Article

A Comparison Study on Nonlinear Dimension Reduction Methods with Kernel Variations: Visualization, Optimization and Classification

Katherine C. Kempfert^{1,†}, Yishi Wang ²*, Cuixian Chen ², and Samuel W.K. Wong ³

- ¹ University of Florida; kkempfert2@ufl.edu
- ² University of North Carolina Wilmington; {wangy, chenc}@uncw.edu
- ³ University of Waterloo; samuel.wong@uwaterloo.ca
- * Correspondence: wangy@uncw.edu; Tel.: +1-910-962-3292
- + Current address: The University of Florida. Gainesville, FL 32611

Version May 23, 2019 submitted to Intelligent Data Analysis

- Abstract: Because of high dimensionality, correlation among covariates, and noise contained in data,
- ² dimension reduction (DR) techniques are often employed to the application of machine learning
- algorithms. Principal Component Analysis (PCA), Linear Discriminant Analysis (LDA), and their
- ⁴ kernel variants (KPCA, KLDA) are among the most popular DR methods. Recently, Supervised Kernel
- ⁵ Principal Component Analysis (SKPCA) has been shown as another successful alternative. In this
- ⁶ paper, brief reviews of these popular techniques are presented first. We then conduct a comparative
- ⁷ performance study based on three simulated datasets, after which the performance of the techniques
- are evaluated through application to a pattern recognition problem in face image analysis. The gender
- classification problem is considered on MORPH-II and FG-NET, two popular longitudinal face aging
- databases. Several feature extraction methods are used, including biologically-inspired features (BIF),
- local binary patterns (LBP), histogram of oriented gradients (HOG), and the Active Appearance
- ¹² Model (AAM). After applications of DR methods, a linear support vector machine (SVM) is deployed
- with gender classification accuracy rates exceeding 95% on MORPH-II, competitive with benchmark
- results. A parallel computational approach is also proposed, attaining faster processing speeds and
- similar recognition rates on MORPH-II. Our computational approach can be applied to practical
- ¹⁶ gender classification systems and generalized to other face analysis tasks, such as race classification
- and age prediction.

Keywords: Dimension Reduction; PCA; LDA; FDA; KPCA; KFDA; SKPCA; SVM; Parameter
 Optimization; Gender Classification; MORPH-II.

20 1. Introduction

Due to advances in data collection and storage capabilities, the demand has been growing 21 substantially for gaining insights into high-dimensional, complex-structured, and noisy data. 22 Researchers from diverse areas have applied DR techniques to visualize and analyze such data 23 [1,2]. DR techniques are also helpful to address the issues of collinearity and " $p \gg n$ " (i.e., number of 24 features exceeding the sample size in a dataset), by projecting the data into a lower dimensional space 25 with less correlation, so that classical statistical methods can be applied [3]. Principal Component 26 Analysis (PCA) [4,5] is a well-studied algorithm with the goal of projecting input features onto a 27 lower dimensional subspace while preserving the largest variance possible; lower dimensionality 28 permits easier visualization, for example via heat maps. While PCA is a fully automatic algorithm, DR 29 techniques that account for domain expertise via user input have also been more recently studied [6,7]. 30 For classification problems, in which the label information as the response variable is available, Linear 31 Discriminant Analysis (LDA) (sometimes referred to as Fisher's Discriminant Analysis (FDA)) can 32 be used for DR by minimizing intra-class variation and maximizing inter-class variation [8,9]. Since 33 PCA only utilizes the correlation or covariance matrices, it is considered an unsupervised approach, 34

³⁶ function. Despite the dissimilarities, both PCA and LDA search for linear combinations of the features

and, therefore, can be applied in linearly separable types of problems [10]. The main challenge is that

many problems in practical applications of machine learning are nonlinear [11,12]. For nonlinear DR,

³⁹ kernel methods are popular choices because of their flexibility [13–15], e.g. Kernel PCA [16], Kernel

LDA for two classes [17], and more generalized Kernel LDA for multiple classes [18]. For kernel methods, it is also possible to design specialized kernels based on domain knowledge of a problem

42 [19,20].

Given the problems in image analysis of high dimensionality and complex correlation structures,

⁴⁴ DR techniques are often a necessary step [21]. Thus, variants of PCA, LDA, and their kernel extensions

have been popular in computer vision with applications of image classification and discrimination

⁴⁶ [22–24]. Studies include Eigenfaces [25], Fisherfaces [26], face recognition with KPCA [27], face

recognition with Kernel Direct LDA [28], 2D-PCA [29], 2D-LDA [30], among many others. When
 there are sufficient labeled face images, LDA is experimentally reported to outperform PCA for face

recognition [26]. In the case of a small training set, the conclusion could be reversed [23]. Studies

⁴⁵ recognition [20]. If the case of a small training set, the conclusion could be reversed [20]. Studies

⁵⁰ comparing classification performance of PCA, LDA, and their kernel variations include [31,32]. The

⁵¹ connections among KLDA, KPCA, and LDA are further discussed in [33]. By incorporating labeling

⁵² information into the construction of the objective function, Supervised Kernel PCA (SKPCA) [34]

has been proposed for visualization, regression, and classification. A modified version of SKPCA for

classification problems can be found in [35]. These studies suggest that SKPCA works well in practice
 among different DR algorithms [36–38]. Moreover, it has been found in [39] that with bounded kernels,

⁵⁶ projections from SKPCA are uniformly converging, regardless of the input features' dimension.

projections from ord errare almorning converging, regaratess of the input

57 2. Associated Work

In recent years, facial demographic analysis has become popular in computer vision, because of its 58 broad applications in human-computer interaction (HCI), security, surveillance, and marketing, which 59 can benefit from the automatic estimation of characteristics like age, gender, and race. Recent surveys 60 on demographic estimation from biometrics are presented in [40,41]. Specifically, a major task is gender 61 classification, aiming to automatically determine if a person is male or female. Beyond computer 62 vision, the topic has been studied extensively by anthropologists, sociologists, and psychologists. 63 Gender can easily be identified by humans, achieving 96% accuracy in an experiment classifying 64 photographs of adult faces [42]. Automating gender classification has been a priority in real-world 65 applications. A number of biometrics have been used to identify gender, including face, voice, gait, 66 handwriting, and even the iris [41]. However, gender classification from faces is the most common, probably because photography of faces is non-intrusive and ubiquitous. Ng et al. provide a survey of 68 gender classification via face and gait [43]. 69 Gender classification with faces launched in 1990, when neural networks were applied directly to 70

pixels from face photographs [44,45]. Many other early studies utilized the geometric-based approach 71 to represent human faces, relying on measurements of facial landmarks [46,47]. Though intuitive, such approaches are sensitive to the placement of landmarks, which can only accommodate frontal 73 representations of the face, and may omit some important information from the face (such as texture of 74 the skin). In recent years, the appearance-based methods have been more commonly adopted, which 75 rely on a transformation of an image's pixels [48–50]. Such methods capture both the geometric 76 relationships of the face and texture information. However, a drawback is their sensitivity to 77 illumination and viewpoint variations. Other issues are associated with the high dimensionality 78 of the transformed pixels, which will be discussed further in the next paragraph. Some most recent 79 gender classification studies involve convolutional neural networks (CNN) [51–54]. Though CNNs 80 have reached state-of-the-art accuracy rates, they are known to be less interpretable than some other 81

82 methods.

Pixels often contain high redundancy and noise, which cannot be removed completely by 83 pre-processing steps. Hence, the vectors resulting from *appearance-based* feature extraction methods 84 genetically inherit redundancy and noise. Popular image feature extraction methods include local texture techniques such as local binary patterns (LBP) [55–58], Gabor filters [59], biologically-inspired 86 features (BIF) [60,61], and histogram of oriented gradients (HOG) [60]. Such methods could lead 87 to a high dimension of extracted features, thwarting practical applications by increasing runtime 88 and memory consumption. When " $p \gg n$ ", for which the dimension of the feature space exceeds 89 the sample size of the dataset, a fundamental assumption of many standard statistical procedures is violated. Additionally, collinearity of features can cause numerical problems, while noisy features 91 can obscure true relationships with the response variable and hinder predictive performance. These 92 significant issues motivate the use of DR techniques. The fundamental goal of DR is to extract and 93 retain information in a lower dimensional space. Many of these methods fall under manifold learning, identifying a low-dimensional manifold embedded in a high-dimensional ambient space [62]. Even though PCA and LDA have been widely considered as popular and effective approaches Qŕ

for DR in machine learning, their kernel versions are much less investigated. To our best knowledge,
KPCA, KLDA, and SKPCA have never before been directly compared on visualization and classification
performance through simulations and practical applications to face image analysis problems.

Our main contributions in this study can be summarized as follows. (1) The nonlinear manifold 100 learning projections for KPCA, KLDA, and SKPCA are directly compared with visualization through 10: simulated datasets. (2) Motivated by the nonlinear nature of soft-biometric analysis problems, we 102 utilize KPCA, KLDA, and SKPCA for dimension reduction on four types of appearance-based extracted 103 features (BIF, HOG, LBP, and AAM) for the gender classification task. Moreover, the classification 104 performance is compared systematically on parameter optimization. (3) For applications to practical 105 large-scale systems, we propose an additional parallel computational framework that can decrease 106 runtime while maintaining similar classification rates. 107

The remainder of the paper is structured as follows. In Section 3, we review the theory of KPCA, SKPCA, and KLDA. In Section 4, we conduct simulation studies to visualize projections. We propose our machine learning methods for gender classification on Morph-II in Section 5. The comparative performance of KPCA, SKPCA, and KLDA on Morph-II is presented and discussed in Section 6. The performance of these DR methods is further compared in Section 7 through application to the FG-NET dataset. The computational framework for large-scale practical systems is proposed in Section 8 and investigated on Morph-II. Finally, we conclude and offer future directions of research in Section 9.

