

Supervised machine learning and heterotic classification of maize (*Zea mays* L.) using molecular marker data

Leonardo Ornella^{a,1}, Elizabeth Tapia^{a,b}

^a*CIFASIS-CONICET, Av. 27 de Febrero 210bis, Rosario, Argentina*

^b*FCEIA-UNR, Department of Electronic, Riobamba 210 bis, Rosario, Argentina*

Abstract

The development of molecular techniques for genetic analysis has enabled great advances in cereal breeding. However, their usefulness in hybrid breeding, particularly in assigning new lines to heterotic groups previously established, still remains unsolved. In this work we evaluate the performance of several state-of-art multiclass classifiers onto three molecular marker datasets representing a broad spectrum of maize heterotic patterns. Even though results are variable, they suggest supervised learning algorithms as a valuable complement to traditional breeding programs.

Key words: Maize, Supervised Learning, Heterotic Groups,

1. Introduction

2 Since the first maize hybrid was bred and produced in USA, hybrid breed-
3 ing has become one of the primary goals in any maize breeding programs
4 ([Hallauer and Miranda, 1988](#)); however, varietal development has become
5 more competitive and costly. For example, in USA, development of one va-
6 riety of maize or soybean requires 0.5 – 7.0 million dollar. The lifetime of a
7 variety is usually 3-6 years before it succumbs to the challenges of the pro-
8 duction environment (biotic and abiotic stress) and demands of consumers
9 ([Lee, 1998](#)). Consequently, grouping parent lines into heterotic groups is
10 fundamental in both private and public breeding programs in order to re-
11 duce the number of crosses, and therefore field tests, necessary to evaluate

Email address: ornella@cifasis-conicet.gov.ar (Leonardo Ornella)

¹Corresponding Author. Tel.: +54 341 4821771 Int.104 - Fax: +54 341 4821772

12 potential high-yielding hybrids (Hallauer and Miranda, 1988). By heterotic
13 groups we mean a population of genotypes that, when crossed with indi-
14 viduals from another heterotic group or population, consistently outperform
15 intra-population crosses (Hallauer and Miranda, 1988). Molecular markers,
16 such as RAPD (random amplified polymorphic DNA), AFLP (amplified frag-
17 ment length polymorphism) and microsatellites, among others, have facili-
18 tated the development of new varieties by reducing the time required for the
19 detection of specific traits in progeny plants and the identification of disease
20 resistance genes (Korzun, 2003). Even though they have been proposed to
21 assign new inbred to heterotic groups previously established (dos Santos Dias
22 et al., 2004; Xia et al., 2004), their usefulness in this task still remains un-
23 certain (dos Santos Dias et al., 2004). Machine-learning techniques, such as
24 decision trees and artificial neural networks, are increasingly used in agricul-
25 ture to deal with classification, prediction, and modeling problems (Kirch-
26 ner et al., 2004; Mitchell et al., 1996); however, we found no reports about
27 machine learning algorithms (Kotsiantis, 2007; Witten and Frank, 2005) and
28 heterotic group assignment using molecular marker data. We conjecture that
29 traditional distance-based methods (Reif et al., 2005) currently available for
30 assigning new inbreds to heterotic groups in corn do not capture the possi-
31 ble non-linear relation between parental data and progeny performance (dos
32 Santos Dias et al., 2004; Springer and Stupar, 2007) and that such type of
33 non linearity may be easily captured by supervised machine learning models.

34 In this paper, we evaluate the performance of several state-of-art super-
35 vised learning algorithms on molecular marker data for heterotic assignation,
36 and delineate perspectives for further research.

37 2. Multiclass Classifiers

38 The goal of supervised learning is to build a concise model of the distri-
39 bution of class labels in terms of predictor features, the resulting classifier
40 is then used to assign class labels to the testing instances where the values
41 of the predictor features are known, but the value of the class label is un-
42 known (Kotsiantis, 2007). There are numerous learning algorithms reported
43 in the bibliography (Kotsiantis, 2007; Witten and Frank, 2005), for this intro-
44 ductory work we considered four well-known supervised learning algorithms
45 implemented in Weka workbench (Hall et al., 2009): i) Naive Bayes (John
46 and Langley, 1995), ii) Bayes Net (Friedman et al., 1997), iii) Simple Logistic

47 (Landwehr et al., 2005) and iv) Support Vector Machines (SVMs) with linear
48 and radial basis function kernels (Burges, 1998).

49 2.1. Naive Bayes

50 NB learns from training data the conditional probability of each attribute
51 A_i given the class label C . Classification is then done by applying Bayes rule
52 to compute the probability of C given the particular instance of A_1, \dots, A_n ;
53 and then predicting the class with the highest posterior probability. This
54 computation is rendered feasible by making a strong independence assumption:
55 all the attributes A_i are conditionally independent given the value of
56 the class C . Independence means probabilistic independence, i.e., A is inde-
57 pendent of B given C whenever $P(A|B, C) = P(A|C)$ for all possible values
58 of A, B and C , whenever $P(C) > 0$ (Friedman et al., 1997). Even though the
59 above assumption is clearly unrealistic, its predictive performance is compet-
60 itive with state-of-the-art classifiers (Friedman et al., 1997; Kohonen et al.,
61 2008).

62 2.2. Bayes Net

63 A Bayesian network is an annotated directed acyclic graph that encodes
64 a joint probability distribution over a set of random variables U (Friedman
65 et al., 1997). The graph G encodes independence assumptions: each variable
66 X_i is independent of its nondescendants given its parents in $G(\Pi_{x_i})$:

$$p(x_1, x_2, \dots, x_n) = \prod_{i=1}^n p(x_i | \Pi_{x_i}) \quad (1)$$

67 To use a BN as classifier, a search algorithm find a network B ,
68 $P_B(A_1, A_2, \dots, A_n, C)$, that best matches a training set D according to some
69 scoring function (Cooper and Herskovits, 1992; Friedman et al., 1997). Once
70 a network is learned, B returns the label c that maximizes the posterior
71 probability $P_B(c/a_1, \dots, a_n)$ (Cooper and Herskovits, 1992; Friedman et al.,
72 1997). Naive Bayes can be considered a Bayes Net in where the structure of
73 the graph is constrained (Friedman et al., 1997).

