

# Fusing Face Recognition Algorithms and Humans

Alice J. O'TOOLE, Hervé ABDI, Fang JIANG, and, P. Jonathon PHILLIPS\*

**Abstract**—It has been demonstrated recently that state-of-the-art face recognition algorithms can surpass human accuracy at matching faces over changes in illumination. The ranking of algorithms and humans by accuracy, however, does not provide information about whether algorithms and humans perform the task comparably or whether algorithms and humans can be fused to improve performance. Here, we fused humans and algorithms using partial least squares (PLS) regression. In the first experiment, we applied PLS to face-pair similarity scores generated by seven algorithms participating in the Face Recognition Grand Challenge (FRGC). PLS produced an optimal weighting of the similarity scores, which we tested for generality with a jackknife procedure. Fusing algorithms' similarity scores using the optimal weights produced a two-fold reduction of error rate over the most accurate algorithm. Next, human subject-generated similarity scores were added to the PLS analysis. Fusing humans and algorithms increased performance to near-perfect classification accuracy. These results are discussed in terms of maximizing face recognition accuracy with hybrid systems consisting of multiple algorithms and humans.

**Index Terms**— face and gesture recognition, performance evaluation of algorithms and systems, human information processing

## I. INTRODUCTION

THE field of automatic face recognition algorithms has expanded in the past decade from consisting of simple algorithms that operate on highly controlled images of faces to more sophisticated algorithms aimed at operating in the natural conditions that characterize most security applications. One particularly difficult challenge in advancing algorithms from controlled environments to natural environments has

been the problem of operating over substantial changes in illumination. The computational difficulties posed by the *illumination problem* have been well documented in the automatic face recognition [cf. 1-3] and human perception literatures [4-6].

In more practical terms, the performance of face recognition algorithms in controlled and uncontrolled illumination environments was assessed recently in the Face Recognition Grand Challenge (FRGC), a U.S. Government-sponsored test of face recognition algorithms, with the goal of fostering algorithm development [7,8]. The FRGC (2004-2006) included academic, industrial, and research lab competitors. Competitors participated in the program by volunteering to have their algorithms tested on one or more of six face matching experiments, varying in difficulty. The set of experiments included both a controlled illumination face matching experiment and a more difficult experiment where algorithms matched face identity in images taken under different illumination conditions. Because the FRGC tested multiple algorithms simultaneously using a standardized evaluation protocol and a common image set, it provides a useful, time-locked look at the performance of state-of-the-art face recognition algorithms.

The difficulty of the illumination problem can be seen clearly by comparing the performance of algorithms in the controlled and uncontrolled illumination experiments of the FRGC. In both cases, the task of the algorithms was to decide for each of a large number of face pairs ( $> 128$  million), whether the images were of the same person or of different people. In the controlled illumination experiment, the illumination conditions were the same for both images in the pair. In the uncontrolled illumination experiment, one image was taken under controlled illumination conditions and the other was taken under uncontrolled illumination (see Fig. 1 for an example image pair).

Twenty algorithms competed in the controlled illumination experiment and achieved an average verification rate of .91 at the .001 false accept rate. By contrast, in the uncontrolled illumination experiment, only 7 algorithms participated, achieving an average verification rate of .41 at the .001 false accept rate. The difference in participant numbers and average performance in these experiments is evidence that the illumination problem continues to challenge face recognition algorithms.

Manuscript received May 25, 2006. This work was supported by a contract to A. J. O'Toole and H. Abdi from TSWG. P. J. Phillips was supported in part by funding from the National Institute of Justice.

A. J. O'Toole is at the School of Behavioral and Brain Sciences, GR4.1 The University of Texas at Dallas Richardson, TX 75083-0688, USA, (e-mail: [otoole@utdallas.edu](mailto:otoole@utdallas.edu)).

