}\left(C^j\right)$,

162 where $\overline{L}_{CV}$ denotes an averaged form of the cross-validation error computed using

163 different random replicates of the original database. As $K$ augments, $\overline{L}_{CV}$ will use more

164 samples from the database than $L_{CV}$, which is based on a single replicate, and can thus

165 presumably obtain a better estimation. Following this rationale, two well-known stable

166 induction algorithms, SVMs and NNs, have been employed as base classifiers in

167 wrappers.

168 SVMs (Vapnik, 1998), which has been developed in accordance with main results of

169 statistical learning theory, have also obtained a practical success in a range of practical

170 problems that makes them an appreciated part of many practitioners' toolbox. Multi-

171 class SVMs (Hsu and Lin, 2002) are a required extension of two-class SVMs that deal

172 with R-class classification problems, with R>2. In the experiments, we used two multi-

173 class SVMs implemented in the Spider library (Weston et al., 2006): 1) 1-vs-R ("one-

174 against-all") SVMs (Steinwart and Christmann, 2008), and 2) 1-vs-1 ("one-against-

175 one") SVMs (Schölkopf and Smola, 2001). Other SVM algorithms also implemented in

176 the library were ruled out in a previous round of experiments, since the results obtained

177 with them were outperformed by both 1-vs-R and 1-vs-1 SVMs.

178      Nearest-neighbor classifiers (Duda et al., 2001; pp.161-214) remain one of the

179    simplest yet most valuable nonparametric classification procedures. Given a set of

180    labeled prototypes $P$, the $k$-NN algorithm assigns the test point $x$ to that class majority

181    among its $k$ nearest neighbors belonging to $P$. In the experiments reported, the 1-NN,

182    also simply denoted as the NN rule, was used, since it has less computational burden

183    than the $k$-NN rule. Although the NN rule is sub-optimal with respect to the $k$-NN rule

184    in terms of the asymptotic error probability (i.e. with an unlimited number of

185    prototypes), its error rate is never worse than twice the Bayes error (Devroye et al.,

186    1996; pp. 61-90).

187

188    *2.4. Statistical assessment of experiments*

189      As pre-processing, whitening –i.e. mean removal and scaling by the variance of each

190    feature– was performed on the dataset so as to prevent the negative effect of their very

191    different scaling on the SVMs and NNs, and thus improving dramatically their

192    classification accuracy (see e.g. Ali and Smith-Miles, 2006). In (Bermejo, 2014), the

193    positive effect of such standardization is specifically discussed for this dataset.

194      A previous round of simple experiments was done to limit the set of values for the

195    parameters of the multi-class SVMs. According to the results obtained, radial basis

196    function (RBF) kernels $K(\mathbf{x}, \mathbf{x}_i) = \exp\left(\|\mathbf{x} - \mathbf{x}_i\|^2 / 2\sigma\right)$ were selected with a kernel width of

197    $\sigma = \{5, 10, 15, 20, 25\}$, while the rest of the parameters involved were the default values

198    defined in the Spider library (Weston et al., 2006).

199      The whole training set was chosen as nearest-neighbor prototypes in order to reduce

200    the computational burden due to the use of the learning algorithm. This brute-force

201    strategy, which usually works better than significant condensing and editing, achieves

202    competitive results with learning algorithms that compute a reduced number of

203    prototypes (see e.g. Bermejo, 2000).

204  Since the datasets here are medium- and small-sized, it was considered preferable to

205 maximize the learning set size in order to get enough training data. Thus, test sets were

206 formed containing only 25% of the database the test set size according to common

207 practices found in the literature; in particular, test sets ranged from 50% to 25% of the

208 complete database in fourteen datasets from the STATLOG project (Michie et al.,

209 1994). Accordingly, the datasets were first randomly divided using stratification into a

210 test set $T_i$ (25%) and a learning set $D_i$ (75%) for each split i=1,…,I of the database (with

211 I=75 when SVMs are used as the base classifiers and $K$=100 for NNs). Then, $D_i$ was

212 divided using stratification into five equal-sized parts or folds (i.e. n=5) that maintained

213 approximately the original proportion of data belonging to each class; in order to reduce

214 variance in the estimates of classification accuracy, this random division of $D_i$ was

215 repeated ten times, forming a sequence of folds. Thus, steps 5-13 of Algorithms 1 and 2

216 were repeated ten times and results conveniently averaged; in the case of SVMs, a

217 sequence of classifiers using a kernel width of σ={5,10,15,20,25} was also generated

218 for each split i, each feature set j and fold, and only those classifiers with parameters

219 obtaining, on average, the best results on the validation set were retained for testing.

