ARTICLE IN PRESS

[m5G;August 31, 2016;13:48]

Pattern Recognition Letters xxx (2016) xxx-xxx

Contents lists available at ScienceDirect



Q1 Q2 Q3

Pattern Recognition Letters



journal homepage: www.elsevier.com/locate/patrec

Bayes covariant multi-class classification[☆]

Ondrej Šuch^{a,b,*}, **Santiago** Barreda^c

^a University of Žilina, Žilina, Slovakia

^b Mathematical Institute of Slovak Academy of Sciences, Banská Bystrica, Slovakia

^c Department of Linguistics, University of California, Davis, CA, United States

ARTICLE INFO

Article history: Received 9 December 2015 Available online xxx

Keywords: Multi-class classification Bradley—Terry model Bayes classifier Combining binary classifiers TIMIT Vowel classification

ABSTRACT

We consider multi-class classification models built from complete sets of pairwise binary classifiers. The Bradley–Terry model is often used to estimate posterior distributions in this setting. We introduce the notion of Bayes covariance, which holds if the multi-class classifier respects multiplicative group action on class priors. As a consequence, a Bayes covariant method yields the same result whether new priors are considered before or after combination of the individual classifiers, which has several practical advantages for systems with feedback. In the paper, we construct a Bayes covariant combining method and compare it with previously published methods in both Monte Carlo simulations as well as on a practical speech frame recognition task.

© 2016 Published by Elsevier B.V.

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

1 1. Introduction

Statistical and machine-learning classification methods have 2 found widespread applications in industry, as well as in scientific 3 research. Successful applications include optical character recog-4 nition [1], speech recognition systems [2], automated medical di-5 agnoses [3] and credit-risk scoring [4]. Although in some practi-6 7 cal applications binary decisions may be sufficient (e.g. cancer/no 8 cancer decision), most applications require correct classification 9 among multiple classes.

10 Broadly speaking, multi-class classification will pose a more 11 challenging problem than binary classification. One reason for this 12 is that the set of boundaries among multiple classes may be more complex and thus may be harder to learn than the boundary be-13 tween two classes. Another reason is that several powerful ma-14 chine learning methods for classification of two classes have no 15 direct analogues for multiple classes, making these methods inap-16 plicable for those faced with a multi-class problem. Important ex-17 amples of such methods include support vector machines [5,6] and 18 19 Adaboost [7].

There are many ways to reduce the multi-class classification of K classes to binary classification subproblems. One common approach is one-vs-all classification when one trains K classifiers to distinguish each class from all of the rest [8]. Another common

E-mail address: ondrej.such@fri.uniza.sk, ondrej.such@gmail.com (O. Šuch).

http://dx.doi.org/10.1016/j.patrec.2016.08.014 0167-8655/© 2016 Published by Elsevier B.V. approach is all-vs-all when one trains $\binom{K}{2}$ pairwise classifiers [9]. Other approaches have been proposed based on error correcting coding theory [10–14] and on training statistical meta-classifiers [15].

In our work, we consider the question of combining the output of binary classifiers in an all-vs-all setting. Some reasons to consider this approach rather than the one-vs-all approach [8] include:

- larger number of parameters allow for more powerful models,
- simpler and faster training of individual classifiers compared to one-vs-all ([8, pp. 123–124]),
- when samples are densely packed in Euclidean space, the allvs-all boundaries should be simpler, and thus easier to learn than one-vs-all boundaries; for example English vowels lie essentially in a 2-dimensional space [16],
- larger number of binary models allows for some tolerance of imprecision of individual classifiers ([17], [8, p. 102], [18]). Imprecise computation is typical for neuromorphic circuits for classification problems, which on the other hand are highly parallel and highly energy efficient [19,20].

Bayes theorem provides a rigorous foundation of classification. 44 The theorem explains the crucial role played by class priors on 45 the outcome of classification (cf. (2)). Usually, class priors are a 46 fixed quantity during classification. However, in multi-tiered sys-47 tems with feedback, one may desire to reevaluate evidence with 48 different priors based on feedback from other tiers. For instance, 49 a typical automated speech recognition system consists of three 50 parts - an acoustic model, a lexicon and a language model [21]. 51

^{*} This paper has been recommended for acceptance by Maria De Marsico.

^{*} Corresponding author at: University of Žilina, Žilina, Slovakia, Fax: +421 415134312.