Hervé Abdi is at the School of Behavioral and Brain Sciences, GR4.1 The University of Texas at Dallas Richardson, TX 75083-0688, USA, (e-mail: [herve@utdallas.edu](mailto:herve@utdallas.edu)).

Fang Jiang is at the School of Behavioral and Brain Sciences, GR4.1 The University of Texas at Dallas Richardson, TX 75083-0688, USA, (e-mail: [fxj018100@utdallas.edu](mailto:fxj018100@utdallas.edu)).

P.J. Phillips is with the National Institute of Standards and Technology, 100 Bureau Dr., MS 8940 Gaithersburg MD 20899 [jonathon@nist.gov](mailto:jonathon@nist.gov)

A rather different perspective on the relatively poor performance of algorithms in the uncontrolled illumination experiment comes from comparing the algorithms to humans performing a comparable task. In a recent study [9], human face matching performance was compared to the performance of the seven algorithms participating in the uncontrolled illumination matching experiment of the FRGC. We describe this previous study in some detail here, because it forms the base of the present work.

Algorithms in the FRGC uncontrolled illumination experiment, (Experiment 4 in FRGC nomenclature), matched face identities in all possible pairs of 16028 *target* images and 8014 *probe* images, with target images taken under controlled illumination conditions and probe images taken under uncontrolled illumination conditions (see Fig. 1 for a sample pair). The output for each algorithm was a matrix of similarity scores for all possible pairs of faces. For each algorithm, a Receiver Operating Characteristic (ROC) curve was generated from the similarity score matrix. The performance of the seven algorithms was compared using these ROC curves [cf. 9 for complete results].



**Fig. 1** A sample pair of face images from a “match” trial (a) Participants responded by rating the likelihood that the pictures were of the same person using a five-point scale ranging from “1.) sure they are the same person” to “5.) sure they are not the same.

The primary difficulty in comparing the performance of humans to algorithms in the FRGC is the implausibly large number of face pair comparisons required for an exhaustive comparison. Therefore, to compare the performance of humans to algorithms, face pairs were sampled from the matrix by selecting a set of the easiest and most difficult pairs [9]. In the present work, we concentrate on the most difficult image pairs. In both cases, however, the sampling was done with the help of a control algorithm based on a principal components analysis (PCA) of the aligned and scaled face images. Using this algorithm, *easy match pairs* were defined based on similarity scores that were substantially greater than the mean for the distribution of matched face pairs, i.e., highly similar images of the same person. *Difficult match pairs* were those with similarity scores substantially lower than the match mean, i.e., highly dissimilar images of the same person. Easy

and difficult non-match pairs were defined inversely.

Human subjects matched the identity of 240 sample face pairs, by rating their certainty that the pairs were of the same person. Human responses ranged on a five-point scale from “certain the two images are of the same person” to “certain that two images are not of same person”. The rating data allowed for the generation of an ROC curve for human performance that was comparable to the ROC curves derivable from the performance of the algorithms.

The human-machine comparison was conducted by extracting the algorithms’ similarity scores for same face pairs tested in the human face matching experiment. These were plotted on ROC curves along with human match accuracy data [9]. The results demonstrated clearly that three algorithms [10-12] surpassed human performance on the difficult face pairs. Of these, the algorithm from The New Jersey Institute of Technology [10] and the algorithm from Carnegie Mellon University [11] have been published. Details on the third algorithm, from the Viisage Corporation<sup>1</sup>, are only partially available [12].

In addition to the finding that three algorithms were competitive with humans on the difficult pairs of faces, all but one algorithm surpassed human performance on the easy face pairs. Combined, these findings suggest that although algorithm performance on the uncontrolled illumination experiment in the FRGC may be poor in absolute terms, it is nonetheless competitive with human performance. This comparison is of interest due to the fact that humans are currently performing this task in most applied situations.