115 3. Kernel-Based Dimension Reduction Methods

The nonlinearity in a classification problem can often be addressed by kernel-based DR methods, 116 with the appropriate choice of kernels. The driving reasons are the nonlinearity of chosen kernels, 117 flexibility of tuning parameter selection, and most importantly, the kernel tricks. Mercer's theorem 118 guarantees that a symmetric positive-definite function can be written as the sum of a convergent 119 sequence of product functions, which potentially project the data into infinite-dimensional space [63]. 120 Thus, it is feasible to separate the data in the new space. On the other hand, Representer Theorem 121 shows that the solution for certain kernel methods lies in the finite-dimensional span of the training 122 data [63,64]. This is very helpful, since we do not need to compute the coordinates of the projected 123 data in the infinite-dimensional space, but only the inner products between all pairs of data in the 124 feature space. 125

126 3.1. Notations

With the goal of emphasizing the connections between KPCA, SKPCA, and KLDA, we define thefollowing notations for classification problems.

Let \mathcal{X} be the feature space, a non-empty subset in \mathbb{R}^p with p as the number of covariates for each subject. Let \mathcal{Y} be the space for the response variable, a subset in \mathbb{R} . Let $\{(x_1, y_1), \dots, (x_n, y_n)\} \subset$

 $\mathcal{X} \times \mathcal{Y}$ be a series of *n* independent observations following a joint probability measure $P_{\mathcal{X},\mathcal{Y}}$. Let $Y = [y_1, y_2, \dots, y_n]^T$ denote the outcomes of the response variable. Let *X* be an $n \times p$ feature matrix, with x_i^T as the *i*-th row for $i = 1, \dots, n$, and $x^{(l)} \in \mathbb{R}^n$ for $l = 1, \dots, p$ as its *l*-th column. Thus, the *X* matrix can be written as:

$$X = \left[x_1, x_2, \cdots, x_n\right]^T = \left[x^{(1)}, x^{(2)}, \cdots, x^{(p)}\right].$$

Without loss of generality, we may assume that each column of the *X* matrix is normalized, such that the mean of $x^{(l)}$ is 0 and standard deviation is 1, for $l = 1, \dots, p$.

Let Σ be the sample covariance matrix of *X*. We then have

$$\sum_{p \times p} = \frac{1}{n-1} X^T X = \frac{1}{n-1} \sum_{i=1}^n x_i x_i^T.$$
 (1)

Let \mathcal{F} be a reproducing kernel Hilbert space on \mathcal{X} from a kernel function $k(\cdot, \cdot)$, which is a Mercer kernel (symmetric and positive-definite), and \mathcal{G} be a reproducing kernel Hilbert space on \mathcal{Y} from a kernel function $l(\cdot, \cdot)$.

For the kernel $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$, its associated space \mathcal{F} may be infinite-dimensional, but with some additional conditions, the minimizer of a regularized risk function lies in the finite span of the training observations [63]. Additionally, it has been shown [63] that there exists a function

$$\phi: \mathcal{X} \to \mathcal{F} \tag{2}$$

such that for all $x, x' \in \mathcal{X}$,

$$k(x, x') = \langle \phi(x), \phi(x') \rangle,$$
 (3)

where $\langle \cdot \rangle$ is the dot product. Let *K* be a matrix such that its *ij*-th element is $k(x_i, x_j)$. We then have

$$K = \{k(x_i, x_j)\}_{ij} = \{\langle \phi(x_i), \phi(x_j) \rangle\}_{ij} = \Phi(X)\Phi(X)^T,$$
(4)

where $\Phi(X) = [\phi(x_1), \phi(x_2), \cdots, \phi(x_n)]^T$. Here, the kernel matrix *K* is the Gram matrix of the $\phi(x_i)$'s.

135 3.2. Principal Component Analysis and Kernel Principal Component Analysis

In standard PCA, we seek an orthogonal transformation matrix A satisfying

$$T_{n \times d} = X_{n \times p} A_{p \times d},$$
(5)

where $T = [t_1, t_2, \dots, t_d]$ for some $d \le p$, such that each column vector t_i successively inherits maximal proportion of variance from the column vectors $x^{(l)}$'s, while ensuring the projection directions are perpendicular. The solutions can be expressed as the eigenvalue problem

$$\Sigma a_i = \lambda_i a_i,\tag{6}$$

where a_i is the *i*-th column of A, for i = 1, ..., d.

Following the work of [65], PCA can be extended to KPCA by first choosing a Mercer kernel k, with which x_i is transformed to $\phi(x_i)$. This maps the features in X to $\Phi(X)$. Assume that $\sum_{i=1}^{n} \phi(x_i)$ is a vector with 0 in each entry. With the Gram matrix $K = \Phi(X)\Phi(X)^T$ as defined in (4) and through the kernel trick from (3), we have the eigenvalue problem

$$Ka_i^* = \lambda_i^* a_i^*, \tag{7}$$

where *d* is the desired dimension and a_i^*, \dots, a_d^* are the eigenvectors of *K*, with associated eigenvalues $\lambda_1^* \ge \lambda_2^* \ge \dots \ge \lambda_d^*$. Hence, the advantage of the kernel-based approach is to calculate the Gram matrix *K* without an explicit expression for ϕ . Without the centralization assumption on ϕ , the *K* matrix in (7) can be replaced by

$$K^* = H_n K H_n, \tag{8}$$

where *d* is the desired dimension, $H_n = I_n - \frac{1}{n} \mathbb{1}_n$, I_n is an identity matrix with dimension $n \times n$, and $\mathbb{1}_n$ is a matrix of 1's with dimension $n \times n$.

We note that H_n is idempotent, since it is a square matrix satisfying $H_n = H_n H_n$. For any square matrix *S* with dimension $n \times n$, the average of each column of the matrix $H_n S$ is 0, as is the average of each row of the matrix SH_n . Thus, the K^* matrix is the "centralized" version of the original *K* matrix.

146 3.3. Supervised Kernel Principal Component Analysis

PCA and KPCA are unsupervised methods, since they do not consider the response variable, only considering directions of maximum variability in the covariates. If the goal is classification, this may not be ideal, since the principal components may be unrelated to the class difference. SKPCA is a supervised generalization of KPCA, which aims to find the principal components with maximal dependence on the response variable. Drawing from [34] and [35], we formulate SKPCA as follows.

In SKPCA, class information is incorporated by maximizing the Hilbert Schmidt independence criterion (HSIC) [66]. With the aforementioned reproducing kernel Hilbert spaces \mathcal{F} on \mathcal{X} and \mathcal{G} on \mathcal{Y} and related kernel functions $k(\cdot, \cdot)$ and $l(\cdot, \cdot)$ respectively, the HSIC can be expressed as

$$HSIC(P_{\mathcal{X},\mathcal{Y}},\mathcal{F},\mathcal{G}) = E_{x,x',y,y'}[k(x,x')l(y,y')] + E_{x,x'}[k(x,x')]E_{y,y'}[l(y,y')] - 2E_{x,y}(E_{x'}[k(x,x')]E_{y'}[l(y,y')]),$$
(9)

where $E_{x,x',y,y'}$ represents the expectation on independent pairs of (x, y) and (x', y') (with respect to $P_{\mathcal{X},\mathcal{Y}}$) and $E_{x,x'}$ and alike are the expectations based on various marginal distributions from $P_{\mathcal{X},\mathcal{Y}}$.

With the results from [66], an empirical estimator of (9) is

$$HSIC(X,Y,\mathcal{F},\mathcal{G}) = \frac{1}{(n-1)^2} tr(KH_n LH_n),$$
(10)

where *K* and H_n are defined as before for KPCA and $L = \{1(y_i = y_j)\}_{ij}$ is a link matrix with dimension $n \times n$, where $1(\cdot)$ is an indicator function with value 1 if the event is true and 0 otherwise.

Similarly as for KPCA, *K* and *L* can be adjusted to satisfy the centralization assumption. As discussed previously, H_n is an idempotent matrix. Therefore, following from (10),

$$HSIC^{*}(X, Y, \mathcal{F}, \mathcal{G}) = \frac{1}{(n-1)^{2}} tr(KH_{n}H_{n}LH_{n}H_{n})$$
$$= \frac{1}{(n-1)^{2}} tr(H_{n}KH_{n}H_{n}LH_{n})$$
$$= \frac{1}{(n-1)^{2}} tr(K^{*}L^{*}),$$
(11)

where K^* and L^* are the "centralized" versions of the *K* and *L* matrices respectively.

On another note, in the binary gender classification problem, rank(L) = 2 and $rank(KH_nLKH_n) \le 2$ [35]. Therefore, we modify the link matrix according to [35] by

$$L = \{1(y_i = y_j) \times k(x_i, x_j)\}_{ij}.$$
(12)

It can be shown that maximization of (10) is equivalent to solving the generalized eigenvalue problem

$$Av_i = \lambda_i K v_i, \tag{13}$$

where $A = KH_nLH_nK$ and each v_i is an eigenvector with related eigenvalue λ_i for $i = 1, \dots, d$, where *d* is the desired dimension [35]. Therefore, the main advantage of the link matrix in (12) becomes apparent: the rank of KH_nLKH_n may increase, permitting more eigenvalues to be computed.

160 3.4. Linear Discriminant Analysis and Kernel Linear Discriminant Analysis

Given a dataset with finite classes, LDA aims to find the best set of features to discriminate among the classes. We first review standard LDA, then generalize to KLDA. We note that sometimes parametric assumptions for LDA are made, such as that observations from each class are normally distributed with common covariance. Here, we make no such assumptions. Suppose that each observation x_i for $i = 1, \dots, n$ belongs to exactly one of C classes. Define the following feature vectors: $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$ as the overall mean and $\bar{x_c} = \frac{1}{n_c} \sum_{i=1}^{n} x_i 1(x_i \in \text{class c})$ as the mean of the *c*-th class with n_c the size of the *c*-th class in the sample, for $c = 1, \dots, C$.

In standard LDA, we seek to maximize the objective function

$$J(v) = \frac{v^T S_B v}{v^T S_W v},\tag{14}$$

where v is a $p \ge 1$ vector, S_B is the between-class scatter matrix, and S_W is the within-class scatter matrix defined by

$$S_B_{p \times p} = \sum_c n_c (\bar{x}_c - \bar{x}) (\bar{x}_c - \bar{x})^T \text{ and}$$

$$S_W_{p \times p} = \sum_c \sum_{i \in c} (x_i - \bar{x}_c) (x_i - \bar{x}_c)^T.$$
(15)

Hence, maximizing J(v) involves finding some rotation of the scatter matrices such that the "distance" between the groups is maximized relative to the variations within each group.