JID: PATREC

ARTICLE IN PRESS

3

190

215

135 Solving $r'_{12} = r'_{23} = r'_{31}$ results in:

$$A^{3} = \frac{(\frac{1}{r_{23}} - 1)(\frac{1}{r_{31}} - 1)}{(\frac{1}{r_{12}} - 1)^{2}}, \qquad B^{3} = \frac{(\frac{1}{r_{31}} - 1)^{2}}{(\frac{1}{r_{12}} - 1)(\frac{1}{r_{23}} - 1)}.$$
 (15)

From (9) and (11) we conclude that the posterior of the Bayes covariant method *M* must be:

$$M(\mathbf{R}) = \frac{1}{1 + 1/A + 1/B} (1, \frac{1}{A}, \frac{1}{B}).$$
 (16)

Eqs. (15) and (16) define *M* uniquely. The resulting combining method is clearly 3-symmetric and it remains to check for Bayes covariance. Let $\mathbf{s} > 0$ be another reweighing vector. Then for $\mathbf{R}^{\mathbf{s}} =$ 141 (r_{ii}') we have:

$$\frac{s_2}{s_1} \left(\frac{1}{r_{12}} - 1 \right) = \frac{1}{r_{12}''} - 1 \tag{17}$$

$$\frac{s_1}{s_3} \left(\frac{1}{r_{31}} - 1 \right) = \frac{1}{r_{31}^{\prime\prime}} - 1 \tag{18}$$

143

142

$$\frac{s_3}{s_2} \left(\frac{1}{r_{23}} - 1 \right) = \frac{1}{r_{23}''} - 1.$$
(19)

144 It follows that:

$$A\frac{s_1}{s_2}\left(\frac{1}{r_{12}''}-1\right) = \frac{s_3}{s_1B}\left(\frac{1}{r_{31}''}-1\right) = \frac{s_2}{s_3}\frac{B}{A}\left(\frac{1}{r_{23}''}-1\right)$$
(20)

145 and thus:

$$M(\mathbf{R}^{\mathbf{s}}) \propto \left(1, \frac{s_2}{As_1}, \frac{s_3}{s_1B}\right) \propto \left(s_1, \frac{s_2}{A}, \frac{s_3}{B}\right).$$
(21)

This concludes verification of Bayes covariance and the proof of the theorem. \Box

148 5. A general Bayes covariant combining method

We will now construct Bayes covariant classifiers for cases with more than three categories (K > 3). Consider the moduli F of all feasible matrices. Inside F there is a submanifold B of feasible matrices for which (3) is a consistent system, which we shall call the *Bradley*-Terry manifold. Let us denote by P the point on the Bradley-Terry manifold corresponding to $r_{ij} = 1/2$ for $i \neq j$. We will consider only methods M for which

$$M(P) = \left(\frac{1}{K}, \frac{1}{K}, \ldots\right).$$
(22)

There is a natural action of the group **G** of reweighing vectors $\mathbf{q} > 0$ on *F* given by (6). Since the action is simply transitive on the Bradley–Terry manifold *B*, it follows from (9) that a Bayes covariant method is uniquely determined on the Bradley–Terry manifold.

Going back to proof of Theorem 1, we see that we took advantage of a set *S* of matrices satisfying (10). The set contained a single representative of each orbit under the action of reweighing vectors on *F* and by 3-symmetry we knew the exact value of the combining method on *S*.

In general, one may desire *S* to be a manifold of codimension K - 1 inside the variety of all feasible matrices and prescribe that

$$M(s) = \left(\frac{1}{K}, \frac{1}{K}, \ldots\right) \quad \text{for } s \in S.$$
(23)

167 If we express any feasible point $f \in F$ as $f = s^{\mathbf{q}}$ for reweighting vec-168 tor $\mathbf{q} > 0$ and $s \in S$, we will have from (9):

$$M(f) \propto \mathbf{q}.\tag{24}$$

Since the group of reweighing vectors \mathbf{G} acts transitively on the Bradley–Terry manifold, S has to have a single point intersection

with the Bradley–Terry manifold. To arrive at easily computable

expressions, we propose to take for *S* the set of points *Q* such that P is orthogonal to the tangent space of the Bradley–Terry manifold at the point *P*. This set is dependent on parameterization of the variety of feasible matrices, and may not contain a single representative of each orbit. Since we propose a linearly defined set *S*, it is natural to consider a parameterization in which the group action is linear. 172

Theorem 2. There exists a Bayes covariant combining method for every $K \ge 3$.