This previous study forms the base of the present work. Although the quantitative ranking of human performance relative to a set of algorithms provides a useful benchmark, this ranking does not offer any insight into whether algorithms recognize faces in ways that are similar to humans. The FRGC showed that none of the algorithms performed face recognition in uncontrolled illumination environments perfectly. Our previous work showed the same result for humans. If algorithms and humans take diverse approaches to the problem of face matching, it is possible that an appropriate fusion of algorithms and humans can yield better performance than a single algorithm or the fusion of multiple algorithms. Indeed, previous work has shown that fusing the multiple face recognition algorithms improves performance over a single algorithm [cf. 13-15]. However, no previous studies have fused human and algorithm performance.

In the majority of applications for face recognition, a human operator is present and involved in the decision process. Thus, it is natural to optimize system performance by explicitly incorporating human face recognition capabilities into the decision process. Towards this end, we present a methodology for fusing algorithm and human performance.

<sup>1</sup> See Acknowledgement.

In this study, we asked two questions. First, can performance be improved by fusing algorithms from the FRGC uncontrolled illumination experiment? Second, does fusing humans and algorithms improve performance above the level achieved by the algorithm fusion? The availability of multiple algorithm estimates of face similarity, in conjunction with analogous human estimates of similarity, offers the possibility of exploring these questions in a more systematic way than is generally possible. Here, we investigated the possibility of fusing face similarity estimates from algorithms and humans to improve face-matching performance.

Fusion was performed by partial least squares (PLS) regression, a statistical technique that generalizes and combines features from principal component analysis and multiple regression [16,17]. The technique is used to predict a set of dependent variables from a set of independent variables (predictors). Though less known in the pattern recognition literature, PLS is widely used in chemometrics, sensory evaluation, and neuroimaging data analysis [cf. 16,18,19,21].

In the present work, algorithm and human estimates of face similarity were the predictors and the match status of individual face pairs (i.e., same person or different people) was the dependent variable. PLS regression gives a set of orthogonal factors, sometimes called latent vectors  $\{t_1 \dots t_l\}$ , from the covariance matrix of predictors and dependent variables. These can be used to predict the dependent variable(s), by appropriately weighting the predictors. This set of weights is called  $\mathbf{B}_{pls}$  in the PLS-regression literature [16]. To fuse algorithms, the weights prescribed in the latent vector(s) are used to combine the similarity scores from each of the seven algorithms to produce an estimate of the match status for the face pairs. When fusing humans and algorithms, there are eight predictors, seven from the algorithms and one from the averaged human data.

The predictive power of these factors is generally assessed with cross-validation techniques such as a bootstrap or jackknife procedure. All factors, or only a subset of them, can be used to compute the prediction of the dependent variable(s), which are obtained as a weighted combination of the original predictors given by  $\mathbf{B}_{pls}$ . The larger the number of factors kept, the better the prediction of the “learning set” but, in general, a smaller number of factors is optimal for robust prediction (i.e., for test set predictions).

In the first experiment, we applied PLS to the similarity scores generated by seven algorithms that participated in the FRGC uncontrolled illumination experiment. We tested the generality of the optimal weights found in the analysis for predicting face match status using a jackknife procedure. In the second experiment, we added human-generated similarity scores to the algorithms’ scores and measured the contribution human estimates make to the fusion.

## II. PROCEDURE

### A. Stimuli

Face stimuli were chosen from a large database developed for the FRGC study [7,8]. The uncontrolled illumination probe faces had a resolution of 2272 x 1704 pixels. The controlled illumination target faces had a resolution of 1704 x 2272 pixels. For the present analyses, we used the same set of difficult face pairs sampled for the previous quantitative comparison between humans and algorithms [9]. These were sampled from the 128,448,392 pairs available, which included 407,352 (0.32%) match pairs (i.e., image pairs of the same person, and 128,041,040 (99.68%) no-match pairs (i.e., image pairs of different people). To eliminate the possibility that humans could base identity comparisons on surface facial characteristics associated with race or age, all images in the study were of faces of Caucasian males and females in their twenties. All pairs were matched by sex.