74 2.3. Simple Logistic

75 Landwehr et al. (2005) proposed Logistic Model Trees or LMT, trees that
76 contain linear logistic regression functions at the leaves. In that work they
77 reported that at low number of training instances ($n \leq 100$), Simple Logistic

78 (SL), a logistic model tree of size one, performs as well as more complex LMT
 79 and better than decision tree C4.5 (Quinlan, 1993), with less computational
 80 requirements (Landwehr et al., 2005).

81 Linear logistic regression models the posterior class probabilities $Pr(C =$
 82 $c|\mathbf{X} = \mathbf{x})$ for the J classes via functions linear in \mathbf{x} and ensures that they
 83 sum to one and remain in $[0, 1]$ (Sumner et al., 2005). The model is:

$$P(C = c|X = x) = \frac{e^{F_c(x)}}{\sum_{k=1}^C e^{F_k(x)}} \quad (2)$$

84 where $F_j(\mathbf{x}) = \sum_{m=1}^M f_{mj}(x) = \beta_j^T \cdot \mathbf{x}$. Estimates of β_j^T are obtained
 85 by numeric optimization algorithms that approach the maximum likelihood
 86 solution iteratively (Sumner et al., 2005). In Simple Logistic such iterative
 87 method is the LogitBoost algorithm (Landwehr et al., 2005). In each it-
 88 eration, it fits a least-squares regressor to a weighted version of the input
 89 data with a transformed target variable. y_{ij}^* are the binary pseudo-response
 90 variables which indicate group membership of an observation like this:

$$y_{ij}^* = \begin{cases} 1 & \text{if } y_i = j \\ 0 & \text{if } y_i \neq j \end{cases} \quad (3)$$

91 By constraining f_{mj} to be a linear function of only the attribute that
 92 results in the lowest squared error, we lead to an algorithm that performs
 93 automatic attribute selection (Sumner et al., 2005); also, by using cross-
 94 validation (5 folds) to determine the best number of LogitBoost iterations,
 95 only those attributes that improve the classification performance on unseen
 96 instances are included (Landwehr et al., 2005; Sumner et al., 2005).

97 2.4. Support Vector Machines

98 The support vector machine (SVM) algorithm is based on the statistical
 99 learning theory and the Vapnik-Chervonenkis (VC) dimension introduced by
 100 Vladimir Vapnik and Alexey Chervonenkis (Cortes and Vapnik, 1995); the
 101 underlying idea is to calculate a maximal margin hyperplane (the decision
 102 function) separating two classes of the data (Cortes and Vapnik, 1995), such
 103 decision function is fully specified by a usually small subset of the data (the
 104 support vectors) which defines the position of the separator. New samples
 105 are classified according to the side of the hyperplane they belong to (Cortes
 106 and Vapnik, 1995; Devos et al., 2009).

107 In the case of non separable data, the “ideal boundary” must be adapted
 108 to tolerate errors for some objects i:

$$\text{minimize } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \zeta_i \quad (4)$$

109 under the constraints $\zeta_i \geq 0$, $\zeta_i + y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 \geq 0$, \mathbf{w} and b
 110 are respectively the normal vector and the bias of the hyperplane, and each
 111 ζ_i corresponds to the distance between the object i and the corresponding
 112 margin hyperplane (Devos et al., 2009).

113 The parameter C is a regularization meta-parameter, when C is small,
 114 margin maximization is emphasized whereas when C is large, the error min-
 115 imization is predominant (Cortes and Vapnik, 1995; Devos et al., 2009).

116 To learn non-linearly separable functions, data are implicitly mapped
 117 to a higher dimensional space by means of mercer kernels which can be
 118 decomposed into a dot product, $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(x_i) \cdot \phi(x_j)$ (Burges, 1998).
 119 Examples of kernels are the linear kernel $K = (\mathbf{x}_i \cdot \mathbf{x}_j - 1)^{p=1}$ and the radial
 120 basis function kernel $K = e^{-\gamma(\mathbf{x}_i - \mathbf{x}_j)^2}$.

121 2.5. ECOC codes

122 SVMs have particular high generalization abilities and have become very
 123 popular in the recent years; nevertheless, they are inherently binary classifiers
 124 and a combination scheme is necessary to extend SVMs for problems with
 125 more than two classes (Rifkin and Klautau, 2004). In this work, the One
 126 Against All (OAA) (Rifkin and Klautau, 2004) and the Error Correcting
 127 Output Coding (ECOC) (Dietterich and Bakiri, 1995) combination schemes
 128 are used.

129 Briefly, OAA classifiers rely on the discrimination of individual classes
 130 against the others while ECOC codes are defined by a more general decom-
 131 position or “coding matrix” $M \in \{0, 1\}^{L \times N}$, which converts a L-multiclass
 132 problem into N binary tasks (Dietterich and Bakiri, 1995). There are several
 133 coding matrices reported in the bibliography (Allwein et al., 2000; Dietterich
 134 and Bakiri, 1995; Rifkin and Klautau, 2004). In particular, we work with ran-
 135 dom ECOC codes, each entry of the coding matrix chosen to be 0 or 1 with
 136 equal probability and N limited by the maximum number of different and
 137 non-complementary binary vectors that can be generated for dichotomization
 138 (Dietterich and Bakiri, 1995).

139 The original approach to ECOCs predicts the class whose corresponding
140 row vector has minimum Hamming distance to the vector of 0/1 predictions
141 obtained from the N classifiers (Dietterich and Bakiri, 1995). (Allwein et al.,
142 2000) presented an alternative, loss-based decoding, which notices the mag-
143 nitude of the predictions, sometimes interpreted as a measure of “confidence”
144 of a prediction. Several authors verified that Loss-decoding indeed produces
145 more accurate classifiers than the Hamming distance (Allwein et al., 2000;
146 Frank and Kramer, 2004; Rifkin and Klautau, 2004).