220 Finally, the relative frequency with which the rule $\bar{\varphi}_{EW}$ outperforms or equals $\varphi_W$

221 defined by $\gamma = \sum_i 1\left(Err_i\left(\bar{\varphi}_{EW};\mathbf{T}_i\right) \leq Err_i\left(\varphi_W;\mathbf{T}_i\right)\right)/I$ was computed in order to compare

222 Algorithms 1 and 2.

223

224 **3. Results and discussion**

225  As Table 2 shows, on average, the use of $\bar{\varphi}_{EW}$ improves accuracy, since

226 $Err\left(\bar{\varphi}_{EW}\right) < Err\left(\varphi_W\right)$ for all the classifiers (see also Fig. 1). Also, for each data split i,

227 feature selection done by averaging mainly improves the results achieved by classifiers

228 based on feature sets selected using cross-validation, since $\gamma \in [.75,.96]$ (see also Fig. 2).

229  While the feature selection rule $\overline{\varphi}_{EW}$ generates a single feature set (see Table 2), $\varphi_W$

230  generates a population of feature sets, which only sometimes coincides with $\overline{\varphi}_{EW}$ (these

231  cases are shown as points along the line depicted in Fig. 2). On the other hand, feature

232  sets obtained by $\overline{\varphi}_{EW}$ are not unique with respect to the problem, but depend on the

233  wrapper's base classifier. However, although there is not a total consensus among the

234  classifiers, features set obtained by the selection rule $\overline{\varphi}_{EW}$ are particularly coherent with

235  biological findings, since fish weight (W) and otolith weight (OW) –i.e. the features

236  selected when 1-vs-R SVMs are used as base classifiers– and fish length (L), which is

237  also included when NN classifiers are used, are known to be highly correlated with age

238  and are often used in automatic fish age estimation or classification (Lou et al., 2005,

239  2007; Metin and Ilkyak, 2008; Ochwada et al., 2008; Pino et al., 2004), although other

240  researchers have proposed the use of other features, such as otolith growth rings (Fablet

241  and Le Josse, 2005; Guillaud et al., 1999, 2000; Rodin et al., 1996) or otolith shape

242  (Bird et al., 1986; Campana and Casselman, 1993; Castonguay et al., 1991).

243  Additionally, and more importantly, the feature set obtained by the selection rule $\overline{\varphi}_{EW}$

244  (based only on OW and W) in combination with 1-vs-R SVMs achieves an average test

245  error (20,93%) that outperforms best results computed with previous SVM experiments

246  (Bermejo, 2014) with the same dataset in which feature set selection was performed

247  manually (21,79%).

248    The feature selection rule $\overline{\varphi}_{EW}$ makes it possible to compute a single feature set with

249  the additional information obtained by generating different splits of the original

250  database. Since the repetition of experiments for different splits seems to be

251  recommended to reduce variance in test results (at least for small databases), $\overline{\varphi}_{EW}$ can

252  be used in this context at no extra computational cost. In order to extend this procedure

253  to datasets with a greater number of features, the brute-force search can be replaced

254 with the inspection of a pool of candidates obtained by ordering the feature set space by

255 leave-one-out error, since the minimum leave-one-out errors are obtained for feature

256 sets quite similar to those computed by $\overline{\varphi}_{EW}$ (see Table 1). Also, search strategies

257 (Guyon, 2006; pp.119-136) applied to large dimensionality domains in the context of

258 wrappers (Gheyas and Smith, 2010) are useful for obtaining a feature set subspace

259 where $\overline{\varphi}_{EW}$ and the experimental procedure suggested above were run with moderate

260 computational resources.