Proof. Consider parameterization of *F* given by $s_{ij} = \log(\frac{1}{r_{ij}} - 1)$ 181 for $1 \le i < j \le K$. From (6) we have for the action of $\mathbf{q} = 182$ (q_1, q_2, \dots, q_K) :

$$s_{ij}^{\mathbf{q}} = s_{ij} + \log q_j - \log q_i \tag{25}$$

Thus in parameterization by s_{ij} , the group of reweighing vectors 184 acts via translations. Note that since s_{ij} is a function of r_{ij} only, any 185 symmetry properties of pairwise matrix **R** are preserved in passing 186 to s_{ij} coordinates. Moreover, in s_{ij} coordinates: 187

- the set of feasible matrices is the $\binom{k}{2}$ dimensional real vector 188 space, 189
- point *P* is just the origin of the vector space.

i-1

For $i \leq K$ consider the one-parameter subgroup G_i of G 191 parameterized as: 192

$$\mathbf{q} = (\overbrace{1,1,\ldots,q}^{\sim},1,1,\ldots)$$
(26)

Let $h = (h_1, h_2)$ be a bijective mapping of the set $\{1, 2, ..., \binom{K}{2}\}$ 193 onto the set of pairs $\{(i, j) | 1 \le i < j \le K\}$. Such a function induces 194 ordering of coordinates s_{ij} , which will allow us to express tangent 195 vectors to the Bradley–Terry manifold at *P*. Namely, for $m \le \binom{K}{2}$ the 196 *m*th component of the tangent vector \mathbf{m}_k to *P* along the action of 197 \mathbf{G}_k is given by 198

$$(\mathbf{m}_k)_m = \begin{cases} -1 & \text{if } h_1(m) = k\\ 1 & \text{if } h_2(m) = k\\ 0 & \text{otherwise} \end{cases}$$
(27)

Subgroups \mathbf{G}_i generate \mathbf{G} and therefore vectors \mathbf{m}_k generate the 199 tangent space of the Bradley–Terry manifold at *P*. However vectors 200 \mathbf{m}_k are not linearly independent, because the action of \mathbf{q} and $c \cdot \mathbf{q}$ 201 is the same. Omitting one of the vectors, say \mathbf{m}_1 , from \mathbf{m}_i we arrive at an explicit basis $\mathbf{M} = (\mathbf{m}_2, \mathbf{m}_3, \dots, \mathbf{m}_K)$ of the tangent space of the Bradley–Terry manifold. 204

All that remains at this point is to solve the normal equations. 205 Let ${\bf N}$ be a basis of the orthogonal complement to the tangent 206 space, and let ${\bf s}$ represent a feasible matrix. We have: 207

$$\mathbf{s} = \mathbf{M}\mathbf{u} + \mathbf{N}\mathbf{v} \tag{28}$$

$$\mathbf{M}'\mathbf{s} = \mathbf{M}'\mathbf{M}\mathbf{u} + (\mathbf{M}'\mathbf{N})\mathbf{v} = \mathbf{M}'\mathbf{M}\mathbf{u}$$
(29)²⁰⁸

$$(\mathbf{M}'\mathbf{M})^{-1}\mathbf{M}'\mathbf{s} = \mathbf{u} \tag{30}$$

If we write $\mathbf{u}' = (u_1, \dots, u_{K-1})$, then by (25), (24) and (28) the proposed Bayes covariant combining method is given by: 211

$$\hat{p}_i = \frac{\exp(u_{i-1})}{\sum_{i=1}^{K} \exp(u_{i-1})}, \text{ where } u_0 = 0.$$
 (31)

Finally, note the special form of **M'M**. From (27) it follows that 212 it has (K - 1) on the diagonal and -1 elsewhere, which means it 213 equals $K \cdot \mathbf{I} - \mathbf{j}\mathbf{j}'$, where $\mathbf{j} = (1, 1, ...)'$. Therefore: 214

$$(\mathbf{M}'\mathbf{M})^{-1} = \frac{1}{\mathcal{K}}(\mathbf{I} + \mathbf{j}\mathbf{j}'), \tag{32}$$

allowing for simple computation of \mathbf{u} via (30).