147 3. Materials and Methods

148 3.1. Datasets

149 We compiled three molecular marker datasets representing a broad spec-
150 trum of temperate and tropical germplasm. The **Liu Data** (Liu et al.,
151 2003) comprises 197 inbreeds (instances) of both temperate and tropical
152 germplasm characterized by 188 attributes derived from 94 microsatellites.
153 The number of distinct values per attribute ranges from 4 to 48 with a
154 mean of 18.18. Missing data represents a 4.75 % of the total, ranging
155 from 0% to 25.38%, depending on the attribute. Instances are distributed
156 into 10 heterotic groups (classes) and the number of instances per group
157 is {61, 13, 11, 8, 9, 13, 28, 17, 29, 8}. The **Morales Data** (Morales Yokobori
158 et al., 2005) comprises 26 temperate inbreeds of germoplasm characterized
159 by 42 attributes derived from 21 microsatellites. The number of distinct
160 values per attribute ranges from 2 to 13 with a mean of 4.72. Missing data
161 represents a 8.60% of a total, ranging from 0% to 42% of missing data per
162 attribute. Instances are distributed into 4 heterotic groups and the num-
163 ber of instances per group is {4, 8, 6, 8}. The **Xia Data** (Xia et al., 2004)
164 comprises 73 inbreeds of tropical germplasm characterized by 166 attributes
165 derived from 83 microsatellites. The number of distinct values per attribute
166 ranges from 2 to 14 with a mean of 5.93. Missing data represents the 8.02%
167 from the total, ranging from 0% to 43.84% of missing data per attribute.
168 Instances are grouped into 8 heterotic groups and the number of instances
169 per group is {22, 17, 7, 5, 5, 5, 5, 7}.

170 3.1.1. Classifiers

171 Simple Logistic, Naive Bayes and Bayes Nets were all implemented with
172 defaults parameters of Weka (Witten and Frank, 2005). SVMs were evaluated
173 using linear and radial basis function (RBF) kernels, both also with default

174 parameters ($C = 1$ for linear kernels and $C = 1, \gamma = 0.01$ for radial basis
175 function kernels). In both SVM alternatives, we choose the option “to fit Lo-
176 gistic regression models” of Weka’s SMO (Sequential Minimal Optimization)
177 algorithm for SVMs, which allows to emit an estimate of the confidence for
178 the binary prediction instead of (0,1) hard outputs.

179 Concerning the implementation of ECOC classifiers, in a preliminary re-
180 search we evaluated the data with variable length codes and we did observed
181 a positive correlation between ECOC accuracy and code length. As a trade
182 off between classifier’s performance and computational complexity we choose
183 random codes of length $N = 6$ for Morales Data, $N = 55$ for Xia data and
184 $N = 75$ for Liu data. Therefore, 75 SVMs were used for the ECOC classifi-
185 cation of Liu data, 55 for Xia data, and 6 for Morales data. The multiclass
186 schemes were implemented as a new WEKA classifier and integrated into the
187 original package (Witten and Frank, 2005).

188 3.1.2. Evaluation of classifier’s performance

189 The predictive power of supervised learning algorithms on molecular
190 marker data was evaluated by means of the error rate (Borra and Ciac-
191 cio, 2005) and the Cohen’s Kappa coefficient (Cohen, 1960) exhibited across
192 30 Montecarlo runs of stratified 10-Fold Cross Validation (CV) experiments
193 (Kirchner et al., 2004; Kohavi, 1995). At each Montecarlo run, the data was
194 split into 10 different segments of almost the same size and containing app-
195 proximately the same proportion of categories as the original dataset. For each
196 segment, classifiers were respectively trained and evaluated on the samples
197 derived by omitting the selected segment and on selected segment. At the end
198 of this procedure, the average classification error and the average Kappa co-
199 efficient were reported. The choice of the Kappa coefficient was motivated by
200 its ability to better measure the agreement between binary inter-annotators
201 than the traditional classification error. In particular, the Kappa coefficient
202 takes into account chance agreements (Cohen, 1960; Kirchner et al., 2004)
203 and it is well suited for unequal class distribution datasets.

204 Two main classification scenarios were considered: i) NB, BN, SL, OAA-
205 rbf (SVM with radial basis function), ECOC-rbf, OAA-linear (SVM with
206 linear kernel) and ECOC linear classifiers on full molecular marker data, and
207 ii) the same classifiers evaluated on reduced data derived by the application
208 of feature selection algorithms.

209 *3.1.3. Missing data*

210 Regarding missing data, all associated to nominal attributes, imputation
211 depends on the classifier evaluated (X.Su et al., 2008). In Weka, Naive Bayes
212 ignores the missing values whereas SMO globally replaces all missing values
213 by a default value, e.g., “unknown” (X.Su et al., 2008). Finally, in Bayes
214 Net and Simple Logistic classification, missing values of training and test set
215 are filled in using the mode of the corresponding attribute valuated on the
216 training data (Bouckaert, 2008; Landwehr et al., 2005).

217 *3.1.4. Statistical comparison among classifiers*

218 It is important to assess whether the observed difference in classification
219 performance is statistically significant or simply due to chance (Luengo et al.,
220 2009). Comparisons of arithmetic means and visual inspection of Kappa
221 boxplots was supplemented with Kolmogorov-Smirnov (KS-test) provided
222 by the R² environment (stats package). KS is a nonparametric test and it
223 has the advantage of making no assumption about the distribution of data
224 (Luengo et al., 2009). For each dataset and condition evaluated (Full and
225 reduced data derived by the application of feature selection algorithms), all
226 possible pairs of (A,B) Kappa coefficients distributions were assessed under
227 the alternative hypothesis “distribution B is greater than distribution A”
228 (The R Development Core Team, 2009)

229 *3.2. Feature Selection*

230 Reducing the feature space to non-redundant features results in improved
231 classification accuracy and helps avoid overfitting of the classifiers. In this
232 study, we mainly experimented with Correlation-based Feature Subset selec-
233 tion (CFS) (Hall, 2000). The CFS strategy uses a correlation-based heuristic
234 to evaluate the merit of feature subsets with respect to classification cat-
235 egories and the correlation between features. CFS selection implemented
236 in WEKA is fully automatic and does not require a priori specification of
237 the number of features to be included in the final subset (Hall, 2000). In
238 addition, we applied a second feature selection method, Relief (Kononenko,
239 1994), to Morales Data. This method ranks the worth of an attribute by
240 repeatedly sampling an instance and considering the value of the given at-
241 tribute for the nearest instance of the same and different class (Kononenko,

²<http://www.rproject.org/>

242 1994). In other words, Relief assigns more weight to those attributes that
243 have the same value for instances from the same class and differentiate be-
244 tween instances from different classes (Witten and Frank, 2005). The Relief
245 algorithm was calibrated in order to retain 25, 50 and 75% of the original
246 number of attributes.