In the present study, only “difficult face pairs” were included. These were chosen using a control algorithm based on principal components analysis (PCA) of the aligned and scaled images. Specifically, difficult match face pairs ( $n = 60$ ) were sampled randomly from match pairs that had similarity scores less than 2 standard deviations below the match mean. Difficult non-match face pairs ( $n = 60$ ) were sampled randomly from non-match pairs that had similarity scores greater than 2 standard deviations above the non-match mean.

The validity of PCA as a pre-screening algorithm for humans and algorithms was supported in the previous study [9]. The PCA algorithm reliably predicted “easy” and “difficult” sets of face pairs for humans in three experiments [9]. All seven algorithms were likewise more accurate on the PCA-screened easy face pairs than on the PCA-screened difficult faces [9]. The PCA, therefore, can serve as a useful sampling tool, even though it is not considered “state-of-the-art”.

### B. Human Subject Judgments of Face Similarity

The human subject data for this experiment were collected in an experiment in which subjects viewed the image pairs and rated the likelihood that the images were of the same person or of different people [9]. For completeness, we sketch out the methods used in that study. Forty-nine subjects (25 male and 24 female) viewed the 120 pairs of faces for 2 seconds each and responded by rating each pair on the following scale: 1.) sure the pictures are of the same person; 2.) think the pictures are of the same person; 3.) don’t know; 4.) think the pictures are not of the same person; 5.) sure the pictures are not of the same person. Of the 120 pairs, half were match pairs and half were non-match pairs. Equal numbers of male and female pairs were included in the match and non-match conditions.

For each pair of faces, the average rating, was computed across the 49 subjects. This average served as the human similarity score for that pair of faces in the PLS regression.

### C. Algorithms' Judgments of Face Similarity

The similarity scores of the 120 difficult face pairs presented to participants in the human experiment were extracted from each algorithm's 16028 x 8014 similarity matrix. These scores served as the algorithm data for the PLS regression.

## III. RESULTS

### A. Experiment 1 – Algorithm Fusion by PLS

The similarity scores for the seven algorithms for the 120 difficult face pairs (60 match and 60 non-match) were combined in a column-wise matrix. The dependent variable was a 120-element vector containing the match status (+1 for match and -1 for non-match) for each face pair. PLS was applied simultaneously to the similarity scores and match status matrices.

We varied the number of PLS factors retained from 1 to 5 and found a three-factor solution to be optimal. Retaining three factors indicates that the first three latent vectors, ordered according to the proportion of variance explained in the covariance matrix, are combined linearly to specify the weights for combining the similarity scores.

A robust performance estimate was determined with a jackknife simulation. We started with the 120 face pairs available, and systematically deleted each face pair in turn, re-computing the PLS with the remaining 119 pairs of faces. We tested the match status predictions for the PLS solutions derived from 119 pairs of faces on the “left-out” face pair. This yielded 120 generalized match prediction tests. The error rate we report is the fraction of left-out face pairs incorrectly classified according to match status.

TABLE I  
WEIGHT MATRIX FOR ALGORITHM FUSION DIFFICULT FACE PAIRS

	<i>NJIT</i>	<i>Viisage</i>	<i>CMU</i>	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>
<i>weights</i>	-2.2	-1.81	-.05	.00	-.15	.00	.16
<i>error rates</i>	.12	.20	.14	.37	.23	.25	.36

**Table I.** The table shows weight vector for combining algorithm similarity estimates for optimal match performance in the top row. Large absolute values indicate the most useful predictors, which in this case are the NJIT algorithm [10] and Viisage [12]. The bottom row gives the proportion of classification errors for the algorithms *individually*. The fusion cuts the best algorithm's error rate by a factor of two, from .12 to .059.