261

262 **4. Conclusions**

263 A metric based on averaging, a well-known method employed in ensemble learning for

264 stabilizing, has been proposed to reduce the instability of the feature subset selection

265 process performed by wrappers and has been tested on an Atlantic cod dataset using

266 SVMs and NN classifiers as base classifiers. As shown, a single feature subset can be

267 obtained in such a form of ensemble of wrappers and used to reverse engineer or better

268 explain data. Features selected in fish age classification are particularly noticeable in

269 view of current biological findings and practices in fishery research and outperforms

270 SVM classification accuracies obtained with manual feature selection (Bermejo, 2014).

271

276

277

278

279

11

280 **References**

281 Ali, S., Smith-Miles, K.A., 2006. Improved support vector machine generalization using

282      normalized input space. In: Sattar, A., Kang, B.H. (Eds.), Advances in Artificial

283      Intelligence Lecture Notes in Computer Science, vol. 4304. Springer-Verlag, Berlin,

284      pp. 362–371.

285 Bauer, E., Kohavi, R., 1999. An empirical comparison of voting classification

286      algorithms: bagging, boosting, and variants. Mach. Learn. 36, 105-139.

287 Bermejo, S., 2000. Learning with Nearest Neighbour Classifiers. Ph.D. Dissertation,

288      Departament d'Enginyeria Electrònica, Universitat Politècnica de Catalunya.

289 Bermejo, S., Monegal, B., Cabestany, J., 2007. Fish age categorization from otolith

290      images using multi-class support vector machines. Fish. Res. 84, 247-253.

291 Bermejo, S., 2014. The benefits of using otolith weight in statistical fish age

292      classification: A case study of Atlantic cod species. Comp. Elec. Agr. 107, 1–7.

293 Bird, J.L., Eppler, D.T., Checkley, D.M., 1986. Comparisons of hearing otoliths using

294      Fourier series shape analysis. Can. J. Fish. Aquat. Sci. 43, 1228–1234.

295 Blum, A. L., Langley, P., 1997. Selection of relevant features and examples in machine

296      learning. Artif. Intell. 97, 245–271.

297 Breiman, L., 1996a. Heuristics of instability and stabilization in model selection. Ann.

298      Statist. 24, 2350-2383.

299 Breiman, L., 1996b. Bagging Predictors. Mach. Learn. 24, 123-140.

300 Burke, N., Brophy, D., King, P.A., 2008. Shape analysis of otolith annuli in Atlantic

301      herring (*Clupea harengus*); a new method for tracking fish populations. Fish. Res.

302      91, 133–143.

303 Campana, S.E., Casselman, J.M., 1993. Stock discrimination using otolith shape

304      analysis. Can. J. Fish. Aquat. Sci. 50, 1062–1083.

305    Castonguay, M., Simard, P., Gagnon, P., 1991. Usefulness of Fourier analysis of otolith

306        shape for Atlantic mackerel (*Scomber scombrus*) stock discrimination. Can. J. Fish.

307        Aquat. Sci. 48, 296–302.

308    Devroye, L., Györfi, L., Lugosi, G., 1996. A Probabilistic Theory of Pattern

309        Recognition. Springer-Verlag, Berlin.

310    Dietterich, T. G., 2003. Ensemble learning. In: M. A. Arbib (Ed.), The Handbook of

311        Brain Theory and Neural Networks, second edition. The MIT Press, Cambridge,

312        MA, pp. 405–408.

313    Doering-Arjes, P., Cardinale, M., Mosegaard, H., 2008. Estimating population age

314        structure using otolith morphometrics: a test with known-age Atlantic cod (*Gadus*

315        *morhua*) individuals. Can. J. Fish. Aquat. Sci. 65, 2342–2350.

316    Duda, R. O., Hart, P. E., Stork, D. G., 2001. Pattern Classification (2nd Edition). Wiley-

317        Inter-science, New York.

318    Efron, B., Tibshirani, R.J., 1994. An Introduction to the Bootstrap. Chapman & Hall,

319        London.

320    Fablet, R., Le Josse, N., 2005. Automated fish age estimation from otolith images using

321        statistical learning. Fish. Res. 72, 279 –290.