JID: PATREC

ARTICLE IN PRESS

O. Such, S. Barreda/Pattern Recognition Letters xxx (2016) xxx-xxx

4

216 6. Comparison to previously published methods

Let us compare the Bayes covariant method described above to previously suggested methods for combining pairwise classifiers. Two kinds of combining methods have been suggested, both of them derived from the Bradley–Terry model.

The first group of methods is based on minimization of a functional, for which the location of the optimum coincides with the solution of (3) when the latter exists. This underlies the HT method suggested by Hastie and Tibshirani [17], who propose to minimize the Kullback–Leibler divergence between r_{ij} and $\hat{p}_i/(\hat{p}_i + \hat{p}_j)$. Two further methods were suggested by Wu et al. [18], who propose to minimize the quadratic forms:

$$\min_{\mathbf{p}} \sum_{i=1}^{k} \left[\sum_{j: j \neq i}^{k} (r_{ij} p_j - r_{ji} p_i) \right]^2, \text{ for the WLW1 method}$$
(33)

228

$$\min_{\mathbf{p}} \sum_{i=1}^{k} \sum_{j:j \neq i}^{k} (r_{ij}p_j - r_{ji}p_i)^2 \quad \text{for the WLW2 method}$$
(34)

The second group of methods is characterized by attempting to solve (3) directly. The PKPD method of Price et al. [26] solves a consistent subset of (3) to arrive at \hat{p}_i for each *i*, whereas the SBT method of Šuch et al. imposes a self-consistency condition on such estimates [27]. Let us abbreviate the Bayes covariant method described in Sections 3 and 4 as SB. All of these methods satisfy the 3-symmetry condition.

Theorem 3. None of the HT, WLW1, WLW2, PKPD and SBT methods is Bayes covariant.

Proof. We will use uniqueness of Bayes covariant classifier proved in Theorem 1. We can see from (15) that the ratio $A = \hat{p}_1/\hat{p}_2$ of the Bayes covariant classifier is not rational over the field generated by r_{ij} . Since all estimates of the methods except the HT method are rational over the field, it follows that the methods are not Bayes covariant.

To exclude the possibility that HT is Bayes covariant we numerically computed the Kullback–Leibler divergence for estimates given by both the HT and SB methods. When the matrix of pairwise likelihoods is [17, p. 452]:

$$\mathbf{R} = \begin{pmatrix} \cdot & 0.9 & 0.4 \\ 0.1 & \cdot & 0.7 \\ 0.6 & 0.3 & \cdot \end{pmatrix}$$
(35)

then the estimate of the posteriors made by HT method is $\hat{p}_{HT} \approx (0.481, 0.242, 0.277)$ and the Bayes covariant estimate is $\hat{p}_{SB} \approx (0.548, 0.192, 0.26)$. The Kullback–Leibler divergence is ≈ 0.376 for the former, whereas it is ≈ 0.397 for the latter. It follows from this that our method does not minimize the Kullback– Leibler divergence, and hence it is distinct from the HT method. \Box

254 To gain further insight into the relationship between the methods outlined here, we prepared multidimensional scaling (MDS) 255 256 plots. For each $K \leq 9$ we sampled 10⁵ values for r_{ij} uniformly from the $\binom{K}{2}$ -dimensional cube. Then we defined the distance be-257 tween methods as the proportion of samples for which the MAP 258 259 estimates of two methods differ. This distance is symmetric and respects the triangular inequality. We used classical MDS [28] to 260 obtain the canonical two dimensional representation as shown in 261 Fig. 1. As seen in the figure, the results for K = 3 turned out quite 262 differently from the others in that all the methods are arranged es-263 sentially on a single line. Moreover, the position of the PKPD and 264 SBT methods were indistinguishable. In looking at the rest of the 265

Table 1

Predicted posteriors for pairwise comparison matrix \mathbf{R} given by (36).