Maximization of J(v) in (14) is equivalent to solving the generalized eigenvalue problem

$$S_B v_i = \lambda_i S_W v_i, \tag{16}$$

where each v_i is an eigenvector with corresponding eigenvalue λ_i , for $i = 1, \dots, d$, where *d* is the desired dimension.

LDA is generalized to KLDA using the kernel representation from (3). Analogously to LDA above, we seek a solution v^* that will result in the maximization of the objective function

$$J(v) = \frac{v^T S_B^* v}{v^T S_w^* v'},$$
(17)

where now $v \in \mathcal{F}$ and S_B^* and S_W^* are the between-class and within-class scatter matrices in \mathcal{F} defined by

$$m^{\phi} = \frac{1}{n} \sum_{i=1}^{n} \phi(x_{i}),$$

$$m_{c}^{\phi} = \frac{1}{n_{c}} \sum_{i=1}^{n} \phi(x_{i}) 1(x_{i} \in \text{class c}),$$

$$S_{B}^{*} = \sum_{c} n_{c} (m_{c}^{\phi} - m^{\phi}) (m_{c}^{\phi} - m^{\phi})^{T}, \text{ and}$$

$$S_{W}^{*} = \sum_{c} \sum_{i \in c} (\phi(x_{i}) - m_{c}^{\phi}) (\phi(x_{i}) - m_{c}^{\phi})^{T}.$$
(18)

Version May 23, 2019 submitted to Intelligent Data Analysis

The above expressions involve knowledge of ϕ , which is often not available. It can be shown that equation (17) is equivalent to

$$I(u) = \frac{u^T M u}{u^T N u},\tag{19}$$

where

$$M_{c} = (M_{cj})_{j} = \left(\frac{1}{n_{c}}\sum_{h=1}^{n}k(x_{j},x_{h})\mathbf{1}(x_{h} \in \text{class c})\right)_{j},$$

$$\bar{M} = (\bar{M}_{j})_{j} = \left(\frac{1}{n}\sum_{h=1}^{n}k(x_{j},x_{h})\right)_{j},$$

$$M = \sum_{c}n_{c}(M_{c} - \bar{M})(M_{c} - \bar{M})^{T},$$

$$K_{c} = K \times outer(X, X_{c}),$$

$$N = \sum_{c}K_{c}H_{n_{c}}K_{c}^{T},$$

(20)

 X_c is a matrix of dimension $n_c \times p$ with rows being features from the *c*-th class, and *outer*(X, X_c) is

an $n \times n_c$ matrix with its *ij*-th element as $1(x_i$ is the *j*-th observation in class c). A full discussion of KLDA can be found in [17].

Maximization of J(u) in equation (19) is equivalent to solving the generalized eigenvalue problem

$$Mu_i = \lambda_i Nu_i, \tag{21}$$

- where each u_i is an eigenvector with associated eigenvalue λ_i , for $i = 1, \dots, d$ with d as the desired dimension.
- ¹⁷⁷ Comparing the generalized eigenvalue problems in (16) and (21), the structures of matrices S_B ¹⁷⁸ and *M* are similar, since both "measure" the variation between different classes.

Let $W_c = [w_{c,1}, \dots, w_{c,n_c}] = K_c H_{n_c}$, a matrix of dimension $n \times n_c$. Due to the centralization function of H_{n_c} , W_c has row-sum equal to zero for every row. Besides, $K_c H_{n_c} (K_c H_{n_c})^T = W_c W_c^T = \sum_{i=1}^{n_c} w_{c,i} w_{c,i}^T$. For the matrix N, due to the idempotent property of H_{n_c} ,

$$N = \sum_{c} K_{c} H_{n_{c}} H_{n_{c}} K_{c}^{T} = \sum_{c} \sum_{i=1}^{n_{c}} w_{c,i} w_{c,i}^{T}.$$
(22)

Thus, the matrix N has an identical structure to the S_W and S_W^* matrices, which "measure" the overall variation within groups.

181 4. Visualization on Simulation Studies

To visualize and improve understanding of the manifold learning methods KPCA, SKPCA, and KLDA, we apply them in three simulation studies. For comparison, the linear techniques PCA and LDA are also considered. Each dataset contains nonlinear patterns, and the goal is to transform the data to be linearly separable. For this reason, the radial basis function (RBF)

$$k(x_i, x_j) = e^{-\delta ||x_i - x_j||_2^2}$$
(23)

is chosen as a kernel for each pair of observed vectors x_i, x_j . For each DR method, a range of values for the tuning parameter δ are tested and selected to visually separate the classes. A full discussion of

the RBF kernel, among others, can be found in [67]. Figures 1, 2, and 3 compare the original data to

¹⁸⁵ 2-dimensional projections from each DR method. In each plot, color corresponds to the true class to

186 which an observation belongs.



Figure 1. Wine Chocolate Simulation Study.

In the first simulation study, the original data are plotted in 3D in Figure 1(a); the green sphere 187 is embedded within the magenta group, necessitating nonlinear manifold learning. The KLDA 188 projections in (b) are linearly separable with very good variation between the classes and a fair amount 189 of variation within the classes. KPCA and SKPCA projections in (c) and (d) are at least approximately 190 linearly separable, as it is not clear whether there is a linear boundary that perfectly separates the 191 two classes. In (e), PCA fails to linearly separate the groups, rotating the wine chocolate in 2D. The 192 maximum dimension LDA can retain is p - 1; with 2 classes, the projections must be plotted on a 1D 193 number line, given in (f). Points from the two classes overlap considerably in plots (e) and (f). 194



Figure 2. Apple Tart Simulation Study.

In the second simulation study, the original data in Figure 2(a) follow a nonlinear pattern. In (b), KLDA produces groups which are linearly separable. The KPCA projections are approximately linearly separable in (c); however, there is some overlap between groups, especially the green and pink groups in the third quadrant. In (d), SKPCA produces almost linearly separable groups. In plots (e) and (f), PCA and LDA simply rotate the original data in 2D space, as expected.



Figure 3. Swiss Roll Simulation Study.

For the third simulation study, the original data in Figure 3(a) are in 3D and follow a swirling, nonlinear pattern. In (b), KLDA yields favorable results; the groups are well-separated linearly. KPCA and SKPCA in (c) and (d) also produce good results, although in (c) more separation between the purple and bright green groups would be ideal. In (e) and (f), respectively, PCA and LDA merely rotate the original data projected in 2D space; there is no linear separation between the magenta and purple groups, nor between the two green groups.

For all three simulation studies, KLDA, KPCA, and SKPCA are effective to transform the data 206 into linearly separable groups. In all cases, the projected data become approximately linearly separable 207 after applying KLDA, KPCA, or SKPCA. In general, KLDA and SKPCA perform the best here. Their 208 success over KPCA is expected, since KLDA and SKPCA are supervised techniques. On the other hand, 209 results indicate that KPCA and SKPCA are more sensitive than KLDA to different choices of tuning 210 parameters. Hence, SKPCA and KPCA may perform better for alternative choices of parameters. In 211 all our studies, the nonlinear techniques outperform linear PCA and LDA. These preliminary studies 212 suggest the radial kernel is appropriate for our face analysis experiments. 213

214 5. Kernel-based Dimension Reduction Optimization and Classification on Morph-II

Motivated by the nonlinear nature of facial demographic analysis, we propose and implement a novel machine learning process for the Morph-II dataset. We consider the kernel-based DR methods KPCA, SKPCA, and KLDA on three types of appearance-based extracted features (LBP, BIF, and HOG) for the gender classification task. We illustrate parameter optimization and compare the performance of these methods on Morph-II.

220 5.1. Longitudinal Face Database

MORPH [68] is one of the most popular face databases available to the public, especially for age estimation, race classification, and gender classification. Multiple versions of MORPH have been released, and the version adopted in this work is the 2008 MORPH-II non-commercial release (referred to as Morph-II in this paper). Morph-II includes over 55,000 mugshots with longitudinal spans and metadata such as date of birth, race, gender, and age.

In addition to its size, Morph-II presents challenges because of disproportionate race and gender ratios. About 84.6% of images are of males, while only about 15.4% of images are of females. Imbalanced classes are known to negatively affect certain classification algorithms [69]. Moreover, Morph-II is skewed in terms of race, with approximately 77.2% of images picturing black subjects. Guo et al. found that age, gender, and race interact for demographic analysis tasks including gender classification, race classification, and age prediction [48,60,70], so both race and gender imbalance in Morph-II can hamper gender classification.

233 5.2. Subsetting Scheme

To overcome the uneven race and gender distributions in Morph-II, Guo et al. proposed a subsetting scheme [48]. Since then, many studies on Morph-II have adopted such an evaluation protocol. Based on discussions in Guo et al. [48], a new automatic subsetting scheme is proposed in [71], aiming to automatically ensure independent training and testing sets. Additionally, inconsistencies in age, gender, and race in Morph-II have been identified and corrected in [71]. After following the steps to clean MORPH-II outlined in [71], we apply the automatic subsetting scheme, summarized in Figure 4 and described below.

Let W be the Whole Morph-II dataset, S the selected training/testing set, and R the remaining set. 241 We further divide S into even subsets S_1 and S_2 . Separately within each subset S_1 and S_2 , we fix the ratios of white (W) to black (B) images as 1:1 and male (M) to female (F) images as 3:1. Further, S_1 and 243 S_2 have been selected such that the age distributions within each set are similar (details shown in [71]). 244 The gender and race summaries for the subsetting scheme are shown in Table 2. Most importantly, 245 the sets R, S_1 , and S_2 are independent; no sets share images from the same subject. We use S as an 246 alternating training and testing set. First, we train on S_1 and test on $S_2 \cup R$, then we train on S_2 and test on $S_1 \cup R$. The final classification accuracy is obtained by averaging the classification accuracies 248 from the alternations. 249



Figure 4. Flowchart representing our subsetting scheme [71] for MORPH-II, which improves the one from [48].

	WF	BF	WM	BM	dF	dM	Overall	F	М
S1	1,285	1,285	3855	3,855	0	0	10,280	2570	7,710
S2	1,285	1,285	3,855	3,855	0	0	10,280	2,570	7,710
R	0	3,150	220	28,980	144	1,850	34,344	3,294	31,050
Overall	2,570	5,720	7,930	36,690	144	1,850	54,904	8,434	46,470

Table 2. Number of Images in Subsets by Race and Gender

Race-gender combinations are abbreviated, e.g., BF represents the black female group. Abbreviations dF and dM represent those who are neither black nor white in race.