Error rates for classification with one through five factors were .067, .075, .059, .067, and .083, respectively. These error rates are all lower than the minimum error rate achieved by any single algorithm operating alone (cf. Table I for error rates for each individual algorithm). Specifically, the data indicate that fusion, following the optimal weighting derived with the PLS regression, cut the error rate of the best performing algorithm (NJIT [10] with a .12 error rate) by a factor of two.

For purposes of interpretation, the weights for combining

similarity scores appear in Table I. These weights are used to combine similarity scores from the seven algorithms to achieve a maximal separation between the match and non-match face pair distributions. Algorithms with weights that have large absolute values are the most useful for improving performance with fusion.

TABLE II  
WEIGHT MATRIX FOR HUMAN-ALGORITHM FUSION

	<i>Human</i>	<i>NJIT</i>	<i>Viisage</i>	<i>CMU</i>	<i>B</i>	<i>D</i>
<i>weights</i>	.47	-1.29	-.71	-.03	-.12	.20

**Table II.** The table shows weight vector for combining human and algorithm similarity estimates for optimal match performance. Algorithms A and C had weights of zero and are not included in the table. The addition of humans to the PLS decreased the error rate from .059, for algorithm fusion, to .008 for human-algorithm fusion.

Using this as an interpretation guide, it is clear that most of the improvement in accuracy comes from combining just two algorithms, NJIT [10] and Viisage [12], whose weights have the largest absolute values. This might be due to these algorithms having maximally diverse strategies for computing face similarity. This interpretation that seems likely given that the CMU algorithm [11] performed somewhat better than the algorithm of Viisage [12]. Thus, more benefit can be derived from combining lesser-performing algorithms that operate in different fashions than by combining higher performing similar algorithms.

### B. Experiment 2 – Human and Algorithm Fusion by PLS

Can fusing humans and algorithms add to the accuracy of the match estimates and further improve classification over that obtained with the fused algorithms? In this experiment, we added human similarity estimates to the PLS model. The analysis proceeded as before but with a column vector containing the averaged human similarity data appended to the predictor matrix.

Again, we varied the number of PLS factors we retained from 1 to 5. In this case, we found a two-factor solution to be most robust, using the jackknife procedure described previously. The weights for combining human and algorithm similarity estimates appear in Table II. Performance with one factor through five factors yielded classification error rates of .042, .008, .033, .033 and .042, respectively.

These results illustrate that it is possible to obtain nearly perfect classification, when humans are added into the predictor matrix. This suggests that human strategies for assigning similarities to faces add usefully to those employed by the best algorithms.

It is worth noting from previous work [9], that the accuracy of humans was found to be below that of NJIT [10], CMU [11], and Viisage [12], but above the accuracy of algorithms A, B, C, and D. In that study, similarity ratings from individual



subjects were collapsed across the 120 face pairs to create an ROC curve for each subject. These individual ROC curves were then averaged to give an overall estimate of human accuracy. Here, we averaged the similarity ratings for 120 face pairs, collapsing across the individual subjects. Interestingly, though perhaps not surprisingly, we found that by averaging across the 49 human subjects' estimates of face similarity for each face pair individually, human error rate was .12, comparable to NJIT, the best algorithm. This suggests that individual subjects, like algorithms, may employ diverse strategies for judging the similarity of face pairs. By consequence, combining the similarity estimates of individual subjects by fusion could likewise benefit accuracy.

#### IV. DISCUSSION

Fusing humans and algorithms is a reasonable goal for face recognition researchers and corporations with hopes of applying their systems to real applications. Knowing how accurately algorithms and humans are by themselves is a start in trying to estimate how well combinations of algorithms and humans will work. But, quantitative measures of accuracy for individual algorithms and humans are not sufficient for guiding the development of hybrid systems. The present study illustrates that the most useful fusions of algorithms and humans is likely to come from combining face recognition systems (algorithms or humans) with diverse face recognition strategies.