Combining	Predicted posteriors			
method	\hat{p}_1	\hat{p}_2	\hat{p}_3	
WLW2	0.499863	0.499862	0.000275	
SBT	0.499975	0.499750	0.000275	
PKPD	0.500067	0.499842	0.000091	
WLW1	0.500088	0.499638	0.000275	
HT	0.500312	0.499413	0.000275	
SB	0.682950	0.316902	0.000147	

panels in Fig. 1 we can see that the Bayes covariant method is closest to the method of Hastie and Tibshirani in its classification. 267

7. Practical considerations

Let us analyze the applicability of the Bayes covariant method 269 in practice. 270

First and foremost, by construction it gives the correct solution 271 of the Bradley-Terry equations whenever these equations have a 272 solution. This is true of all other methods (WLW1, WLW2, HT, SBT, 273 PKPD) and hence in this case all methods agree in predicted poste-274 riors as well as in MAP classification. Second, the implementation 275 of the Bayes covariant method described in Section 5 is straight-276 forward and the resulting method is fast. Its computation requires 277 computing $< K^2$ logarithms and exponentials and computation of 278 two matrix-vector multiplications. The only methods which are 279 simpler to implement are PKPD and the fast classification variant 280 of HT method [17, Theorem 1]. Methods SBT and WLW1 require 281 finding the eigenvector for eigenvalue $\lambda = 1$ of a Markov matrix, 282 which can be done quickly. The algorithm for WLW2 method is 283 based on a fast iterative scheme. The full HT method is the slow-284 285 est.

A weaker point of the Bayes covariant method is numerical 286 instability near the boundary of the subset of feasible matrices 287 within the set of nonnegative matrices satisfying $\mathbf{A} + \mathbf{A}' = \mathbf{j}\mathbf{j}' - \mathbf{I}$. 288 It may well happen in practice that a nondiagonal entry in ma-289 trix **R** approaches 0 or 1. In these cases the entry $s_{ij} = \log(\frac{1}{r_{ij}} - 1)$ 290 approaches infinity or negative infinity. This may cause instabil-291 ity in numerical computations, since a computer represents a real 292 number with a finite mantissa and a finite exponent. A simple 293 workaround is to eliminate from combining all classes C_i for which 294 there exists $j \neq i$ such that r_{ij} is zero, or r_{ij} is very small, e.g. below 295 10^{-5} . 296

Another important consideration is the required precision of 297 the binary predictions r_{ij} . Consider the following pairwise posterior matrix 299

$$\mathbf{R} = \begin{pmatrix} \cdot & 1/2 & 1 - 10^{-4} \\ 1/2 & \cdot & 1 - 10^{-3} \\ 10^{-4} & 10^{-3} & \cdot \end{pmatrix}$$
(36)

The third row contains very small entries, hence one may expect 300 that class 3 is quite unlikely to be predicted (it lost binary comparisons with classes 1 and 2 by large margins). Furthermore, the binary decision between classes 1 and 2 ended up in a tie. As a result, one would guess that final posteriors would be approximately (1/2, 1/2, 0). The results of predictions by all considered methods are shown in Table 1.

All methods give preference to class 1, but only the Bayes covariant method does so by a large margin. The way to understand the preference for class 1 is as follows. On the one hand, pairwise comparison between classes 1 and 2 ended in a tie so one would expect that posteriors for classes 1 and 2 are approximately the 311

268

ARTICLE IN PRESS

O. Šuch, S. Barreda/Pattern Recognition Letters xxx (2016) xxx-xxx



5





ARTICLE IN PRESS

6	
U	
_	

Table 2	
Merged	TIMI

Merged	TIMIT	classes	of	monophthongs	fol-
lowing	30].				

class	TIMIT	IPA	example
1	/iy/	/i/	beet
2	/uh/	/ʊ/	b oo k
3	/eh/	/ɛ/	bet
4	/ae/	/æ/	bat
5	/aa/	/α/	dark
5	/ao/	/ɔ/	all
	/ax/	/ə/	about
6	/ah/	$/\Lambda/$	b u t
	/ax-h/	/ə/	suspect
7	/er/	/3~/	bird
/	/axr/	/3~/	butter
8	/ix/	/1/	debit
	/ih/	/1/	b i t
0	/ux/	/ʉ/	toot
9	/uw/	/u/	boot /

312 same $\hat{p}_1 pprox$

$$_{1} \approx \hat{p}_{2}.$$
 (37)

On the other hand, we expect posterior \hat{p}_3 for the third class to be small. But Bradley–Terry equations require that

$$r_{31} = \frac{\hat{p}_3}{\hat{p}_3 + \hat{p}_1} \approx \frac{\hat{p}_3}{\hat{p}_1} \tag{38}$$

315

$$_{2} = \frac{\hat{p}_{3}}{\hat{p}_{3} + \hat{p}_{2}} \approx \frac{\hat{p}_{3}}{\hat{p}_{2}} \tag{39}$$

316 from which

 r_3

$$\frac{\delta_1}{\delta_2} \approx \frac{r_{32}}{r_{31}} = 10$$
 (40)

Any combining method has to resolve the tension between requirements (37) and (40). Compared to the other methods the Bayes covariant method puts much larger importance on relative ratios of even very small entries in matrix **R** and thus the resulting posteriors may differ significantly from other methods. As we will show in the following section, this may be a very desirable behavior.