5.3. Facial Feature Extraction 250

In computer vision, image preprocessing is often an essential first step to reduce unnecessary 251 variation, decrease pixel dimension, and simplify pixel encoding. Despite the standard format of police 252 photography in mugshots, Morph-II photographs vary in head-tilt, camera distance, occlusion, and 253 illumination. We address this variation as follows. Images are first converted to grayscale. Next, faces 254 are automatically detected, eliminating the background and hair, so that no external cues can be used 255 to classify gender. The resulting images are centered and scaled with respect to the center of the irises. 256 Finally, the images are cropped to be 70 pixels tall by 60 pixels wide. Full methodological details are 257 given in [72] and align with standard preprocessing protocols from face analysis. 258

After preprocessing, pixel-related features are extracted from the face images in Morph-II. 259 As discussed previously, there are numerous approaches for feature extraction. In this study on 260 Morph-II, we incorporate domain expertise by choosing three well-established appearance-based 261 models from image analysis: local texture techniques such as local binary patterns (LBP) [55-262 58], biologically-inspired features (BIF) [60,61], and histogram of oriented gradients (HOG) [60]. 263 Additionally, these model-based approaches provide "robust interpretation ... by constraining solutions to be face-like" [73]. Detailed documentation of our feature extraction process can be 265 found in [72,74]. 266

	LBP	s = 10, 12, 14, 16, 18, 20
Features		r = 1, 2, 3
	HOG	s = 4, 6, 8, 10, 12, 14
		o = 4, 6, 8
	BIF	s=7 - 37, 15 - 29
		$\gamma = 0.1, 0.2, \dots, 1.0$
	КРСА	$\delta = \pm 0.1, \pm 1, \pm 5, \pm 10, \pm 100$
Dimension Reduction	SKPCA	$\delta = -0.0001, -0.001$
		$\eta = 0.001, 0.01, 0.1, 1$
	KLDA	$\delta = \pm 0.01, \pm 0.1, \pm 1, \pm 5, \pm 10, \pm 100$
Classifier	Linear SVM	$c = 10^{-8}, \dots, 10^{-1}, 1, 10, \dots, 10^{8}$

Table 3. Parameter Summary

5.4. Kernel-Based Dimension Reduction Optimization 267

Tuning parameter selection is essential for kernel-based methods in order to achieve good results. 268 Within the framework of feature extraction, dimension reduction, and the classification model, there 269 are many combinations of parameters to be considered. The main parameters and tested values are 270 summarized in Table 3 and discussed as follows. LBP features have two main parameters: block size s 271 and window radius r. For HOG, the two main parameters are block size s and number of orientations 272 o. For BIF, we consider the block size s and the parameter γ , which represents the spatial aspect ratio; 273 there is also a choice of pooling operation, which we select here as the standard deviation operation. 274

For each dimension reduction method, the radial kernel

$$k(x_i, x_j) = e^{\delta ||x_i - x_j||_2^2}$$
(24)

is used for each pair of observation vectors x_i , x_j , based on the results from our simulation studies. In the kernel, we must select the tuning parameter δ , which scales the extent of similarity between pairs of vectors. This parameter must be chosen with particular care, since a poor choice can result in transformed features with little to no variability. Empirically, we observed that SKPCA was more sensitive than KLDA and KPCA to the choice of δ ; values of δ at or above 1 resulted in a rank deficient matrix and failure to compute all requested eigenvalues. For SKPCA, we consider an additional scaling parameter η in the modified link function proposed by Wang et al. [35]:

$$l(y_i, y_j) = e^{\eta \delta ||x_i - x_j||_2^2},$$
(25)

for all observed responses y_i, y_j in the same class. The scale parameter η enables the weighing of dependence between the covariates and response.

Finally, we choose a linear SVM to classify gender based on the dimension-reduced, transformed features. The motivation for this classifier is discussed in the next section. The main parameter for linear SVM is the cost c, which measures the extent to which misclassification in training will be permitted. We consider values of c from 10^{-8} to 10^8 .

Method	Feature	Parameters	Accuracy
	BIF $s = 7 - 37$, $\gamma = 0.1$	$\delta = -1, c = 10$	0.882
KPCA	BIF $s = 7 - 37$, $\gamma = 0.6$	$\delta = -1, c = 10$	0.882
	BIF $s = 15 - 29, \gamma = 0.1$	$\delta = -1, c = 100$	0.882
	BIF $s = 15 - 29, \gamma = 0.6$	$\delta = -1, c = 10$	0.882
	HOG $s = 4, o = 4$	$\delta = -100, c = 0.1$	0.917
	HOG $s = 4, o = 4$	$\delta = -5, c = 0.001$	0.919
	HOG $s = 4, o = 4$	$\delta = -1, c = 0.001$	0.917
	HOG $s = 4, o = 4$	$\delta = -0.1, c = 0.1$	0.917
	LBP $s = 10, r = 1$	$\delta = -100, c = 0.1$	0.912
	LBP $s = 10, r = 1$	$\delta = -5, c = 0.1$	0.912
	LBP $s = 10, r = 1$	$\delta = -1, c = 0.001$	0.912
	LBP $s = 10, r = 1$	$\delta = -0.1, c = 0.1$	0.912
	BIF $s = 7 - 37$, $\gamma = 0.2$	$\delta = 0.0001, \eta = 0.1, c = 1$	0.899
SKPCA	BIF $s = 7 - 37$, $\gamma = 0.8$	$\delta = 0.0001, \eta = 0.1, c = 1$	0.899
	BIF $s = 15 - 29, \gamma = 0.5$	$\delta = 0.0001, \eta = 0.1, c = 1$	0.899
	HOG $s = 6, o = 6$	$\delta = 0.0001, \eta = 0.001, c = 1$	0.931
	HOG $s = 6, o = 6$	$\delta = 0.0001, \eta = 0.01, c = 0.001$	0.931
	HOG $s = 6, o = 6$	$\delta = 0.0001, \eta = 0.1, c = 0.001$	0.931
	LBP $s = 14, r = 2$	$\delta = 0.0001, \eta = 0.001, c = 1$	0.937
	LBP $s = 14, r = 2$	$\delta = 0.0001, \eta = 0.01, c = 1$	0.937
	LBP $s = 14, r = 2$	$\delta = 0.0001, \eta = 0.1, c = 1$	0.938
	LBP $s = 14, r = 2$	$\delta = 0.0001, \eta = 1, c = 1$	0.939
	BIF $s = 7 - 37, \gamma = 0.3$	$\delta = -\overline{1, c} = 10$	0.875
KLDA	BIF $s = 7 - 37$, $\gamma = 0.6$	$\delta = -1, c = 100$	0.875
	BIF $s = 15 - 29$, $\gamma = 0.2$	$\delta = -1, c = 10$	0.875
	BIF $s=15-29, \gamma=0.8$	$\delta = -1, c = 100$	0.875
	HOG $s = 4, o = 4$	$\delta = 1, c = 1$	0.917
	HOG $s = 4, o = 6$	$\delta = 1, c = 1$	0.917
	HOG $s = 12, o = 8$	$\delta = -1, c = 100$	0.904
	LBP $s = 10, r = 1$	$\delta = -0.1, c = 1$	0.906
	LBP $s = 10, r = 1$	$\delta = 1, c = 1$	0.908
	LBP $s = 14, r = 1$	$\delta = 0.1, c = 10$	0.898

Table 4. Tuning Results on a Subset of MORPH-II

We tune on small subsets of Morph-II to reduce runtime, memory consumption, and risk of over-fitting. 1000 images from S_1 and 1000 images from S_2 are randomly selected. The standard method of grid search is followed for tuning on these subsets. For each combination of parameters, a model is trained on the subset from S_1 and then tested on the subset from S_2 . For each dimension

reduction method paired with each feature type (BIF, HOG, and LBP), the best three or four accuracy 285 rates from tuning are obtained. (Except in the case of ties, we choose only the top three accuracy 286 rates.) The tuning results for these top-performing parameters are given in Table 4. The parameters corresponding to these maximum accuracy rates are applied to the full dataset through the previously 288 discussed evaluation protocol. Although this protocol involves testing on images from S_1 and S_2 , any 289 overlap of images is minor (in each testing set, less than 2.3% of images have been used in tuning) and 290 has little impact on the reported accuracy (discussed in Section 6). Regardless, the tuning parameters 291 are selected through the same procedure, so the classification performances can be fairly compared among all considered DR methods. 293