In this paper, we demonstrated that fusing algorithms and humans substantially improved performance on a difficult face-matching task. The use of PLS regression to fuse the algorithms and humans also yielded a precise indication of how to combine the individual components of the fusion optimally. This weight vector serves simultaneously as a recipe for fusing systems and as an indicator of the similarity of algorithm and human strategies for making similarity judgments.

Given that neither algorithms nor humans perform face recognition perfectly in uncontrolled environments and that a majority of applications have a human operator in the loop, a reasonable goal of researchers should be to design face recognition strategies that optimally combine algorithms and humans. Fusion of algorithms and humans to create good hybrids can, therefore, be a useful and practical approach to improving face matching performance in important applications.

#### REFERENCES

- [1] R. Gross, S. Baker, I. Matthews, and T. Kanade, "Face recognition across pose and illumination," in *Handbook of Face Recognition*, S. Z. Li and A.K. Jain, Eds. Springer, pp. 193-216, 2005.
- [2] P.J. Phillips, H. Moon, P. Rizvi, and P. Rauss, "The FERET evaluation method for face recognition algorithms," *IEEE Transactions on Pattern Analysis and Machine Learning*, vol. 22, pp. 1090-1104, 2000.
- [3] P.J. Phillips, P. Grother, R. Micheals, D. Blackburn, E. Tabassi, and J.M. Bone, "FRVT 2002 evaluation report," Tech. Rep. NISTIR 6965 <http://www.frvt.org>, 2003.
- [4] W.J. Braje, "Illumination encoding in face recognition: effect of position shift," *Journal of Vision*, vol. 3, pp. 161-170, 2003.
- [5] W.J. Braje, D. Kersten, M.J. Tarr, and N.F. Troje, "Illumination effects in face recognition," *Psychobiology*, vol. 26, pp. 371-380, 1999.
- [6] W.J. Braje, G.E. Legge, and D. Kersten, "Invariant recognition of natural objects in the presence of shadows," *Perception*, vol. 29, pp. 383-398, 2000.
- [7] P.J. Phillips, P.J. Flynn, T. Scruggs, K.W. Bowyer, J. Chang, K. Hoffman, J. Marques, J. Min, and W. Worek, "Overview of the face recognition grand challenge," *Proc. IEEE Computer Vision & Pattern Recognition*, vol. 1, pp. 947-954, 2005.
- [8] P.J. Phillips, P. J. Flynn, T. Scruggs, K. W. Bowyer, W. Worek, "Preliminary Face Recognition Grand Challenge Results," in *Proceedings of the Seventh International Conference on Automatic Face and Gesture Recognition*, 15-24, 2006.
- [9] A. O'Toole, P.J. Phillips, F. Jiang, J. Ayyad, N. Pénard, and H. Abdi, "Face recognition algorithms surpass humans," Tech. Rep. NISTIR, <http://face.nist.gov>, 2005.
- [10] C. Liu, "Capitalize on dimensionality increasing techniques from improving face recognition Grand Challenge performance," *IEEE Transactions on Pattern Analysis and Machine Learning*, 2006, in press.
- [11] C.M. Xie, M. Savvides, and V. Kumar, "Kernel correlation filter based redundant class-dependence feature analysis (KCFA) on FRGC2.0 Data," *IEEE International Workshop Analysis & Modeling Faces & Gestures*, pp. 32-43, 2005.
- [12] M. Husken, B. Brauckmann, S. Gehlen, and C. von der Malsburg, "Strategies and benefits of fusion of 2D and 3D face recognition," *Proceedings of the 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, vol. 3, pp. 174, 2005.
- [13] P. Grother, "Face Recognition Vendor Test 2002 Supplemental Report," Tech. Rep. NISTIR 7083 <http://www.frvt.org>, 2004.
- [14] O. Melnik, Y. Vardi, C.-H. Zhang, "Mixed Group Ranks: Preference and Confidence in Classifier Combination," *IEEE Transactions on Pattern Analysis and Machine Learning*, vol. 26, pp. 973-981, 2004.
- [15] J. Czyz, J. Kittler, L. Vanderdorpe, "Combining face verification experts," in *Proceedings 16th International Conference on Pattern Recognition II* pp. 28-31, 2002.
- [16] H. Abdi Partial least squares regression (PLS-regression). In M. Lewis Beck, A. Bryman, T. Futing (Eds.) *Encyclopedia for Research Methods for the Social Sciences*. Thousand Oaks: CA, Sage, 2003. pp. 792-795.
- [17] T. Naes, T. Isaksson, T. Fearn and Davis T. *Multivariate calibration and classification*. Chisester (UK) NIR Publications. 2004.
- [18] H. Martens, and M. Martens. *Multivariate analysis of quality*. London: J. Wiley. 2001.
- [19] A.R. McIntosh, and N. Lobaugh. Partial least squares analysis of neuroimaging data: applications and advances. *Neuroimage*, vol. 23, pp. 250-263. 2004.
- [20] H. Abdi Partial least squares regression. In N. J. Salkind (Ed.) *Encyclopedia of measurement and statistics*. Thousand Oaks: CA, Sage, 2007. pp. 740-744.
- [16] H. Abdi Multivariate analysis. In M. Lewis Beck, A. Bryman, T. Futing (Eds.) *Encyclopedia for Research Methods for the Social Sciences*. Thousand Oaks: CA, Sage, 2003. pp. 699-702.