323 8. Example

To gain insight into the practical performance of the Bayes covariant method we conducted experiments classifying speech frames from the benchmark TIMIT corpus [29]. Of particular interest are diphthongs in which the quality of sound rapidly changes among multiple categories. To that end we trained classifiers differentiating frames among 9 classes of monophthongs as described in Table 2.

Fig. 2 shows a representative result of such classification. The six panels in Fig. 2a present multi-class posteriors obtained by different combining methods, and the two panels in Fig. 2b posteriors obtained by two noncombining methods.

Qualitatively, all methods capture the expected acoustic dynamics of the /ai/ diphthong:

- 1. the vowel starts from the /aa,ao/ phoneme,
- 338 2. the middle part of the diphthong lies near /eh/ or /ae/,
- 339 3. the sound ends near the /iy/ phoneme.

However, the Bayes covariant method differs in the following aspects.

1. The posterior of the MAP class is often much higher. This can be seen mostly in /aa,ao/ segments and /iy/ segments. In the



Fig. 2. Posteriors of individual frames in the vowel /ai/ in the word-like' spoken by the speaker MJTC0 in the sentence SA2 from TIMIT [29]. In each plot the horizontal axis represents time in seconds; the grey line indicates p = 0.5.

404 405^{Q5}

407

408

409

410

411

414

415

417

418

419

420

421

422

423

424

425

426

427

428

429

430

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

460

465

466

467

468

469

470

471

472

473

474

475 476

479

- 344 former, the posterior is above 0.9 for the SB method, whereas 345 it drops below 0.5 early on for other methods. In the latter case, the posterior is near 1.0 for the SB method whereas it does not 346
- 347 exceed 0.8 for any of the other methods.
- 2. The overall posterior for |eh| + |ae| segments is similar to that 348 of the other methods. This indicates that the method does not 349 uniformly overshoot the posterior estimate. One can also ob-350 serve that although p(/ae/) + p(/eh) for the SB method closely 351 352 mirrors the other methods, the p(|eh|) component does not.
- 3. The posteriors oscillate less. Again, this can be seen both for 353 354 /aa,ao/ and /iy/ segments. It is even more impressive given that 355 the individual pairwise classifiers oscillate quite a bit as can be 356 inferred from the posterior predictions from the different com-357 bining methods.
- 4. The duration of gaps where none of the three groups (/aa,ao/, 358 /eh/ + /ae/ and /iy/) has posterior > 0.5 is much shorter for 359 the SB method. This is mostly due to shortening of "uncertainty 360 gap" between /eh/+/ae/ segments and /iy/ segments. 361

We can examine the last finding in more detail by analyzing 362 the pairwise likelihood plots shown in Fig. 2c. We can see that the 363 SB method gives a higher posterior (≈ 0.5) at the onset of /iy/ than 364 365 the pairwise LDA model between /eh/ and /iy/. The other methods 366 do quite the opposite, they indicate lower posterior of /iy/ at the onset of /iy/ compared to the pairwise LDA model. This means that 367 SB method extracted evidence for the presence of /iy/ from other 368 pairwise models. One model, which detects earlier onset of /iy/ is 369 370 the pairwise model classifying between /iy/ and /ae/. As we can see in Fig. 2c the model detects the onset of /iy/ earlier that /iy/-/eh/ 371 model. Overall, extracting information from multiple comparisons 372 r_{ii} , as demonstrated for the SB method, is a very desirable behavior 373 374 for a combining method.

375 Methods

We considered 9 classes of monophthongs as described in the 376 work of Lee and Hon [30]. For each class we extracted 10,000 sam-377 ples from the train section of TIMIT sentences of male speakers. 378 Sentences SA1 and SA2 were excluded. Then we trained 379

- a binary LDA model for each of the $\binom{9}{2} = 36$ pairs of classes, 380
- a multi-class LDA model, 381
- 382 • a multinomial model.