247 3.2.1. SVM parameters optimization

248 Optimization of the meta-parameters, C (regularization parameter) of
249 linear kernel and C and γ (RBF kernel), is the key step in SVM performance
250 (Devos et al., 2009). Globally, when C is small the margin maximization
251 is emphasized leading to large margin and smooth boundary. The number
252 of support vectors included in the solution depends on this parameter and,
253 usually, if the number of support vectors is high the solution is unstable and
254 leads to poor classification performance. (Devos et al., 2009; Forman and
255 Cohen, 2004). Also, when the value of γ is large, the separating boundary
256 has a large number of support vectors and can become tortuous. Again,
257 this risks overfitting the training set data to yield an SVM model that is
258 not robust. In contrast, a small value of γ can lead to separating boundaries
259 described with a small number of support vectors but that may be too smooth
260 to classify the training set examples with sufficient accuracy (Devos et al.,
261 2009; Jorissen and Gilson, 2005). In RBF kernels it has been reported that
262 different combinations of C and γ lead to similar classification rates (Devos
263 et al., 2009). To perform the optimization we implemented an exhaustive
264 grid search: 30 points ($C = 0.25, 0.5, 1, 2, 4$ and $G = 0.0001, 0.001, 0.01,$
265 $0.1, 1, 10$) for radial basis function kernel and 5 points ($C = 0.25, 0.5, 1,$
266 $2, 4$) for the linear kernel. This approach enables to visualize directly the
267 effect of both parameters and provides useful information about core SVM
268 classifiers. In order to minimize the risk of overfitting, all parameters were
269 estimated by external leaving out one Cross Validation (Morales) or 10 fold
270 Cross Validation (Liu and Xia datasets) over the training data (Ambroise
271 and McLachlan, 2002).

272 4. Results and Discussion

273 Three native multiclass classifiers plus Support Vector Machines classi-
274 fiers under the OAA (Rifkin and Klautau, 2004) and ECOC frameshifts (Di-
275 etterich and Bakiri, 1995) were evaluated on three molecular marker datasets
276 representing a broad spectrum of maize heterotic patterns. Generalization

277 error of classifiers in this domain was estimated by means of the error-rate
278 and the Kappa Cohen’s Coefficient. Error-rate, defined as the ratio between
279 the number of misclassified cases and the total number of cases examined,
280 is the common measure used in nonparametric classification models (Borra
281 and Ciaccio, 2005). However, it does not compensate for classifications that
282 might have been due to chance. Hence, we also used the Cohen’s Kappa
283 as a statistically robust alternative, especially in datasets with an unequal
284 distribution of classes. Both statistics were determined by 30 runs of Mon-
285 tecarlo 10-Fold CV experiments. Arithmetic means of these statistics, with
286 and without feature selection, are shown in Table 2. It can be observed that
287 results according to mean error-rate and Kappa values do not always agree.
288 For example, in Liu Full data, SL and NB display identical error rates and
289 different kappa values; in Liu CFS reduced data the four SVM ensembles
290 rank different either we consider kappa or error rate values; also in Xia CFS
291 data OAA schemes rank different whatever we choose error rate or kappa
292 (Table 2). Overall, classification results seem to be problem-dependent, in-
293 definite and not always normal. Therefore arithmetics means may be not
294 always provide representative measures of classification performance. Conse-
295 quently, comparison of means and visual inspection of Kappa boxplots was
296 supplemented with Kolmogorov-Smirnov (KS) tests (Luengo et al., 2009).
297 We recall that KS is a nonparametric test which does not rely on an assump-
298 tion of normality (Luengo et al., 2009).

299 4.1. Results on Full Data

300 Bayes Net exhibited the best *mean* performance on full Liu Data (Table
301 2). Visual inspection of Kappa boxplots and KS test agreed with this result
302 (Figure 1). All KS tests were significant when comparing the rest of classifiers
303 to BN. For example, $p\text{-value} = 6.55e-05$ when comparing ECOC-rbf and
304 OAA-rbf (the closest classifiers according to kappa coefficient) to BN.

305 In Xia Data, ECOC-rbf significantly exceeds the rest of classifiers (Table
306 2 and Figure 2). In all KS tests (any classifier vs. ECOC-rbf) the null
307 hypothesis was rejected; as an example, $p\text{-value} = 0.0015$ when comparing
308 ECOC-linear (the second ranked classifier) against this ensemble.

309 Finally, Simple Logistic exhibited the best mean performance on full
310 Morales data (Table 2), a fact that was confirmed by corresponding Kappa
311 boxplots (Figure 3). Moreover, when comparing the rest of the classifiers
312 with Simple Logistic using KS, the highest $p\text{-value}$ obtained was 0.0006, i.e.,
313 all null hypotheses were rejected. Concerning SL, our results are in agreement

314 with Landwehr et al. (2005). When evaluating Liu and Xia data, which are
315 more complex with respect to the number of classes, the number of attributes
316 and the number of instances, the classifier displayed the worst performance
317 (Figures 1 and 2). Even though, we included this classifier in the analysis
318 because its good performance on Morales data, and this dataset is similar,
319 with regard to number of instances and/or attributes, to most works reported
320 in the literature, specially those from development countries (dos Santos Dias
321 et al., 2004).

322 4.2. Impact of Feature Selection

323 The genetic basis of heterosis has been debated for nearly a century
324 without a clear resolution. The two main hypotheses that advanced to
325 explain this phenomenon are dominance and overdominance (Hallauer and
326 Miranda, 1988; Springer and Stupar, 2007). It is also well documented that
327 not all markers will be linkage to genes or QTL (quantitative trait locus)
328 associated with heterosis (Austin et al., 2000). Moreover, the diploid nature
329 of data and the characteristics of the instances (homozygous lines) allow
330 us to infer the existence of some redundancy in attributes. Therefore, we
331 implemented CFS (Correlation-based Feature Selection) in order to remove
332 attributes not related to the class. The number of CFS selected attributes
333 was variable, depending on the dataset; extreme values ranged from 13.83 to
334 47.62 % of the initial number of features (Table 1).