322    Flach, P. A., 2012. Machine Learning: The Art and Science of Algorithms that Make

323        Sense of Data. Cambridge University Press, Cambridge.

324    Galley, E. A., Wright, P. J., and Gibb, F. M., 2006. Combined methods of otolith shape

325    analysis improve identification of spawning areas of Atlantic cod. ICES J. Mar. Sci. 63,

326        1710-1717.

327    Girdler, A., Wellby, I., Welcomme, R., 2010. Fisheries Management: A Manual for

328        Still-Water Coarse Fisheries. Wiley-Blackwell, Oxford.

329    Gheyas, I. A., Smith, L. S., 2010. Feature subset selection in large dimensionality

330        domains, Pattern Recognit. 43, 5-13.

331    Guillaud, A., Ballet, P., Troadec, H., Rodin, V., Benzinou, A., Le Bihan, J., 1999. A

332        multiagent system for edge detection: an application to growth ring detection on fish

333        otoliths, in: Image Processing and its Applications, 1999 Seventh International

334        Conference Publication No. 465, Vol. 1, pp. 445-449.

335    Guillaud, A., Troadec, H., Benzinou, A., Rodin,  V., Le Bihan, J., 2000. Continuity

336        perception using a multiagent system: an application to growth ring detection on

337        fish otoliths, In: Pattern Recognition, 2000. Proceedings. 15th International

338        Conference, Vol. 2, pp. 519-522.

339    Guyon, I., Elisseeff, A., 2003. An introduction to variable and feature selection. J.

340        Mach. Learn. Res. 3, 1157-1182.

341    Guyon, I., Elisseeff, A., 2006.  An Introduction to Feature Extraction. In: I.Guyon, M.

342        Nikravesh, S. Gunn, L. A. Zadeh (Eds.), Feature Extraction Foundations and

343        Applications. Studies in Fuzziness and Soft Computing, Vol. 207, Springer, New

344        York, 2006, pp.1-26.

345    Guyon, I., Nikravesh, M., Gunn, S., Zadeh, L. A., (Eds.), 2006. Feature Extraction

346        Foundations and Applications. Vol. 2007, 1-778. Springer, New York.

347    Hua, J., Li. D., Duan, Q., Han, Y., Chen, G., Si, X., 2012. Fish species classification by

348        color, texture and multi-class support vector machine using computer vision.

349        Comput. Electron. Agr. 88, 133–140.

350    Hsu, C.-W., Lin, C.-J., 2002. A comparison of methods for multi-class support vector

351        machines. IEEE Trans Neural Netw. 13, 415-25.

352    Kohavi, R., 1995. Wrappers for performance enhancement and oblivious decision

353        graphs. Ph.D Thesis, Department of Computer Science, Stanford University.

354    Liu, T.-Y., Li, G.-Z., Yang, J.Y., Yang, M.Q., 2008. Feature selection for the

355        imbalanced QSAR problems by using EasyEnsemble. Int. J. Comp. Biol. Drug Des.

356        1, 334-346.

357    Long, P. M., Vega, V. B., 2003. Boosting and microarray data, Mach. Learn. 52, 31–44.

358    Lou, D.C., Mapstone, B.D., Russ, G.R., Davies, C.R., Begg, G.A., 2005. Using otolith

359        weight–age relationships to predict age-based metrics of coral reef fish populations

360        at different spatial scales. Fish. Res. 71, 279–294.

361    Lou, D.C., Mapstone, B.D., Russ, G.R., Davies, C.R., Begg, G.A., Davies, C.R., 2007.

362        Using otolith weight–age relationships to predict age based metrics of coral reef fish

363        populations across different temporal studies. Fish. Res. 83, 216–227.

364    Metin, G., Ilkyak, G.M., 2008. Use of otolith length and weight in age determination of

365        poor cod (*Trisopterus minutus* Linn., 1758). Turk J. Zool. 32, 293–297.

366    Michie, D., Spiegelhalter, D.J., Taylor, C.C., 1994. Machine Learning, Neural and

367        Statistical Classification, Prentice Hall, New York. Available online at:

368        http://www.amsta.leeds.ac.uk/~charles/statlog/.