#### ACKNOWLEDGMENT

This work was supported by a contract to A. O'Toole and H. Abdi from TSWG. P. J. Phillips was supported in part by funding from the National Institute of Justice. This work was performed for the Department of Justice in accordance with section 303 of the Border Security Act, codified as 8 U.S.C. 1732. Specific hardware and software products identified in this report were used in order to perform the evaluations described in this document. In no case does such identification imply recommendation or endorsement by the National Institute of Standards and Technology, nor does it imply that the products and equipment identified are necessarily the best available for the purpose. The primary goal of the FRGC is to encourage and facilitate the development of face recognition algorithms. To provide the face recognition research community with an unbiased assessment of state-of-the-art algorithms, research groups voluntarily submit similarity scores from prototyped experiments to the National Institute of Standards and Technology (NIST) for analysis. The

results of the analysis by NIST are anonymous, unless otherwise agreed to by the participating algorithm developers. All participating groups were given the choice of remaining anonymous or being identified in this report. Performance results are from Jan. 2005 for all algorithms except Xie et al., 2005, where results are from Aug., 2005.

#### A. Appendix: PLS Regression

In this appendix we give a brief description of Partial Least Square Regression (PLSR), more complete presentations can be found in [16,20]. MATLAB programs can be downloaded from [www.utdallas.edu/~herve](http://www.utdallas.edu/~herve). PLSR generalizes and combine features from principal component analysis and multiple regression. Its goal is to optimally predict a set of dependent variables from a set of predictors. Specifically, PLSR searches for a set of components (called *latent vectors*) that performs a simultaneous decomposition of  $\mathbf{X}$  and  $\mathbf{Y}$  with the constraint that these components explain as much as possible of the *covariance* between  $\mathbf{X}$  and  $\mathbf{Y}$ . This step is followed by a regression step where the decomposition of  $\mathbf{X}$  is used to predict  $\mathbf{Y}$ .