335 Almost none of the classifiers improve their performance with filtered data
336 (Table 2 and boxplots). The only exception were Naive Bayes and Bayes net
337 evaluated on Xia Data (Figure 2). Even though, ECOC-rbf was still the best
338 classifier, all ks tests were statistically significant when comparing the rest
339 of classifiers to this ensemble.

340 In Morales reduced data and according to arithmetic means (Table 2 and
341 boxplot of Figure 3) SL was still the best classifier. However, when ECOC
342 linear (with default parameters) was compared to SL, the p-value was 0.0672.
343 The rest of classifiers did show significant p-values in KS test. Finally, in Liu
344 Data, though Naive Bayes degraded its performance with CFS filtering, like
345 the rest of the classifiers, it ranked second after Bayes Net (p-value < 0.05).

346 Theory suggests that interactions between genes associated with molec-
347 ular markers could play an important role in the generation of the observed
348 heterosis (Dudley and Johnson, 2009). Hence, it may be possible that us-
349 ing filters that contemplate interactions between attributes could lead to
350 improved classification performance.

351 *4.3. Data Complexity*

352 Molecular marker data showed to be complex enough to require the care-
353 ful exploration of non-trivial multiclass classifiers: the attribute-class rela-
354 tionship is possibly non-linear (dos Santos Dias et al., 2004; Springer and Stu-
355 par, 2007) and datasets present noisy and/or missing features (Jones et al.,
356 1997). Also, the dimensionality of molecular marker data is between that
357 of the classic Machine Learning setting ($n/p > 10$) (Asuncion and Newman,
358 2007; Kohavi, 1995) and that posed by recent challenging microarray data
359 classification problems ($n/p \ll 1$) (Mukherjee et al., 2003), where n is the
360 number of instances and p the number of attributes. Actually, the number
361 of classes ranges from 4 to 10 and the number of instances per class is gen-
362 erally less than 30, which is a very low number of training instances (dos
363 Santos Dias et al., 2004; Liu et al., 2003; Morales Yokobori et al., 2005; Xia
364 et al., 2004).

365 When comparing classifiers performance on full data scenarios we did
366 observe significant differences between Liu, Xia and Morales data results
367 (Table 2). Kappa values ranging between 0.61-0.80 indicate a substantial
368 agreement between observed and predicted data whereas values below 0.20
369 indicate only a slight agreement (Landis and Koch, 1977).

370 From a genetic point of view, differences of methods used to established
371 the heterotic groups could be reflecting differences between mechanisms relat-
372 ing attributes (molecular markers) with classes (heterotic groups): heterotic
373 groups of Xia and Morales data where established on the basis of field essays
374 (topcross or diallel) and, according to Xia et al. (2004), the mixed genetic
375 constitution of the populations and pools of Cymmit germplasm (Xia data)
376 made the task of assigning them to genetically diverse and complementary
377 heterotic groups difficult. A similar situation was reported for Morales data
378 (Eyh rabide et al., 2006). Liu data clusters, on the other side, were estab-
379 lished on the basis of genetic origin (Liu et al., 2003) so it was easy to assign
380 new lines to groups solely on molecular data.

381 From a Machine Learning point of view, these differences could be due
382 to a challenging ratio between the number of instances (n) and the number
383 of attributes (p) of training data (Kohavi, 1995; Mukherjee et al., 2003). For
384 example, for microarray data (extremely low n/p ratios) achieving error rates
385 around 0.1-0.2% requires in the order of 75-100 training samples (Mukherjee
386 et al., 2003), whereas Kohavi (1995) reported error rates from 5.8 to 53.2%
387 when working with datasets comprising a number of instances and a number
388 of attributes similar to those used in this work. However, if the modest

389 classification performance for Morales and Xia databases is only due to the
390 n/p ratios (specially for Morales data set), a good feature selection method
391 should be able to improve the results. It can be seen from figures 2 and 3 that
392 attribute CFS selection didn't improve the accuracy of the classifiers. We
393 performed an additional experiment on Morales dataset using another filter
394 method implemented in Weka, Relief (Kononenko, 1994), and selecting 25,
395 50 and 75% of the original number of attributes. Filtered data was evaluated
396 with Simple Logistic and the four SVM ensembles as stated in Materials and
397 Methods. It can be seen from Figure 4 that, except a few and non-significant
398 exceptions, all classifiers degraded their performance at increasingly higher
399 n/p ratios.

400 It has been reported that SVM classifiers are quite sensitive to meta-
401 parameters (Devos et al., 2009; Rifkin and Klautau, 2004). However, we
402 couldn't observe a significant enhancement of ensembles performance with
403 the optimization of the meta-parameters (C in linear kernel and C and γ in
404 radial basis function kernel). None of the optimized linear SVM-ensembles
405 significantly outperformed their standard counterparts (Table 3). In Xia data
406 both, OAA and ECOC, optimized RBF ensembles outperformed classifiers
407 with default values provided by Weka (Table 4). In Morales data, only OAA-
408 RBF showed a significant improvement with optimized parameters (Table 4).
409 With respect to Morales data, this is reasonable because with small training
410 sets optimization of parameters, even by cross-validation, may only lead to
411 over fitting the training set (Forman and Cohen, 2004). Surprisingly, in
412 Liu data none of the optimized SVM ensembles (significantly) outperformed
413 their counterparts with default parameters. This could be attributed to the
414 number of missing data and the imputation technique of SMO (X.Su et al.,
415 2008), or to the robustness of ensembles to base classifier error (Dietterich
416 and Bakiri, 1995).

417 Overall, we should assume that despite the specific relation between pa-
418 rameters n , p , L and the specific relationship between attributes and classes,
419 if we apply the incorrect model, classification performance will be poor. In
420 this sense, above results shed light on how to process molecular marker in-
421 formation to be useful in the problem of assigning new lines to previously
422 established heterotic groups.