369    Ochwada, F.A., Scandol, J.P., Gray, C.A., 2008. Predicting the age of fish using general

370        and generalized linear models of biometric data: a case study of two estuarine

371        finfish from New South Wales, Australia. Fish. Res. 90, 187–197.

372    Pino, C.A., Cubillos, L.A., Araya, M., Sepúlveda, A., 2004. Otolith weight as an

373        estimator of age in the Patagonian grenadier, *Macruronus magellanicus*, in central-

374        south Chile. Fish. Res. 66, 145–156.

375    Raudys, S. J., Jain, A. K., 1991.  Small sample size effects in statistical pattern

376        recognition: recommendations for practitioners. IEEE T. Pattern. Anal. 13, 252-264.

377    Robotham, H., Bosch, P., Gutiérrez-Estrada, J. C., Castillo, J. & Pulido-Calvo, I., 2010.

378        Acoustic identification of small pelagic fish species in Chile using support vector

379        machines and neural networks. Fish. Res. 102, 115–122.

380    Rodin, V., Troadec, H., de Pontual, H., Benzinou,  A., Tisseau, J., Le Bihan,  J., 1996.

381        Growth ring detection on fish otoliths by a graph construction, in: Proceedings of

382        the International Conference on Image Processing, Vol. 1, pp. 685-688.

383 Schölkopf, B., Smola, A. J., 2001. Learning with Kernels. The MIT Press, Cambridge.

384 Steinwart, I., Christmann, A., 2008. Support Vector Machines. Springer, Berlin.

385 Somol, P., Grim, J., Pudil, P., 2009. Criteria Ensembles in Feature Selection, in: J.A.

386   Benediktsson, J. Kittler, F. Roli (Eds.): MCS 2009, LNCS 5519, Spinger, Berlin, pp.

387   304–313.

388 Theodoridis, S., Koutroumbas, K., 2008. Pattern Recognition, Fourth Edition.

389   Academic Press, New York.

390 Tuv, E., Borisov, A., Runger, G., Torkkola, K. (2009). Feature Selection with

391   Ensembles, Artificial Variables, and Redundancy Elimination, J. Mach. Learn. Res.

392   10, 1341-1366.

393 Vapnik, V., 1998. Statistical Learning Theory, Wiley Inter-science, New York.

394 Webb, A. R., 2002. Statistical Pattern Recognition, 2nd Edition. Wiley, New York.

395 Weston, J., Elisseeff, A., BakIr, G., Sinz, F., 2006. The Spider (matlab toolbox).

396   Available online at: http://www.kyb.tuebingen.mpg.de/bs/people/spider/main.html.

397 Windeatt, T., Prior, M., Effron, N., Intrator, N., 2007. Ensemble-based Feature

398   Selection Criteria, in: Petra Perner (Ed.): Machine Learning and Data Mining in

399   Pattern Recognition (MLDM 2007), IBaI publishing, Leipzig, pp. 168-182.

400

401

**Algorithm 1** Baseline algorithm for wrappers based on internal cross-validation

1: For i=1,…,I

2:  Split database randomly into a test set $\mathbf{T}_i$ and a learning set $\mathbf{D}_i$ using a ratio 1:q where 1:q denotes the sampling ratio between $D_k$ and $T_k$, i.e. % of samples q/(1+q) is sampled for $D_k$ and % 1/(1+q) for $T_k$

3:  For j=1 to $2^m$ combinations of feature sets

4:  Obtain for feature space $j^{th}$ a subset $\mathbf{D}_i^j$ from $\mathbf{D}_i$ where j is a vector in a binary representation $\left(j_1 \cdots j_m\right)$ with $j_k$ denoting whether feature $k^{th}$ is present ('1') or not ('0') and $\mathbf{D}_i^j \in X^p, \mathbf{D}_i \in X^m, \mathbf{D}_i^j \subset \mathbf{D}_i, \mathbf{D}_i^j \in \Re^p, \mathbf{D}_i \in \Re^m, 0 < p \le m$