For training we used MASS and nnet R packages. The fea-383 384 ture set consisted of 256 values in the log-periodogram. The log-periodogram was obtained by taking a 512-point window (at 385 386 16 kHz sampling rate its duration was 32 ms), which was weighted 387 by Hanning window function, and then we took logarithms of 388 magnitudes of individual Fourier coefficients.

9. Conclusion 389

390 We have introduced the notion of Bayes covariance for multiclass combing methods. We have shown that a Bayes covariant 391 combining method exist for $K \ge 3$, and that for K = 3 there is a 392 unique method given a 3-symmetry condition. We have compared 393 the Bayes covariant method with five other combining methods. 394 395 Somewhat surprisingly, these combining methods lack Bayes co-396 variance, although they have the 3-symmetry property. The perfor-397 mance of the newly proposed method differs significantly from the other combining methods, both in Monte Carlo simulations as well 398 as in a speech frame classification task. Crucially, the speech frame 399 classification task demonstrated that the Bayes covariant method is 400 able to extract information from all pairwise classifiers to arrive at 401 more certain, less oscillatory posterior classification compared to 402 previously suggested methods. 403

Acknowledgments

This work was supported by Grants APVV-0219-12 and APVV-14-0560. 406

References

- [1] L. Vincent, Google book search: Document understanding on a massive scale, in: Proceedings of the Ninth International Conference on Document Analysis and Recognition ICDAR 2007, vol. 2, 2007, pp. 819-823, doi:10.1109/ICDAR. 2007.4377029
- [2] A. Graves, A.-R. Mohamed, G. Hinton, Speech recognition with deep recurrent 412 neural networks, in: Proceedings of the 2013 IEEE International Conference on 413 Acoustics, Speech and Signal Processing (ICASSP), 2013, pp. 6645-6649, doi:10. 1109/ICASSP.2013.6638947 416
- [3] P.J. Lisboa, A.F. Taktak, The use of artificial neural networks in decision support in cancer: A systematic review, Neural Netw. 19 (4) (2006) 408-415, doi:10. 1016/j.neunet.2005.10.007.
- [4] W.E.H.D.J. Hand, Statistical classification methods in consumer credit scoring: A review, J. R. Stat. Soc. Ser. A (Stat. Soc.) 160 (3) (1997) 523-541. URL http: www.jstor.org/stable/2983268
- [5] C. Cortes, V. Vapnik, Support-vector networks, Mach. Learn. 20 (3) (1995) 273-297, doi:10.1007/BF00994018.
 - V. Vapnik, The Nature of Statistical Learning Theory, Springer, 1996.
- [7] R. Schapire, The boosting approach to machine learning: An overview, in: D. Denison, M. Hansen, C. Holmes, B. Mallick, B. Yu (Eds.), Nonlinear Estimation and Classification, Lecture Notes in Statistics, vol. 171, Springer, New York, 2003, pp. 149-171, doi:10.1007/978-0-387-21579-2_9.
- [8] R. Rifkin, A. Klautau, In defense of one-vs-all classification, J. Mach. Learn. Res. 5 (2004) 101-141. URL http://dl.acm.org/citation.cfm?id=1005332.1005336
- [9] O. Lzoray, H. Cardot, Comparing combination rules of pairwise neural net-431 works classifiers, Neural Process. Lett. 27 (1) (2008) 43-56, doi:10.1007/ 432 s11063-007-9058-433 434
- [10] T. Dietterich, G. Bakiri, Solving multiclass learning problems via errorcorrecting output codes, J. Artif. Intell. Res. 2 (1995) 263-286, doi:10.1613/jair.
- [11] E.L. Allwein, R.E. Schapire, Y. Singer, Reducing multiclass to binary: A unifying approach for margin classifiers, J. Mach. Learn. Res. 1 (2001) 113-141, doi:10. 162/15324430152733133
- [12] K. Crammer, Y. Singer, On the learnability and design of output codes for multiclass problems, Mach. Learn. 47 (2-3) (2002) 201-233, doi:10.1023/A: 1013637720281.