Method	Feature	Parameters	$Acc^{(1)}$	TPR ⁽²⁾	TNR ⁽³⁾	Mem ⁽⁴⁾	Time ⁽⁵⁾
	BIF $s = 7 - 37$, $\gamma = 0.1$	$\delta = -1, c = 10$	0.9296	0.9473	0.8127	34.04	42.26
КРСА	BIF $s = 7 - 37$, $\gamma = 0.6$	$\delta = -1, c = 10$	0.9297	0.9455	0.8112	34.68	36.94
	BIF $s = 15 - 29$, $\gamma = 0.1$	$\delta = -1, c = 100$	0.9071	0.9377	0.7050	31.74	33.83
	BIF $s = 15 - 29$, $\gamma = 0.6$	$\delta = -1, c = 10$	0.9096	0.9374	0.7266	31.80	35.97
	HOG $s = 4, o = 4$	$\delta = -100, c = 0.1$	0.9391	0.9726	0.7172	34.00	31.54
	HOG $s = 4, o = 4$	$\delta = -5, c = 0.001$	0.9391	0.9727	0.7170	34.00	30.86
	HOG $s = 4, o = 4$	$\delta = -1, c = 0.001$	0.9391	0.9724	0.7192	34.00	32.17
	HOG $s = 4, o = 4$	$\delta = -0.1, c = 0.1$	0.9364	0.9626	0.7634	34.35	31.41
	LBP $s = 10, r = 1$	$\delta = -100, c = 0.1$	0.9391	0.9726	0.7172	34.00	31.54
	LBP $s = 10, r = 1$	$\delta = -5, c = 0.1$	0.9391	0.9726	0.7172	34.00	30.86
	LBP $s = 10, r = 1$	$\delta = -1, c = 0.001$	0.9391	0.9724	0.7192	34.00	32.17
	LBP $s = 10, r = 1$	$\delta = -0.1, c = 0.1$	0.9364	0.9626	0.7634	34.35	31.41
	BIF $s = 7 - 37$, $\gamma = 0.2$	$\delta = 0.0001, \eta = 0.1, c = 1$	0.9507	0.9616	0.8781	35.48	42.04
SKPCA	BIF $s = 7 - 37, \gamma = 0.8$	$\delta = 0.0001, \eta = 0.1, c = 1$	0.9532	0.9639	0.8823	33.04	38.34
	BIF $s = 15 - 29, \gamma = 0.5$	$\delta = 0.0001, \eta = 0.1, c = 1$	0.9260	0.9477	0.7827	20.03	34.58
	HOG s = 6, o = 6	$\delta = 0.0001, \eta = 0.001, c = 1$	0.9467	0.9645	0.8292	36.69	37.39
	HOG s = 6, o = 6	$\delta = 0.0001, \eta = 0.01, c = 0.001$	0.9489	0.9786	0.7528	38.28	53.96
	HOG $s = 6, o = 6$	$\delta = 0.0001, \eta = 0.1, c = 0.001$	0.9488	0.9786	0.7517	39.83	60.55
	LBP $s = 14, r = 2$	$\delta = 0.0001, \eta = 0.001, c = 1$	0.9585	0.9727	0.8641	28.68	25.33
	LBP $s = 14, r = 2$	$\delta = 0.0001, \eta = 0.01, c = 1$	0.9585	0.9764	0.8642	23.22	38.42
	LBP $s = 14, r = 2$	$\delta = 0.0001, \eta = .1, c = 1$	0.9585	0.9730	0.8640	29.87	28.00
	LBP $s = 14, r = 2$	$\delta = 0.0001, \eta = 1, c = 1$	0.9585	0.9727	0.8640	27.92	22.83
	BIF $s = 7 - 37, \gamma = 0.3$	$\delta = -1, c = 10$	0.9415	0.9539	0.8594	24.89	34.50
KLDA	BIF $s = 7 - 37, \gamma = 0.6$	$\delta = -1, c = 100$	0.9426	0.9558	0.8858	24.74	35.46
	BIF $s = 15 - 29, \gamma = 0.2$	$\delta = -1, c = 10$	0.9131	0.9374	0.7532	22.80	26.78
	BIF $s = 15 - 29, \gamma = 0.8$	$\delta = -1, c = 100$	0.9205	0.9421	0.7783	22.83	33.88
	HOG $s = 4, o = 4$	$\delta = 1, c = 1$	0.9369	0.9517	0.8392	36.52	81.71
	HOG s = 4, o = 6	$\delta = 1, c = 1$	0.9398	0.9545	0.8425	52.48	148.24
	HOG $s = 12, o = 8$	$\delta = -1, c = 100$	0.9175	0.9421	0.7542	21.18	21.57
	LBP $s = 10, r = 1$	$\delta = -0.1, c = 1$	0.9418	0.9578	0.8428	24.58	37.17
	LBP $s = 10, r = 1$	$\delta = 1, c = 1$	0.9417	0.9558	0.8486	24.70	36.45
	LBP $s = 14, r = 1$	$\delta = 0.1, c = 10$	0.9392	0.9543	0.8397	20.77	31.12

Table 5. Gender Classification Results on MORPH-II

(1) Acc represents mean accuracy.

(2) TPR represents mean true positive rate (recall/sensitivity): the proportion of male faces correctly classified.

(3) TNR represents mean true negative rate (specificity): the proportion of female faces correctly classified.

(4) Mem represents mean memory in gigabytes.

(5) Time represents mean runtime in hours for training and testing.

5.5. *Gender Classification*

For the classification part of the modeling, linear SVM is adopted. Many face analysis studies have involved SVM, as summarized in [75]. Briefly, SVM identifies a separating hyperplane with maximal margin between the classes. Several popular kernels for SVM include linear, polynomial, and RBF [67]. We select the linear kernel, because directions of variability in the data are expected to be linear after the nonlinear transformations of KPCA, SKPCA, or KLDA. Indeed, Schölkopf et al. observed this to be true for KPCA in their landmark study [65]. The linear kernel for SVM also reduces

³⁰¹ computational cost, compared to nonlinear kernels.

With the parameters in Table 4 that are selected from tuning on subsets, we implement dimension 302 reduction and classification on the full Morph-II dataset, following the subsetting scheme discussed 303 in Section 5.2. The challenges of the large size of Morph-II, the high dimensionality of the features, and the computational complexity of these dimension reduction methods necessitate the use of 305 high-performance computing (HPC). For example, the kernel matrix for each dimension reduction 306 method is 55134×55134 , requiring approximately 23 gigabytes of storage. Thus, we implement the 307 process on the HiPerGator 2.0 supercomputing cluster at the University of Florida. The code is written 308 in R. The R package *rARPACK* is used to optimize the solving of eigenvalue problems [76], and the e1071 package is utilized for training and testing the SVM model [77]. 310

311 6. Experiment Results

The kernel-based DR methods KPCA, SKPCA, and KLDA are applied to three facial feature extraction methods: BIF, HOG, and LBP. The DR methods transform the feature data, then reduce the 313 dimension. In all cases, a dimension of 100 is retained, substantially lower than the dimension of the 314 original feature space. The dimensionality of 100 is selected as a trade-off between computation time 315 and classification accuracy based on our preliminary studies. The transformed and dimension-reduced 316 data serve as input for the linear SVM, which classifies each image subject as male or female. Additionally, these predicted gender classes are mapped to probabilities through a sigmoid function, 318 following [78]. This process is applied to each alternation of the evaluation protocol: 1) train on S_1 , 319 test on $S_2 \cup R$ and 2) train on S_2 , test on $S_1 \cup R$. The classification results are averaged over these 320 two testing sets. The mean classification accuracy over the testing images is chosen as the evaluation 321 criterion for our methods on Morph-II, as it is the usual performance metric for gender classification 322 [60], especially in similar studies [49,51,52,79]. 323

These mean classification results from Morph-II are shown in Table 5. In addition to the accuracy, 324 the true positive rate (also known as sensitivity or recall) and true negative rate (also called specificity) 325 are given. For this study, we define the true positive rate (TPR) as the proportion of male faces correctly 326 classified, while the true negative rate (TNR) as the proportion of female faces correctly classified. 327 The memory and runtime are also listed in Table 5. The runtime is the total time for training and 328 testing on HPC, i.e., the average of time1 (train on S_1 , test on $S_2 \cup R$) and time2 (train on S_2 , test on 329 $S_1 \cup R$). As mentioned in Section 5.4, there is a small overlap between the tuning and testing sets that 330 could contribute to over-fitting. We have assessed the potential impact of over-fitting on our reported 331 accuracy rates and found it to be very small: it is estimated to be (at most) between 0.09% and 0.2% 332 and to monotonically decrease as reported accuracy rates increase. 333



Figure 5. Receiver operating characteristic (ROC) curve and area under the curve (AUC) are compared by method for gender classification on Morph-II. Each color corresponds to a DR method paired with feature type. For each probability threshold, the true and false positive rates are reported as the averages from the testing sets of the alternating evaluation protocol.

The classification performance is further visualized in Figure 5 through receiver operating characteristic (ROC) curves for the nine combinations of DR method and feature extraction type. For each combination, its displayed curve corresponds to the "best" results from Table 5 (the combination of parameters reaching maximum mean classification accuracy or maximum mean true positive rate in the event of ties). For each alternation of the evaluation protocol, the true and false positive rates in testing are calculated for each probability threshold. To construct the ROC curves, each of the resulting rates for each threshold is averaged over the testing sets.

Table 5 shows that for the feature BIF, SKPCA and KLDA outperform KPCA. For the feature HOG, 341 SKPCA achieves higher accuracy than both KPCA and KLDA, while the latter two techniques perform 342 very similarly. Last, for the feature LBP, SKPCA produces better classification accuracy than KPCA and 343 KLDA. In summary, our experiment's results indicate that SKPCA outperforms KLDA consistently, 344 while KLDA outperforms KPCA for all three features BIF, LBP, and HOG. On the other hand, for 345 KPCA, the features HOG and LBP produce approximately the same accuracies, outperforming BIF. 346 For SKPCA, LBP achieves slightly better results than BIF, while LBP and BIF both outperform HOG. 347 Finally, for KLDA, BIF reaches slightly higher accuracy than LBP, while BIF and LBP both exceed HOG. 348 In most cases, the accuracy (in Table 5) and AUC (in Figure 5) metrics agree on the best methods. 349 An exception is that SKPCA with the HOG features achieves slightly higher accuracy (94.89%) than 350 KLDA with the BIF features (94.18%), but SKPCA with HOG has lower AUC than KLDA with BIF. 351 The other exception is that KPCA with the HOG features has the lowest AUC of the nine methods, but 352 its accuracy is higher than KPCA with the BIF features. In summary, the accuracy and AUC results 353 imply that SKPCA generally performs best for gender classification on Morph-II, while KLDA tends to 354 outperform KPCA. Meanwhile, the LBP and BIF features often yield better classification performance, 355 with less memory usage, than the HOG features.

It is interesting that, overall, LBP achieves even slightly better performance than BIF for the dimension reduction method SKPCA on the task of gender classification, since BIF is popular in demographic analysis such as age estimation, gender classification, and race classification [48,49,60, 70,79]. Another interesting fact is displayed by the results of the true positive and negative rates in Table 5: males have a higher probability of correct identification than females, with the biggest margin exceeding 20%. Our finding is consistent with [61]: females are more challenging to correctly classify
than males, both for automatic approaches and human perception. Similarly, for race classification on
Morph-II, Guo and Mu found in [70] that training a model on female faces (and testing on male faces)
contributed to significantly more errors on average than training on male faces (and testing on female
faces), even when controlling for differences in the training sample sizes. Our results also indicate that,
overall, HOG and LBP outperform BIF for males, while BIF works consistently better than LBP and
HOG for females.