5:  Split $\mathbf{D}_i^j$ into n disjoint sets $\left\{\mathbf{D}_i^{j,k}, k = 1,..., n\right\}$, i.e. $\bigcup_{k=1}^{n} \mathbf{D}_i^{j,k} = \mathbf{D}_i^j, \bigcap_{k=1}^{n} \mathbf{D}_i^{j,k} = 0$

6:  For k=1 to n folds

7:  Obtain a training dataset $\mathbf{D}_i^{j,-k} = \bigcup_{m=1, m \ne k}^{n} \mathbf{D}_i^{j,m}$ and a validation set $\mathbf{V}_i^{j,k} = \mathbf{D}_i^{j,k}$

8:  Define a sequence of classifiers' parameters $\left\{\boldsymbol{\sigma}_l, l = 1,..., L\right\}$

9:  For l=1,…,L

10:  Compute classifier $C_l\left(\mathbf{x}^j; \mathbf{D}_i^{j,-k}, \boldsymbol{\sigma}_l\right)$ or, in short, $C_l\left(\mathbf{x}^j; \boldsymbol{\sigma}_l\right)$, i.e. a classifier $C_l\left(\mathbf{x}^j\right)$ working in feature space $X^p$ with $\mathbf{x}^j \in X^p$ using the training data set $\mathbf{D}_i^{j,-k}$ for the classifier's parameters $\boldsymbol{\sigma}_l$

11:  Obtain the cross-validation error for $C_l\left(\mathbf{x}^j; \boldsymbol{\sigma}_l\right)$ as the loss error for this classifier computed using $\mathbf{V}_i^{j,k}$, i.e. $L_{CV}\left(C_l^j\right) = L\left(C_l\left(\mathbf{x}^j; \boldsymbol{\sigma}_l\right), \mathbf{V}_i^{j,k}\right)$

12:  Choose the best classifier $C^k\left(\mathbf{x}^j\right)$ of the sequence $\left\{C_l\right\}$ with optimal parameters $\boldsymbol{\sigma}^k$ as the one that minimizes the cross validation (CV) error, i.e.

$$C^k\left(\mathbf{x}^j; \boldsymbol{\sigma}^k\right) = \arg_C \min_l L_{CV}\left(C_l^j\right) \text{ or } L_{CV}\left(C^{k,j}\right) = \min_{l=1,...,L} L_{CV}\left(C_l^j\right)$$

13:  Obtain mean CV error in $\mathbf{D}_i^j$ for feature space $j^{th}$ as $L_{CV}\left(\mathbf{D}_i^j\right) = \frac{1}{n} \sum_{k=1}^{n} L_{CV}\left(C^{k,j}\right)$

14:  Select the feature subset from which the mean CV error $L_{CV}\left(\mathbf{D}_i^j\right)$ is minimum, i.e. $\varphi_W(i) = \arg \min_j L_{CV}\left(\mathbf{D}_i^j\right)$

15:  Obtain the generation error $Err_i\left(\varphi_W(i); \mathbf{T}_i\right)$ of classifiers in feature space $\varphi_W(i)$

16: Compute the mean generalization error for the baseline wrapper $\varphi_W$ as

$$Err\left(\varphi_W\right) = \sum_{i=1}^{I} Err_i\left(\varphi_W(i); \mathbf{T}_i\right) / I$$

**Algorithm 2** Ensembles of wrappers (as a variation of Algorithm 1)

1: For i=1,…,I

2:    Split database randomly into a test set $\mathbf{T}_i$ and a learning set $\mathbf{D}_i$ using a ratio 1:q

3:    For j=1 to $2^m$ combinations of feature sets

4:       Obtain for feature space j$^{th}$ a subset $\mathbf{D}_i^j$ from $\mathbf{D}_i$ with

$$\mathbf{D}_i^j \in X^p, \mathbf{D}_i^j \in X^m, \mathbf{D}_i^j \subset \mathbf{D}_i, \mathbf{D}_i^j \in \mathfrak{R}^p, \mathbf{D}_i \in \mathfrak{R}^m, 0 < p \leq m$$

5:       Split $\mathbf{D}_i^j$ into n disjoint sets $\left\{\mathbf{D}_i^{j,k}, k=1,...,n\right\}$