423 5. Summary and conclusions

424 The information on germplasm diversity and relationships among elite
425 materials is a fundamental importance in crop improvement (Hallauer and
426 Miranda, 1988). Assigning lines to different heterotic groups would avoid the
427 development and evaluation of many of the crosses that would eventually be
428 discarded (Terron et al., 1997). Our proposal was to complement traditional
429 breeding using molecular markers information and supervised learning al-
430 gorithms. Three well-known multiclassifiers and support vector machine (a
431 binary classifier) with linear and radial basis function kernels and under two
432 decomposition schemes were evaluated using three molecular datasets rep-
433 resenting a broad spectrum of maize heterotic patterns. Morales dataset
434 includes 26 lines, mostly derived from orange flint (temperate) germplasm,
435 clustered in four heterotic groups by topcross field essays (Eyh rabide et al.,
436 2006), Liu data includes 248 inbred lines of importance to temperate breeding
437 and many important tropical and subtropical lines (Liu et al., 2003) and Xia
438 data 73 inbreds of tropical germplasm grouped mainly by diallel (Xia et al.,
439 2004). We also used CFS filtering to improve classifiers performance, but we
440 only obtained a slight improvement in Xia data. We also evaluated Relief
441 filtering on Morales data, with negative results. However, CFS removes noisy
442 attributes non-correlated between them and theory suggest that interactions
443 between genes associated with molecular markers could play an important
444 role in the generation of the observed heterosis (Pea et al., 2008) so filters
445 that contemplates this situation remains to be explored. Finally, although
446 results obtained with heterotic groups established by field essays (top cross
447 or diallel) are modest, there is a strong evidence that using data with more
448 training instances could generate successful classifiers. Also it is necessary
449 to evaluate other algorithms; the potential impact, in time and money, on
450 crop sustainability makes our research worth to try: while traditional genetic
451 breeding requires expensive field tests and a time scale in the order of years
452 for obtaining an heterotic assignment, in our proposed framework costs are
453 significantly lower and the time scale is in the order of weeks, two weeks for
454 growing an small plant plus a week to obtain molecular data and a couple of
455 days for computational analysis.

456 6. Acknowledgments

457 Authors would like to acknowledge Dr. Marilyn Warburton from Cimmyt,
458 Mexico, for permitting use of Xia data. We also thanks J. Coronel and

459 L. Angelone for technical support. Authors also would like to thank the
460 reviewers for their comments that help improve the manuscript. E. Tapia's
461 work is supported by project PICT 11-15132, Agencia Nacional de Promoción
462 Científica y Tecnológica from Argentina . Leonardo Ornella is a postdoctoral
463 fellow of CONICET, Argentina.

464 **References**

465 Allwein, E. L., Schapire, R. E., Singer, Y., 2000. Reducing Multiclass to
466 Binary: A Unifying Approach for Margin Classifiers. *Journal of Machine*
467 *Learning Research* 1, 113–141.

468 Ambroise, C., McLachlan, G. J., 2002. Selection bias in gene extraction on the
469 basis 405 of microarray gene-expression data. *Proceedings of the National*
470 *Academy of Sciences* 99, 6562–6566.

471 Asuncion, A., Newman, D., 2007. UCI Machine Learning Repository, Uni-
472 versity of California, Irvine, School of Information and Computer Sciences.

473 Austin, D. F., Lee, M., Veldboom, L. R., Hallauer, A. R., 2000. Genetic
474 Mapping in Maize with Hybrid Progeny Across Testers and Generations:
475 Grain Yield and Grain Moisture. *Crop Sci* 40 (1), 30–39.

476 Borra, S., Ciaccio, A., 2005. Methods to compare nonparametric classifiers
477 and to select the predictors. In: Vichi, M., Monari, P., Mignani, S., Mon-
478 tanari, A. (Eds.), *New Developments in Classification and Data Analysis*.
479 Springer, pp. 11–19.

480 Bouckaert, R. R., 2008. Bayesian Network Classifiers in Weka for Version
481 3-5-7.

482 Burges, C. J. C., 1998. A Tutorial on Support Vector Machines for Pattern
483 Recognition. *Data Mining and Knowledge Discovery* 2, 121–167.

484 Cohen, J., 1960. A coefficient of agreement for nominal scales. *Educational*
485 *and Psychological measurements* 20, 37–46.

486 Cooper, G. F., Herskovits, E., 1992. A bayesian method for the induction of
487 probabilistic networks from data. *Machine Learning* 9 (4), 309–347.

- 488 Cortes, C., Vapnik, V., 1995. Support-Vector Networks. *Machine Learning*
489 20, 273–297.
- 490 Devos, O., Ruckebusch, C., Durand, A., Duponchel, L., Huvenne, J. P., 2009.
491 Support vector machines (SVM) in near infrared (NIR) spectroscopy: Fo-
492 cus on parameters optimization and model interpretation. *Chemometrics*
493 *and Intelligent Laboratory Systems* 96 (1), 27 – 33.
- 494 Dietterich, T. G., Bakiri, G., 1995. Solving Multiclass Learning Problems via
495 Error-Correcting Output Codes. *Journal of Artificial Intelligence Research*
496 2, 263–286.
- 497 dos Santos Dias, L., de Toledo Picoli, E., Rocha, R. B., Alfenas, A. C., 2004.
498 A priori choice of hybrid parents in plants. *Genet Mol Res* 3, 356–368.
- 499 Dudley, J. W., Johnson, G. R., 2009. Epistatic Models Improve Prediction
500 of Performance in Corn. *Crop Sci* 49 (3), 763–770.
- 501 Eyhétabide, G., Nestares, G., Hourquescos, M., 2006. Development of a
502 heterotic pattern in orange flint maize. In: Lamkey, K., Lee, M. (Eds.),
503 *Plant Breeding: The Arnel R. Hallauer International Symposium*. Black-
504 well Publishing, pp. 352–379.
- 505 Forman, G., Cohen, I., 2004. Learning from Little: Comparison of Classifiers
506 Given Little Training. In: 8th European Conference on Principles and
507 Practice of Knowledge Discovery in Databases (PKDD). pp. 161–172.
- 508 Frank, E., Kramer, S., 2004. Ensembles of Nested Dichotomies for Multi-
509 Class Problems. In: *Proceedings of the 21st International conference of*
510 *Machine Learning (ICML-2004)*. ACM Press, pp. 305–312.
- 511 Friedman, N., Geiger, D., Goldszmidt, M., 1997. Bayesian Network Classi-
512 fiers. *Machine Learning* 29, 131–163.
- 513 Hall, M. and Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., Wit-
514 ten, I. H., 2009. The WEKA data mining software: an update. *SIGKDD*
515 *Explor. Newsl.* 11, 10–18.
- 516 Hall, M. A., 2000. Correlation-based Feature Selection for Discrete and Nu-
517 meric Class Machine Learning. In: *Proc. 17th International Conf. on Ma-*
518 *chine Learning*. Morgan Kaufmann, San Francisco, CA, pp. 359–366.