Next, in Table 6 we compare our results to studies using similar methods on Morph-II, as well as recent state-of-the-art works with deep learning on MORPH-II. With the exception of [61], all studies' results in the table are mean testing classification accuracy from an alternating evaluation protocol based on Guo et al [48]. Hence, our results can be directly compared to these studies. With LBP features, SKPCA, and a linear support vector machine (SVM), our gender classification accuracies approximate 96%, competitive with benchmark results. Interestingly, several reported accuracy rates from human observers of gender range from 96% [42] to 96.9% [61]. The similarity in recognition rates between our methods and human observers can further validate the success of our approach.

ween our methods and numan observers can further variate the success of our approa

Method	Accuracy	Reference	Year
BIF+OLPP	98%	[49]	2011
BIF+PLS	97.34%	[49]	2011
BIF+KPLS	98.2%	[49]	2011
BIF+CCA	95.2%	[79]	2014
BIF+KCCA	98.4%	[79]	2014
BIF+rCCA	97.6%	[79]	2014
Multi-scale CNN	97.9%	[52]	2014
Ranking CNN	97.9%	[51]	2015
BIF+Hierarchical-SVM	97.6%	[61]	2015
Human Estimators	96.9%	[61]	2015
LBP+SKPCA+L-SVM	95.85%	This work	2019

Table 6. Comparison Results for Gender Classification on MORPH-II

377 7. Kernel-based Dimension Reduction Optimization and Classification on FG-NET

For further comparison between KPCA, SKPCA, and KLDA, we apply a modification of our approach from Section 5 to a smaller face dataset, the face and gesture recognition network (FG-NET). FG-NET is a popular, publicly available database used for age estimation, gender classification, face recognition, and other demographic analysis tasks [80]. It contains 1002 images from 82 subjects: 47 males and 35 females with ages varying from 0 to 69 years [80].

For each image, 109 features are extracted using the Active Appearance Model (AAM), a 383 commonly adopted appearance-based approach that models the shape and texture of the face [73,81]. 384 As in Section 5.4, the radial kernel defined in equation (24) is chosen for each of the DR methods KPCA, 385 SKPCA, and KLDA. Additionally, the modified link function from equation (25) is applied in the SKPCA algorithm. Thus, the tuning parameter δ in the radial kernel and η in the modified link function must be selected. As in our experiments on Morph-II, linear SVM is chosen as the classifier for FG-NET. 388 On Morph-II, values of the cost parameter c ranging from 10^{-8} to 10^{8} were tested. On FG-NET, we 389 have observed convergence issues in the SVM algorithm for values of c exceeding 10, so only the 390 values 10^{-8} , 10^{-7} , ..., 10^{-1} , 1, 10 are tested. The considered tuning parameters are summarized in 39: Table 7. 302 For cross-validation, we use leave-one-person-out (LOPO), the most well-accepted scheme for

For cross-validation, we use leave-one-person-out (LOPO), the most well-accepted scheme for FG-NET [80]. LOPO is a variation of *k*-fold cross-validation that produces independent training and testing folds in longitudinal datasets. The number of folds *k* is set equal to the number of subjects in the dataset, so k = 82 here. For i = 1, 2, ... 82, testing fold *i* contains only images of person *i*, while training fold *i* contains all remaining images. Similarly to on Morph-II, we choose the mean classification

³⁰⁸ accuracy over the testing folds to be the evaluation criterion.

	КРСА	$\delta = 3.2, 3.2\overline{6}, 3.\overline{3}, 3.4, 3.4\overline{6}, 3.5\overline{3}, 3.6, 3.\overline{6}, 3.7\overline{3}, 3.8$
Dimension Reduction	SKPCA	$\delta = 0.0098$
		$\eta = 0.001, 0.01, 0.1, 1$
	KLDA	$\delta = 3, 3.\overline{5}, 4.\overline{1}, 4.\overline{6}, 5.\overline{2}, 5.\overline{7}, 6.\overline{3}, 6.\overline{8}, 7.\overline{4}, 8$
Classifier	Linear SVM	$c = 10^{-8}, \dots, 10^{-1}, 1, 10$

Table 8. Gender Classification Results on FG-NET

Table 7. Parameter Summary for FG-NET

			(=)	(2)
Method	Parameters	Acc ⁽¹⁾	TPR ⁽²⁾	TNR ⁽³⁾
	$\delta = 3.2\overline{6}, c = 10$	0.7025	0.7325	0.6621
KPCA	$\delta = 3.\overline{3}, c = 10$	0.6932	0.7233	0.6528
	$\delta = 3.4, c = 10$	0.6801	0.6651	0.7001
	$\delta = 0.0098, \eta = 0.1, c = 10$	0.7154	0.7542	0.6633
SKPCA	$\delta = 0.0098, \eta = 1, c = 0.1$	0.6933	0.7413	0.6289
	$\delta = 0.0098, \eta = 0.01, c = 0.1$	0.6893	0.7701	0.5809
	$\delta = 3, c = 0.01$	0.7225	0.7593	0.6730
KLDA	$\delta = 5.\overline{7}, c = 1$	0.7176	0.7810	0.6324
	$\delta = 8, c = 0.1$	0.7131	0.7431	0.6727

(1) Acc represents mean accuracy.

(2) TPR represents mean true positive rate (recall/sensitivity): the proportion of male faces correctly classified.

(3) TNR represents mean true negative rate (specificity): the proportion of female faces correctly classified.



Figure 6. Receiver operating characteristic (ROC) curve and area under the curve (AUC) are compared by method for gender classification on FG-NET. Each color corresponds to a DR method.

For each fold, we transform and reduce the dimension of the features through each DR method. In all cases, a dimension of 100 is retained to facilitate comparison with the results on Morph-II. The

transformed, dimension-reduced features then predict the gender of the testing fold's images through

- a linear SVM. The predicted classes from SVM are also mapped to probabilities through [78], similarly
- as in Section 6. The gender classification accuracy is calculated for the testing fold. Finally, all such

testing classification accuracies are averaged to compute the mean classification accuracy from testing;
 the testing probabilities are used to form ROC curves.

The optimum gender classification results on FG-NET are presented in Table 8. The maximum classification accuracy of about 72.25% is achieved by KLDA. For other choices of parameters, KLDA 407 reaches above 71% accuracy, which is close to the maximum accuracy attained by SKPCA. Meanwhile, 408 the peak accuracy reached by KPCA is 70.25%. In general here, KLDA is observed to outperform 409 SKPCA and KPCA, while SKPCA tends to surpass KPCA. In most cases, the probability of correctly 410 classifying males (sensitivity/true positive rate) is higher than the probability of correctly classifying females (specificity/true negative rate). For each DR method, an ROC curve (corresponding to the 412 results from Table 8 with maximal mean classification accuracy) is displayed in Figure 6. The area 413 under the curve (AUC) is highest for KLDA, followed by KPCA then SKPCA. 414

Overall, the gender classification results on Morph-II are stronger than on FG-NET. Lower 415 accuracy on FG-NET could be caused by the greater number of minors (aged 0-18), who have been 416 more difficult to classify than adults in some studies [35,82]. Additionally, there are substantially fewer 417 faces for training in FG-NET versus Morph-II (under 1000 versus 10280 images). Another contributor 418 could be the choice of features and its dimension; the AAM features have dimension 109 on FG-NET, 419 while the HOG, LBP, and BIF features have dimensions ranging from 500 to thousands on Morph-II. 420 SKPCA reaches peak performance on Morph-II, while KLDA attains optimal results on FG-NET. 421 However, the results on Morph-II and FG-NET are similar in that the supervised methods KLDA and 422 SKPCA outperform the unsupervised method KPCA for gender classification. Further, both datasets 423 evidence that female faces are more challenging to classify than male faces. 424

425 8. Computational Framework for Practical Systems

To tackle the challenges of high dimensionality and intensive computation for large-scale databases (like Morph-II, as shown in the Time column of Table 5) in real-world applications, we propose a computational framework to substantially decrease runtime.

Our approach involves parallel computing, the bootstrap resampling method, and ensemble 429 learning. Let M1 denote the main training set and M2 the testing set. If M1 is very large, we can save 430 some time by drawing bootstrapped samples from M1. Let S_i denote the *i*th bootstrapped sample 431 from M1. Send S_i to a core (or processor), Core *i*. Train the model on S_i . Test on the full testing set 432 *M*2, obtaining a set of gender predictions corresponding to Core *i* and S_i . Repeat this process for all 433 bootstrapped samples and corresponding cores *i*. The final predictions are obtained by taking the 434 majority rule of the predictions from all *i* cores and samples. Hence, the results from this scheme 435 approximate the results from the full Morph-II. This framework is summarized in Figure 7. 436



Figure 7. Flowchart representing the parallel computational framework for practical systems proposed for Morph-II and other large datasets.

To explore the effectiveness, this framework is applied to Morph-II with a selection of BIF, LBP, and HOG features as preliminary studies. This experiment is implemented through the HiPerGator 2.0 supercomputer at University of Florida with five cores per combination of feature and dimension reduction method. Following the subsetting scheme discussed in Section 5.2, for simplicity, we consider only the case of bootstrapping image samples from S_1 for training, while each image from $S_2 \cup R$ is used for testing.