##### 1) Notations

The  $I$  observations described by  $K$  dependent variables are stored in an  $I \times K$  matrix denoted  $\mathbf{Y}$ , the  $I \times J$  matrix of predictors is denoted  $\mathbf{X}$ . Without loss of generality, both  $\mathbf{X}$  and  $\mathbf{Y}$  are supposed to be centered and normalized. The common set of (orthogonal) latent vectors is stored in the  $I \times L$  matrix  $\mathbf{T}$  (i.e.,  $\mathbf{T}^T \mathbf{T} = \mathbf{I}$ ). PLSR decomposes  $\mathbf{X}$  as

$$\mathbf{X} = \mathbf{T} \mathbf{P}^T$$

where  $\mathbf{P}$  is a  $J \times L$  matrix called the  $\mathbf{X}$ -loadings matrix. The matrix  $\mathbf{Y}$  is estimated as

$$\mathbf{Y} = \mathbf{T} \mathbf{B} \mathbf{C}^T$$

where  $\mathbf{B}$  is a diagonal matrix with the “regression weights” as diagonal elements and  $\mathbf{C}$  is the “weight matrix” of the dependent variables

##### 2) Computations of the Latent vectors, loadings and weights

A latent vector is obtained by finding two sets of weights  $\mathbf{w}$  and  $\mathbf{c}$  in order to create (respectively) a linear combination of the columns of  $\mathbf{X}$  and  $\mathbf{Y}$  such that their covariance is maximum. Specifically, the goal is to obtain a first pair of vectors

$$\mathbf{t} = \mathbf{X} \mathbf{w} \text{ and } \mathbf{u} = \mathbf{Y} \mathbf{c} \quad (1)$$

under the constraints that

$$\mathbf{w}^T \mathbf{w} = 1, \mathbf{t}^T \mathbf{t} = 1 \text{ and } \mathbf{t}^T \mathbf{u} \text{ be maximal.} \quad (2)$$

When the first latent vector has been found, it is *subtracted* from both  $\mathbf{X}$  and  $\mathbf{Y}$  and the procedure is re-iterated until  $\mathbf{X}$  becomes a null matrix (see the algorithm section for more).

##### 3) Algorithm

The different component of PLSR can be found by a series of singular value decomposition followed by a deflation. Specifically, the first weight vectors  $\mathbf{w}$  and  $\mathbf{c}$  are respectively the first right and left singular vector of the matrix  $\mathbf{X}^T \mathbf{Y}$ . Vectors  $\mathbf{t}$  and  $\mathbf{u}$  are then derived using Equation 1. With these vectors, the value of  $b$  is computed as  $b = \mathbf{t}^T \mathbf{u}$  and then used to predict  $\mathbf{Y}$  from  $\mathbf{t}$  as  $\mathbf{Y} = b \mathbf{t} \mathbf{c}^T$ . The factor loadings for  $\mathbf{X}$  are computed as  $\mathbf{p} = \mathbf{X} \mathbf{t}$ . Now subtract (i.e., partial out) the effect of  $\mathbf{t}$  from both  $\mathbf{X}$  and  $\mathbf{Y}$  as follows  $\mathbf{X} = \mathbf{X} - \mathbf{t} \mathbf{p}^T$  and  $\mathbf{Y} = \mathbf{Y} - b \mathbf{t} \mathbf{c}^T$ . The vectors  $\mathbf{t}$ ,  $\mathbf{u}$ ,  $\mathbf{w}$ ,  $\mathbf{c}$ , and  $\mathbf{p}$  are then stored in the corresponding matrices, and the scalar  $b$  is stored as a diagonal element of  $\mathbf{B}$ . If  $\mathbf{X}$  is a null matrix, then the whole set of

latent vectors has been found, otherwise the procedure is re-iterated.

##### 4) Prediction of the dependent variables

The dependent variables are predicted using the multivariate regression formula as

$$\mathbf{Y} = \mathbf{T} \mathbf{B} \mathbf{C}^T = \mathbf{X} \mathbf{B} \mathbf{P}_{\text{PLS}}$$

with

$$\mathbf{B}_{\text{PLS}} = \mathbf{P}^{T+} \mathbf{B} \mathbf{C}^T$$

(where  $\mathbf{P}^{T+}$  is the Moore-Penrose pseudo-inverse of  $\mathbf{P}^T$ ).