- 519 Hallauer, A. R., Miranda, J. B., 1988. Quantitative Genetics in Maize Breed-
520 ing, 2nd Edition. Iowa State University Press, Ames.
- 521 John, G. H., Langley, P., 1995. Estimating Continuous Distributions in
522 Bayesian Classifiers. In: Eleventh Conference on Uncertainty in Artificial
523 Intelligence. Morgan Kaufmann, pp. 338–345.
- 524 Jones, C., Edwards, K., Castaglione, S., Winfield, M., Sala, F., 1997. Repro-
525 ducibility testing of RAPD, AFLP and SSR markers in plants by a network
526 of european laboratories. *Molecular Breeding* 3, 381–390.
- 527 Jorissen, R. N., Gilson, M. K., 2005. Virtual screening of molecular databases
528 using a support vector machine. *Journal of Chemical Information and
529 Modeling* 45, 549–561.
- 530 Kirchner, K., Tölle, K., Krieter, J., 2004. The analysis of simulated sow
531 herd datasets using decision tree technique. *Computers and Electronics in
532 Agriculture* 42, 111–127.
- 533 Kohavi, R., 1995. A Study of Cross-Validation and Bootstrap for Accuracy
534 Estimation and Model Selection. In: *IJCAI*. pp. 1137–1145.
- 535 Kohonen, J., Talikota, S., Corander, J., Auvinen, P., Arjas, E., 2008. A Naive
536 Bayes classifier for protein function prediction. In *Silico Biology* 9, 0003.
- 537 Kononenko, I., 1994. Estimating attributes: Analysis and extensions of relief.
538 In: Bergadano, F., Raedt, L. D. (Eds.), *European Conference on Machine
539 Learning*. Springer, pp. 171–182.
- 540 Korzun, V., 2003. Molecular markers and their application in cereals breed-
541 ing. In: *Marker Assisted Selection: A fast track to increase genetic gain in
542 plant and animal breeding Session I: MAS in plant*. Tech. rep., FAO.
- 543 Kotsiantis, S. B., 2007. Supervised Machine Learning: A Review of Classifi-
544 cation Techniques. *Informatica* 31, 249–268.
- 545 Landis, J. R., Koch, G. G., 1977. The measurement of observer agreement
546 for categorical data. *Biometrics* 33 (1), 159–174.
- 547 Landwehr, N., Hall, M., Frank, E., 2005. Logistic Model Trees. *Machine
548 Learning* 95 (1-2), 161–205.

- 549 Lee, M., 1998. Genome projects and gene pools: New germplasm for plant
550 breeding? Proc. Natl. Acad. Sci. USA 95, 2001–2004.
- 551 Liu, K., Goodman, M., Muse, S., Smith, J., Buckler, E., Doebley, J., 2003.
552 Genetic Structure and Diversity Among Maize Inbred Lines as Inferred
553 From DNA Microsatellites. Genetics 165, 2117–2128.
- 554 Luengo, J., García, S., Herrera, F., 2009. A study on the use of statistical
555 tests for experimentation with neural networks: Analysis of parametric test
556 conditions and non-parametric tests. Expert Systems with Applications
557 36 (4), 7798–7808.
- 558 Mitchell, R. S., Sherlock, R. A., Smith, L. A., 1996. An investigation into the
559 use of machine learning for determining oestrus in cows. Computers and
560 Electronics in Agriculture 15, 195–213.
- 561 Morales Yokobori, M., Decker, V., Ornella, L., 2005. Analysis of heterotic
562 maize (*Zea mays* L.) populations using molecular markers. Maize Genetics
563 Cooperation Newsletters 79, 36.
- 564 Mukherjee, S., Tamayo, P., Rogers, S., Rifkin, R., Engle, A., Campbell, C.,
565 Tr, G., Jp, M., 2003. Estimating Dataset Size Requirements for Classifying
566 DNA Microarray Data. Computational Biology 10, 119–142.
- 567 Pea, G., Ferron, S., Gianfranceschi, L., Krajewski, P., Pè, M. E., 2008. Gene
568 expression non-additivity in immature ears of a heterotic F1 maize hybrid.
569 Plant Science 174 (1), 17 – 24.
- 570 Quinlan, R., 1993. C4.5: Programs for Machine Learning. Morgan Kaufmann
571 Publishers, San Mateo, CA.
- 572 Reif, J., Melchinger, A., Frisch, M., 2005. Genetical and Mathematical Prop-
573 erties of Similarity and Dissimilarity Coefficients Applied in Plant Breeding
574 and Seed Bank Management. Crop Sci 45, 1–7.
- 575 Rifkin, R., Klautau, A., 2004. In Defense of One-Vs-All Classification. Jour-
576 nal of Machine Learning Research 5, 101–141.
- 577 Springer, N. M., Stupar, R. M., 2007. Allelic variation and heterosis in maize:
578 How do two halves make more than a whole? Genome Res 17, 264–275.

- 579 Sumner, M., Frank, E., Hall, M. A., 2005. Proc 9th european conference
580 on principles and practice of knowledge discovery in databases,eeding up
581 logistic model tree induction. In: PKDD. pp. 675–683.
- 582 Terron, A., Preciado, E., Cordova, H., Mickelson, H., Lopez, R., 1997. De-
583 terminación del patrón heterótico de 30 líneas de maíz derivadas de la
584 población 43 SR del CIMMYT. Agron. Mesoamericana 8, 26–34.
- 585 The R Development Core Team, dic 2009. R: A language and environment
586 for statistical computing. reference index. <http://www.r-project.org/>.
- 587 Witten, I. H., Frank, E., 2005. Data Mining: Practical machine learning tools
588 and techniques, 2nd Edition. Morgan Kaufmann, San Francisco.
- 589 Xia, X. C., Reif, J. C., Hoisington, D. A., Melchinger, A. E., Frisch, M.,
590 Warburton, M. L., 2004. Genetic diversity among CIMMYT Maize Inbred
591 Lines Investigated with SSR Markers:I. Lowland Tropical Maize. Crop Sci
592 44, 2230–2237.
- 593 X.Su, Khoshgoftaar, T. M., Greiner, R., 2008. Using imputation techniques
594 to help learn accurate classifiers. Tools with Artificial Intelligence, IEEE
595 International Conference on 1, 437–444.