6:       For k=1 to n folds

7:          Obtain $\mathbf{D}_i^{j,-k} = \bigcup_{m=1, m \neq k}^{n} \mathbf{D}_i^{j,m}$ and $\mathbf{V}_i^{j,k} = \mathbf{D}_i^{j,k}$

8:          Define a sequence of classifiers' parameters $\left\{\mathbf{\sigma}_l, l=1,...,L\right\}$

9:          For l=1,…,L

10:             Compute classifier $C_l\left(\mathbf{x}^j; \mathbf{D}_i^{j,-k}, \mathbf{\sigma}_l\right)$

11:             Obtain $L_{CV}\left(C_l^j\right) = L\left(C_l\left(\mathbf{x}^j; \mathbf{\sigma}_l\right), \mathbf{V}_i^{j,k}\right)$

12:             Choose $C^k\left(\mathbf{x}^j; \mathbf{\sigma}^k\right) = \arg_C \min_l L_{CV}\left(C_l^j\right)$ or

$$L_{CV}\left(C^{k,j}\right) = \min_{l=1,...,L} L_{CV}\left(C_l^j\right)$$

13:          Compute $L_{CV}\left(\mathbf{D}_i^j\right) = \frac{1}{n} \sum_{k=1}^{n} L_{CV}\left(C^{k,j}\right)$

14: For i=1,…,I

15:    Compute the mean CV error for feature space j$^{th}$ as $L_{CV}(j) = \frac{1}{I} \sum_{i=1}^{I} L_{CV}\left(\mathbf{D}_i^j\right)$

16: Select the feature subset from which the mean cross-validation $L_{CV}(j)$ is minimum,

i.e. $\varphi_{EW} = \arg\min_j L_{CV}(j)$

16: For i=1,…,I

17:    Obtain the generation error of classifiers in feature space $\varphi_{EW}$ for $\mathbf{T}_i$ as

$Err_i\left(\varphi_{EW}; \mathbf{T}_i\right)$

18: Compute the mean generalization error for the averaged wrapper $\varphi_{EW}$ as

$Err\left(\varphi_{EW}\right) = \sum_{i=1}^{I} Err_i\left(\varphi_{EW}; \mathbf{T}_i\right)/I$

465

| Size | No. of Features | Features / Feature vector | No. of Classes | Minimum Leave-one-out Error |
|------|-----------------|---------------------------|----------------|------------------------------|
| 145 | 8 | Fish sex (S), fish length (L), fish weigh (W), otolith weight (OW), otolith contour length (C), otolith area (A), otolith maximum internal distance (I), otolith maximum perpendicular distance (P) / (P I A C OW W L S) | 5 [fish age: 2 to 6] | 0.1931 [for feature set $12 = (00001100)_2$] |

466

467                         Table 1. Codfish dataset summary.

468

469

470

| | | $Err(\varphi_W)$ | $Err(\overline{\varphi}_{EW})$ | Feature vector(*) / $\overline{\varphi}_{EW}$ | $\gamma$ |
|---|---|---|---|---|---|
| SVM | 1-vs-1 | .2297 | .2147 | (P I A C OW W L S)/ 175=$(10101111)_2$ | .74567 |
| | 1-vs-R | .2273 | .2093 | (P I A C OW W L S)/ 12=$(00001100)_2$ | .96 |
| NN | | .2459 | .214 | (P I A C OW W L S)/ 14=$(00001110)_2$ | .84 |

471

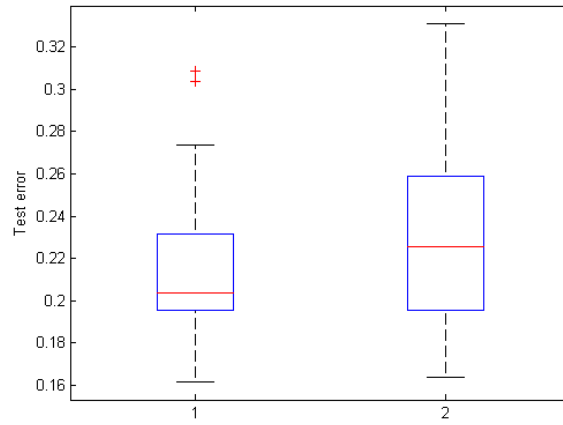472    Table 2. Comparison of feature set selection using averaging and cross-validation.