Method	Feature	Accuracy	Memory (gb)	Time (min)
KPCA	BIF $s = s7 - 37$, $\gamma = 0.4$	0.9330	27.59	90
KICA	HOG $s = 12, o = 8$	0.9178	29.77	101
	LBP $s = 10, r = 1$	0.8927	25.85	37
SKPCA	BIF $s = s7 - 37$, $\gamma = 0.4$	0.9417	53.28	89
SKICA	HOG $s = 12, o = 8$	0.9056	51.43	74
	LBP $s = 10, r = 1$	0.9274	20.33	24
KIDA	BIF $s = s7 - 37$, $\gamma = 0.4$	0.9416	30.99	100
KLDA	HOG $s = 12, o = 8$	0.9133	25.42	102
	LBP $s = 10, r = 1$	0.9118	17.05	26

Table 9. Classification Results Based on Bootstrapping

We evaluate this framework by comparing the approximated results in Table 9 to the results from 443 Table 5. For each combination of feature and dimension reduction method, each of the five cores 444 independently trains a bootstrapped sample of 1000 images from S_1 and tests on $S_2 \cup R$. Then the 445 gender predictions over all five cores are compared with a simple majority rule; e.g., if an image is 446 predicted male for three images and female for two images, the final gender prediction is male. The 447 times in Table 9 are the total runtimes for this process, which include training and testing on HPC. 448 Therefore, the times and memory can be compared between Tables 5 and 9. A distinction is that in 449 Table 5, results are averaged for the alternating scheme, while in Table 9, the results are only from 450 when S_1 is used for training and $S_2 \cup R$ for testing. 451

It is shown in Table 9 that, in many cases, the accuracy rates from the approximations are similar to those from the main approach in Table 5. This is a very good result, especially considering that the bootstrapping approach uses no more than 5000 images total for training, while the main approach used all 10280 images for training. This finding suggests that our methods may perform reasonably well on Morph-II with smaller training sets. The most substantial difference between the ⁴⁵⁷ bootstrapped approach and the main approach is in the runtime. For all combinations of features
⁴⁵⁸ and dimension reduction methods, the bootstrapping approach has decreased the runtime to under
⁴⁵⁹ two hours. Meanwhile, the main approach in Table 5 yields runtimes exceeding 20 hours. Hence, our
⁴⁶⁰ preliminary results indicate the parallel approximation approach can attain similar accuracy rates to
⁴⁶¹ the main approach, while substantially saving time. Such a result is promising for practical gender
⁴⁶² classification systems, where gender predictions must be made in real-time.

463 9. Conclusion

We have performed a comparative study of the nonlinear dimension reduction methods KPCA, SKPCA, and KLDA. These kernel-based methods are first applied to three simulated datasets for visualization and comparison. SKPCA and KLDA outperform KPCA, reinforcing the need for supervised approaches in classification tasks. The radial kernel performed well, encouraging its use for face analysis.

Next, we have proposed and evaluated a new machine learning process for Morph-II. First, 469 we use a novel subsetting scheme that reduces class imbalances while establishing independence 470 between training and testing sets. Then we preprocess Morph-II photographs and extract three 471 appearance-based features: HOG, LBP, and BIF. We transform and reduce the dimension of these 472 features through KPCA, SKPCA, and KLDA. Linear SVM classifies the gender of Morph-II subjects, 473 reaching accuracy rates of 95%. With promising preliminary results on Morph-II, a practical 474 computational framework is offered that reduces runtime through parallelization and approximation. 475 The performance of the dimension reduction methods are further compared through an 476 application to the FG-NET dataset. Images are represented through the appearance-based AAM 477 features; transformed and reduced in dimension through KPCA, SKPCA, and KLDA; and classified as containing a male or female subject through linear SVM. While SKPCA performed optimally 479 on Morph-II, KLDA reached top performance on FG-NET with 72% leave-one-person-out (LOPO) 480 accuracy. 481

Further directions of research involve automatic tuning parameter selection, reduction of 482 computational cost, and application to other face analysis tasks. Our approach could yield improved 483 results with better choices of parameters, but it is impossible to anticipate and try all combinations. Automatic parameter selection for kernels could help identify a good set of parameters more easily. 485 Perhaps the most important future direction of research on Morph-II is to reduce computational cost. 486 For many practical demographic analysis systems, predictions must be made in real-time. For our 487 gender classification methods, our parallel approximation approach substantially reduced runtime 488 while attaining similar accuracy rates to the main approach. Such computational strategies should 489 be further investigated to help bring gender classification and other face analysis tasks to practical 490 implementation. Finally, our machine learning pipeline for Morph-II could be generalized to race 491 classification or even age estimation. 492

493 10. Acknowledgments

This material is based in part upon work supported by the National Science Foundation under Grant Numbers DMS-1659288. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation. The authors would like to thank the reviewers for the helpful comments that significantly improves the presentation of the paper.

499

- Fodor, I.K. A survey of dimension reduction techniques. Technical report, Lawrence Livermore National
 Lab., CA (US), 2002.
- Sorzano, C.O.S.; Vargas, J.; Montano, A.P. A survey of dimensionality reduction techniques. *arXiv preprint arXiv:1403.2877* 2014.

504 505	3.	Izenman, A.J. Modern multivariate statistical techniques. <i>Regression, classification and manifold learning</i> 2008 .
506	4.	Pearson, K. LIII. On lines and planes of closest fit to systems of points in space. <i>The London, Edinburgh, and</i>
507		Dublin Philosophical Magazine and Journal of Science 1901 , 2, 559–572.
508	5.	Hotelling, H. Analysis of a complex of statistical variables into principal components. <i>Journal of educational</i>
500	0.	nsuchology 1933 24 417
509	6	Yang L: Peng W: Ward M Ω : Rundensteiner F A. Interactive hierarchical dimension ordering spacing
510	0.	and filtering for exploration of high dimensional datasets. IEEE Symposium on Information Visualization
511		2003 (IEEE Cat. No. 03TH8714) IEEE 2003 np. 105-112
512	7	Johansson S: Johansson I. Interactive dimensionality reduction through user-defined combinations of
513	7.	quality metrics. <i>IEEE transactions on visualization and computer graphics</i> 2009 , <i>15</i> , 993–1000.
515	8.	Fisher, R.A. The use of multiple measurements in taxonomic problems. Annals of human genetics 1936,
516		7, 179–188.
517	9.	Rao, C.R. The utilization of multiple measurements in problems of biological classification. Journal of the
518		Royal Statistical Society. Series B (Methodological) 1948, 10, 159–203.
519	10.	Lee, J.A.; Verleysen, M. Nonlinear dimensionality reduction; Springer Science & Business Media, 2007.
520	11.	Nhan Duong, C.; Luu, K.; Gia Quach, K.; Bui, T.D. Beyond principal components: Deep boltzmann
521		machines for face modeling. Proceedings of the IEEE Conference on Computer Vision and Pattern
522		Recognition, 2015, pp. 4786–4794.
523	12.	Yin, J.; Liu, Z.; Jin, Z.; Yang, W. Kernel sparse representation based classification. <i>Neurocomputing</i> 2012 , 77, 120, 129
524	10	//, 120–120.
525	13.	Shawe-Taylor, J.; Cristianini, IN. <i>Kernet methods for pattern analysis</i> ; Cambridge university press, 2004.
526	14.	Motal, 1. Kernel association for classification and prediction: A survey. <i>IEEE transactions on neural networks</i>
527	15	Via V. Luu V. Souvidos M. A robust approach to facial attraitive description on large scale face
528	15.	Ate, 1., Luu, K., Savvides, M. A robust approach to factal enflucty classification of farge scale face
529		on IEEE 2012 no. 142, 140
530	16	Schölkonf B. Smola, A. Müller K.P. Kernel principal component analysis. International Conference on
531	10.	Artificial Neural Networks, Springer 1007, pp. 592–598
532	17	Mika S : Ratech C : Waston I : Scholkonf B : Mullars K P. Eishar discriminant analysis with karnals
533	17.	Noural networks for signal processing IV 1000 Proceedings of the 1000 IEEE signal processing society
534		workshop Jose 1000 pp 41 48
535	18	Baudat C : Anouar E Congratized discriminant analysis using a kernel approach. Neural commutation 2000
530	10.	12 2385 2404
537	10	12, 2007–2404. Barzilay O : Brailovsky VI. On domain knowledge and feature selection using a support vector machine
538	17.	Pattern Recognition Letters 1999 , 20, 475–484.
540	20.	Schölkopf, B.; Simard, P.; Smola, A.I.; Vapnik, V. Prior knowledge in support vector kernels. Advances in
541		neural information processing systems, 1998, pp. 640–646.
542	21.	Hinton, G.E.; Salakhutdinov, R.R. Reducing the dimensionality of data with neural networks. <i>science</i> 2006 ,
543		313, 504–507.
544	22.	Zhao, W.; Krishnaswamy, A.; Chellappa, R.; Swets, D.L.; Weng, J. Discriminant analysis of principal
545		components for face recognition. In <i>Face Recognition</i> ; Springer, 1998; pp. 73–85.
546	23.	Martínez, A.M.; Kak, A.C. Pca versus lda. <i>IEEE transactions on pattern analysis and machine intelligence</i> 2001 ,
547		23, 228–233.
548	24.	Yang, J.; Yang, J.y. Why can LDA be performed in PCA transformed space? Pattern recognition 2003,
549		36, 563–566.
550	25.	Turk, M.; Pentland, A. Eigenfaces for recognition. Journal of cognitive neuroscience 1991, 3, 71–86.
551	26.	Belhumeur, P.N.; Hespanha, J.P.; Kriegman, D.J. Eigenfaces vs. fisherfaces: Recognition using class specific
552		linear projection. <i>IEEE Transactions on pattern analysis and machine intelligence</i> 1997 , <i>19</i> , 711–720.
553	27.	Kim, K.I.; Jung, K.; Kim, H.J. Face recognition using kernel principal component analysis. IEEE signal
554		processing letters 2002 , 9, 40–42.
555	28.	Lu, J.; Plataniotis, K.N.; Venetsanopoulos, A.N. Face recognition using kernel direct discriminant analysis
556		algorithms. IEEE Transactions on Neural Networks 2003, 14, 117–126.