- [13] J. Fürnkranz, Round robin classification, J. Mach. Learn. Res. 2 (2002) 721-747, doi:10.1162/153244302320884605
- [14] C.-W. Hsu, C.-J. Lin, A comparison of methods for multiclass support vector machines, IEEE Trans. Neural Netw. 13 (2) (2002) 415-425, doi:10.1109/72. 991427
- [15] Y. Shiraishi, K. Fukumizu, Statistical approaches to combining binary classifiers for multi-class classification, Neurocomputing 74 (5) (2011) 680-688, doi:10. 1016/j.neucom.2010.09.004.
- [16] D.J. Broad, H. Wakita, Piecewise-planar representation of vowel formant frequencies, J. Acoust. Soc. Am. 62 (6) (1977) 1467-1473, doi:10.1121/1.381676
- [17] T. Hastie, R. Tibshirani, Classification by pairwise coupling, Ann. Stat. 26 (2) (1998) 451-471, doi:10.1214/aos/1028144844.
- [18] T.-F. Wu, C.-J. Lin, R.C. Weng, Probability estimates for multi-class classification by pairwise coupling, J. Mach. Learn. Res. 5 (2004) 975-1005. URL http://dl. acm.org/citation.cfm?id=1005332.1016791
- [19] P. O'Connor, D. Neil, S.-C. Liu, T. Delbruck, M. Pfeiffer, Real-time classification 458 459 and sensor fusion with a spiking deep belief network, Front, Neurosci. 7 (2013) 178, doi:10.3389/fnins.2013.00178.
- [20] S. Hussain, S.-C. Liu, A. Basu, Improved margin multi-class classification us-461 ing dendritic neurons with morphological learning, in: Proceedings of the 462 2014 IEEE International Symposium on Circuits and Systems (ISCAS), 2014, 463 pp. 2640-2643, doi:10.1109/ISCAS.2014.6865715, 464
- [21] S. Renals, S. King, Automatic speech recognition, in: W.J. Hardcastle, J. Laver, F.E. Gibbon (Eds.), Handbook of Phonetic Sciences, Wilev-Blackwell, 2013, pp. 804-838.
- [22] T. Nearey, P. Assmann, M.-J. Solé, P.S. Beddor, M. Ohala, Probabilistic "sliding template" models for indirect vowel normalization, Experimental Approaches to Phonology, Oxford University Press, 2007, pp. 246-270.
- [23] T.M. Nearey, P.F. Assmann, Modeling the role of inherent spectral change in vowel identification, J. Acoust. Soc. Am. 80 (5) (1986) 1297-1308, doi:10.1121/ 1.394433
- [24] J.M. Hillenbrand, T.M. Nearey, Identification of resynthesized /hvd/ utterances: Effects of formant contour, J. Acoust. Soc. Am. 105 (6) (1999) 3509-3523, doi:10.1121/1.424676
- [25] R.A. Bradley, M.E. Terry, Rank analysis of incomplete block designs: The 477 478 method of paired comparisons, Biometrika 39 (3-4) (1952) 324-345, doi:10. 1093/biomet/39.3-4.324.
- [26] D. Price, S. Knerr, L. Personnaz, G. Drevfus, G. Tesauro, D. Touretzky, T. Leen, 480 Pairwise neural network classifiers with probabilistic outputs, Advances in 481 Neural Information Processing Systems, vol. 7, MIT Press, 1995, pp. 1109– 482 1116 483

JID: PATREC

O. Šuch, S. Barreda/Pattern Recognition Letters xxx (2016) xxx-xxx

- 8
- **Q6**484 485 [27] O. Such, S. Benus, A. Tinajová, A new method to combine probability estimates from pairwise binary classifiers, Information Technologies Applications and Theory, ITAT, 2015, pp. 194–199. URL http://ceur-ws.org/Vol-1422/ [28] J.C. Gower, Some distance properties of latent root and vector methods used in 486
- [29] J.S. Garofolo, L.F. Lamel, W.M. Fisher, J.G. Fiscus, D.S. Pallett, N.L. Dahlgren, 490 Darpa TIMIT acoustic phonetic continuous speech corpus CDROM, 1993, URL http://www.ldc.upenn.edu/Catalog/LDC93S1.html.
- [30] K.F. Lee, H.W. Hon, Speaker-independent phone recognition using hidden Markov models, IEEE Trans. Acoust. Speech Signal Process 37 (11) (1989) 1641– 1648, doi:10.1109/29.46546.

487 488 multivariate analysis, Biometrika 53 (3-4) (1966) 325-338, doi:10.1093/biomet/ 489 53.3-4.325.