473    (* see Table 1 for further details)
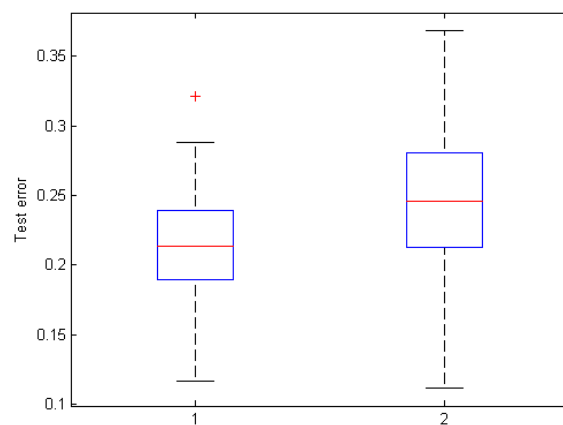
474

475

476

477

478

479           a)

480

481           b)

482
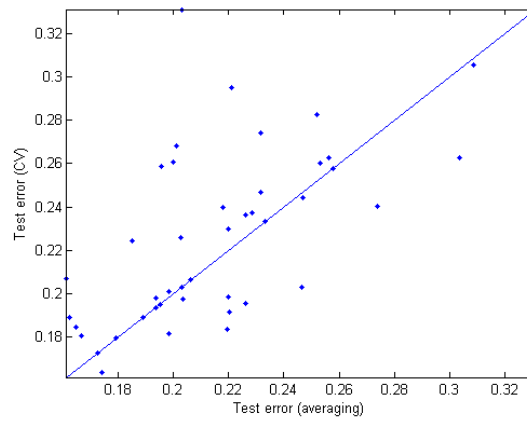
483           c)

484      Fig.1. Box plot of average test errors $Err(\overline{\varphi}_{EW})$ [left] and $Err(\varphi_W)$ [right] using: a) 1-

485           vs-1 SVMs, b) 1-vs-R SVMs and c) NN classifiers.

486
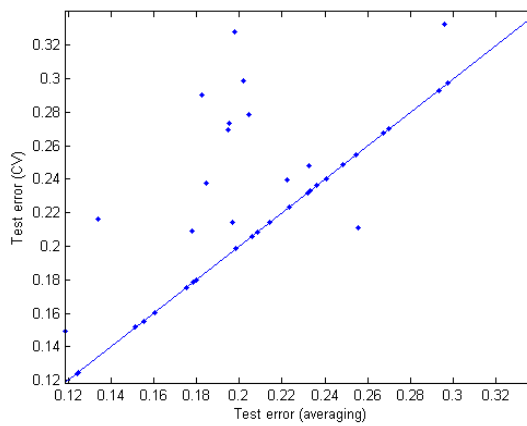
487



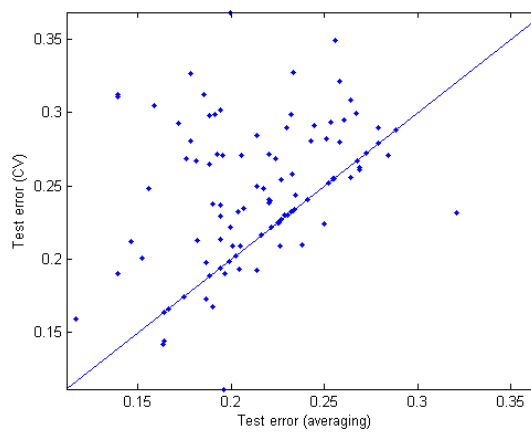488                              a)

489



490                              b)

491



492                              c)

493        Fig.2. Test errors of ensembles of wrappers based on averaging, $Err_i\left(\overline{\varphi}_{EW}\right)$, vs.

494    those based on internal CV, $Err_i\left(\varphi_W\right)$, for different $\mathbf{T}_i$ using a) 1-vs-1 SVMs, b) 1-vs-R

495                        SVMs and c) NN classifiers.