Table 1: Number of features preserved by Correlation-based Feature Selection (CFS). *Liu*, *Xia*, and *Morales* are the original molecular marker datasets. Full data denotes the initial number of features of each dataset. Min and Max are respectively the arithmetic means of the maximum and minimum number of features selected during the 30 Montecarlo runs of 10-Fold CV experiments.

	Dataset		
	Liu	Xia	Morales
Full data	188	166	42
Min	26	29	8
Max	50	42	20

Table 2: Means of the error rate and Kappa values in 30 Montecarlo runs of 10-Fold CV experiments. Native multiclass classifiers: Bayes Net (BN), Naive Bayes (NB), and Simple Logistic (SL). Multiclass extensions of Support Vector Machines: One Against All (OAA) and Error Correcting Output Coding (ECOC). Three molecular marker datasets, namely *Liu*, *Xia*, and *Morales*, are considered. Results on full and Correlation-based Feature Selection (CFS) reduced data are reported. Best results are shown in boldface.

Classifier	Full data						CFS reduced data					
	Liu		Xia		Morales		Liu		Xia		Morales	
	error	kappa	error	kappa	error	kappa	error	kappa	error	kappa	error	kappa
BN	0.205	0.749	0.475	0.368	0.715	0.039	0.280	0.658	0.428	0.455	0.755	-0.032
NB	0.345	0.685	0.472	0.372	0.751	0.000	0.294	0.638	0.432	0.439	0.772	-0.057
ECOC linear*	0.252	0.701	0.435	0.469	0.660	0.087	0.341	0.598	0.459	0.436	0.753	-0.039
ECOC rbf*	0.223	0.730	0.385	0.523	0.681	0.078	0.320	0.613	0.402	0.500	0.786	-0.078
OAA linear*	0.245	0.706	0.415	0.465	0.645	0.116	0.348	0.571	0.460	0.424	0.768	-0.059
OAA rbf*	0.223	0.730	0.429	0.442	0.690	0.043	0.357	0.579	0.462	0.433	0.819	-0.127
SL	0.345	0.576	0.436	0.433	0.572	0.210	0.367	0.552	0.537	0.326	0.703	0.033

* SVM with linear and radial basis function (rbf) kernels were implemented with defaults parameters of the Weka workbench (see Materials and Methods).

Table 3: Means of the error rate and Kappa values in 30 Montecarlo runs of 10-Fold CV experiments of *optimized SVM* with *linear* kernel under two decomposition schemes (OAA and ECOC).

Classifier	One Against All			Random Code		
	error	kappa	KS test (kappa) *	kappa	error	KS test (kappa) *
Morales Data	0.6308	0.1338	p-val = 0.1184	0.6500	0.1021	p-val = 0.3012
Xia Data	0.4160	0.4631	p-val = 0.5866	0.4438	0.4576	p-val = 0.9672
Liu Data	0.2302	0.7160	p-val = 0.9560	0.2330	0.7210	p-val = 0.9354

* Kolmogorov Smirnov test was performed between kappa values of classifier with default parameter (Table 2) and outputs of classifier with optimized parameters (this table) as stated in Materials and Methods.

Figure 1: *Liu* data. Boxplots of the Cohen’s Kappa coefficient in 30 Montecarlo runs of 10-Fold CV experiments. Native multiclass classifiers: Bayes Network (BN), Naive Bayes (NB), and Simple Logistic (SL). Multiclass extensions of Support Vector Machines (SVM): One Against All (OAA) and Error Correcting Output Coding (ECOC). Base classifiers: lin - SVM with linear kernel, rbf - SVM with radial basis function kernel. Results on full (Top) and Correlation-based Feature Selection (CFS) reduced data (Bottom) are shown.

Figure 2: *Xia* data. Boxplots of the Cohen’s Kappa coefficient in 30 Montecarlo runs of 10-Fold CV experiments. Native multiclass classifiers: Bayes Network (BN), Naive Bayes (NB), and Simple Logistic (SL). Multiclass extensions of Support Vector Machines (SVM): One Against All (OAA) and Error Correcting Output Coding (ECOC). Base classifiers: lin - SVM with linear kernel, rbf - SVM with radial basis function kernel. Results on full (Top) and Correlation-based Feature Selection (CFS) reduced data (Bottom) are shown.

Figure 3: *Morales* data. Boxplots of the Kappa coefficient in 30 Montecarlo runs of 10-Fold CV experiments. Native multiclass classifiers: Bayes Network (BN), Naive Bayes (NB), and Simple Logistic (SL). Multiclass extensions of Support Vector Machines (SVM): One Against All (OAA) and Error Correcting Output Coding (ECOC). Base classifiers: lin - SVM with linear kernel, rbf - SVM with radial basis function kernel. Results on full (Top) and Correlation-based Feature Selection (CFS) reduced data (Bottom) are shown.

Table 4: Means of the error rate and Kappa values in 30 Montecarlo runs of 10-Fold CV experiments of *optimized SVM* with *radial basis function* kernel under two decomposition schemes (OAA and ECOC).

Classifier	One Against All			Random Code		
	error	kappa	KS test (kappa) *	kappa	error	KS test (kappa)*
Morales Data	0.6795	0.0509	p-val = 0.0761	0.7556	-0.0410	p-val = 1.0000
Xia Data	0.4201	0.4550	p-val = 0.0357	0.3583	0.5540	p-val = 0.0327
Liu Data	0.2200	0.7350	p-val = 0.9030	0.2430	0.7500	p-val = 0.9350

* Kolmogorov Smirnov test was performed between kappa values of classifier with default parameter (Table 2) and outputs of classifier with optimized parameters (this table) as stated in Materials and Methods.

Figure 4: *Morales* data. Boxplots of the Kappa coefficient in 30 Montecarlo runs of 10-Fold CV experiments. Full and Relief Filtered data: Simple Logistic (SL). Multiclass extensions of Support Vector Machines (SVM): One Against All (OAA) and Error Correcting Output Coding (ECOC). Base classifiers: lin - SVM with linear kernel, rbf - SVM with radial basis function kernel. 42, 33, 21 and 12 indicates the number of attributes retained after filtering