- Yang, J.; Zhang, D.; Frangi, A.F.; Yang, J.y. Two-dimensional PCA: a new approach to appearance-based face representation and recognition. *IEEE transactions on pattern analysis and machine intelligence* 2004, 26, 131–137.
- 30. Li, M.; Yuan, B. 2D-LDA: A statistical linear discriminant analysis for image matrix. *Pattern Recognition* Letters 2005, 26, 527–532.
- Karg, M.; Jenke, R.; Seiberl, W.; Kühnlenz, K.; Schwirtz, A.; Buss, M. A comparison of PCA, KPCA and
 LDA for feature extraction to recognize affect in gait kinematics. Affective computing and intelligent
 interaction and workshops, 2009. ACII 2009. 3rd international conference on. IEEE, 2009, pp. 1–6.
- Ye, F.; Shi, Z.; Shi, Z. A comparative study of PCA, LDA and Kernel LDA for image classification.
 Ubiquitous Virtual Reality, 2009. ISUVR'09. International Symposium on. IEEE, 2009, pp. 51–54.
- Yang, J.; Jin, Z.; Yang, J.y.; Zhang, D.; Frangi, A.F. Essence of kernel Fisher discriminant: KPCA plus LDA.
 Pattern Recognition 2004, 37, 2097–2100.
- Barshan, E.; Ghodsi, A.; Azimifar, Z.; Jahromi, M.Z. Supervised principal component analysis:
 Visualization, classification and regression on subspaces and submanifolds. *Pattern Recognition* 2011, 44, 1357–1371.
- ⁵⁷² 35. Wang, Y.; Chen, C.; Watkins, V.; Ricanek, K. Modified Supervised Kernel PCA for Gender Classification.
 ⁵⁷³ International Conference on Intelligent Science and Big Data Engineering. Springer, 2015, pp. 60–71.
- Fewzee, P.; Karray, F. Dimensionality reduction for emotional speech recognition. Privacy, Security,
 Risk and Trust (PASSAT), 2012 International Conference on and 2012 International Conference on Social
 Computing (SocialCom). IEEE, 2012, pp. 532–537.
- Samadani, A.A.; Ghodsi, A.; Kulić, D. Discriminative functional analysis of human movements. *Pattern Recognition Letters* 2013, 34, 1829–1839.
- Wu, H.; Bowers, D.M.; Huynh, T.T.; Souvenir, R. Biomedical video denoising using supervised manifold
 learning. Biomedical Imaging (ISBI), 2013 IEEE 10th International Symposium on. IEEE, 2013, pp.
 1244–1247.
- Ashtiani, H.; Ghodsi, A. A dimension-independent generalization bound for kernel supervised principal
 component analysis. Feature Extraction: Modern Questions and Challenges, 2015, pp. 19–29.
- Fu, Y.; Guo, G.; Huang, T.S. Age synthesis and estimation via faces: A survey. *IEEE transactions on pattern analysis and machine intelligence* 2010, 32, 1955–1976.
- Sun, Y.; Zhang, M.; Sun, Z.; Tan, T. Demographic analysis from biometric data: Achievements, challenges,
 and new frontiers. *IEEE transactions on pattern analysis and machine intelligence* 2018, 40, 332–351.
- Burton, A.M.; Bruce, V.; Dench, N. What's the difference between men and women? Evidence from facial
 measurement. *Perception* 1993, 22, 153–176.
- ⁵⁹⁰ 43. Ng, C.B.; Tay, Y.H.; Goi, B.M. Vision-based human gender recognition: A survey. *arXiv preprint* ⁵⁹¹ *arXiv:*1204.1611 2012.
- Golomb, B.A.; Lawrence, D.T.; Sejnowski, T.J. Sexnet: A neural network identifies sex from human faces.
 NIPS, 1990, Vol. 1, p. 2.
- 45. Cottrell, G.W.; Metcalfe, J. EMPATH: Face, emotion, and gender recognition using holons. Advances in neural information processing systems, 1991, pp. 564–571.
- 46. Poggio, B.; Brunelli, R.; Poggio, T. HyberBF networks for gender classification, 1992.
- 47. Wiskott, L.; Fellous, J.M.; Krüger, N.; Von der Malsburg, C. Face recognition and gender determination,
 1995.
- 48. Guo, G.; Mu, G. Human age estimation: What is the influence across race and gender? Computer Vision
 and Pattern Recognition Workshops (CVPRW), 2010 IEEE Computer Society Conference on. IEEE, 2010,
 pp. 71–78.
- Guo, G.; Mu, G. Simultaneous dimensionality reduction and human age estimation via kernel partial least
 squares regression. Computer vision and pattern recognition (cvpr), 2011 ieee conference on. IEEE, 2011,
 pp. 657–664.
- 50. Shan, C. Learning local binary patterns for gender classification on real-world face images. *Pattern recognition letters* 2012, 33, 431–437.
- 51. Yang, H.F.; Lin, B.Y.; Chang, K.Y.; Chen, C.S. Automatic age estimation from face images via deep ranking.
 networks 2013, *35*, 1872–1886.

609	52.	Yi, D.; Lei, Z.; Li, S.Z. Age estimation by multi-scale convolutional network. Asian conference on computer
610		vision. Springer, 2014, pp. 144–158.
611	53.	Antipov, G.; Berrani, S.A.; Dugelay, J.L. Minimalistic CNN-based ensemble model for gender prediction
612		from face images. Pattern recognition letters 2016 , 70, 59–65.
613	54.	Antipov, G.; Baccouche, M.; Berrani, S.A.; Dugelay, J.L. Effective training of convolutional neural networks
614		for face-based gender and age prediction. <i>Pattern Recognition</i> 2017 , 72, 15–26.
615	55.	Yang, Z.; Ai, H. Demographic classification with local binary patterns. International Conference on
616		Biometrics. Springer, 2007, pp. 464–473.
617	56.	Lian, H.C.; Lu, B.L. Multi-view gender classification using local binary patterns and support vector
618		machines. International Symposium on Neural Networks. Springer, 2006, pp. 202–209.
619	57.	Mäkinen, E.; Raisamo, R. An experimental comparison of gender classification methods. pattern recognition
620		<i>letters</i> 2008 , <i>29</i> , 1544–1556.
621	58.	Alexandre, L.A. Gender recognition: A multiscale decision fusion approach. Pattern recognition letters 2010,
622		31, 1422–1427.
623	59.	Xia, B.; Sun, H.; Lu, B.L. Multi-view gender classification based on local Gabor binary mapping pattern and
624		support vector machines. Neural Networks, 2008. IJCNN 2008.(IEEE World Congress on Computational
625		Intelligence). IEEE International Joint Conference on. IEEE, 2008, pp. 3388–3395.
626	60.	Guo, G.; Dyer, C.R.; Fu, Y.; Huang, T.S. Is gender recognition affected by age? Computer Vision Workshops
627		(ICCV Workshops), 2009 IEEE 12th International Conference on. IEEE, 2009, pp. 2032–2039.
628	61.	Han, H.; Otto, C.; Liu, X.; Jain, A.K. Demographic estimation from face images: Human vs. machine
629		performance. IEEE Transactions on Pattern Analysis & Machine Intelligence 2015 , pp. 1148–1161.
630	62.	Ma, Y.: Fu, Y. Manifold learning theory and applications: CRC press, 2011.
631	63.	Schölkopf, B.: Herbrich, R.: Smola, A.I. A generalized representer theorem. International conference on
632		computational learning theory. Springer, 2001, pp. 416–426.
633	64	Wabba G. Spline models for observational data: Vol. 59 Siam 1990
624	65	Schölkopf B · Smola A · Müller K R Nonlinear component analysis as a kernel eigenvalue problem
625	00.	Neural computation 1998 10 1299–1319
636	66	Cretton A Bousquet O Smola A Schölkonf B Measuring statistical dependence with Hilbert-Schmidt
030	00.	norms. International conference on algorithmic learning theory. Springer 2005, pp. 63–77
037	67	Stoinwart I: Christmann A Summart machines: Springer Science & Business Modia 2008
038	68	Bicanok K : Tocafavo T. Mornh: A longitudinal image database of normal adult age progression
039	00.	Automatic Ease and Costure Recognition 2006 ECP 2006 7th International Conference on IEEE
640		2006 pp. 241–245
641	60	2000, pp. 541–545.
642	69.	Japkowicz, N.; Stephen, S. The class initialance problem: A systematic study. <i>Intelligent untu unutysis</i> 2002,
643	70	0,429-449.
644	70.	Guo, G.; Mu, G. A study of large-scale ethnicity estimation with gender and age variations. Computer
645		Vision and Pattern Recognition Workshops (CVPRW), 2010 IEEE Computer Society Conference on. IEEE,
646		2010, pp. 79–86.
647	/1.	Tip, B.; Bingnam, G.; Kempfert, K.; Fabish, J.; Kling, T.; Chen, C.; Wang, Y. Preliminary Studies on a Large
648	70	Face Database. arXiv preprint arXiv:1811.06446 2018.
649	72.	Yip, B.; Towner, R.; Kling, T.; Chen, C.; Wang, Y. Image Pre-processing Using OpenCV Library on
650		MORPH-II Face Database. arXiv preprint arXiv:1811.06934 2018.
651	73.	Edwards, G.J.; Taylor, C.J.; Cootes, T.F. Interpreting face images using active appearance models
652		Proceedings Third IEEE International Conference on Automatic Face and Gesture Recognition. IEEE, 1998,
653		pp. 300–305.
654	74.	Kling, T. Morph-II: Feature Vector Documentation: NSF-REU site at UNC Wilmington. http://libres.uncg
655		edu/ir/uncw/f/wangy2018-1.pdf, 2017.
656	75.	Byun, H.; Lee, S.W. Applications of support vector machines for pattern recognition: A survey. In Pattern
657		recognition with support vector machines; Springer, 2002; pp. 213–236.
658	76.	Qiu, Y.; Mei, J.; Qiu, M.Y. Package 'rARPACK' 2016 .
659	77.	Dimitriadou, E.; Hornik, K.; Leisch, F.; Meyer, D.; Weingessel, A. Misc Functions of the Department of
660		Statistics (e1071), TU Wien. <i>R package version</i> 2005 , pp. 1–5.

- 78. Platt, J.; others. Probabilistic outputs for support vector machines and comparisons to regularized
 likelihood methods. *Advances in large margin classifiers* 1999, 10, 61–74.
- ⁶⁶³ 79. Guo, G.; Mu, G. A framework for joint estimation of age, gender and ethnicity on a large database. *Image* ⁶⁶⁴ and Vision Computing 2014, 32, 761–770.
- Panis, G.; Lanitis, A.; Tsapatsoulis, N.; Cootes, T.F. Overview of research on facial ageing using the FG-NET
 ageing database. *IET Biometrics* 2016, *5*, 37–46.
- ⁶⁶⁷ 81. Cootes, T.F.; Edwards, G.J.; Taylor, C.J. Active appearance models. European conference on computer
 ⁶⁶⁸ vision. Springer, 1998, pp. 484–498.
- Wang, Y.; Ricanek, K.; Chen, C.; Chang, Y. Gender classification from infants to seniors. 2010 Fourth IEEE
 International Conference on Biometrics: Theory, Applications and Systems (BTAS). IEEE, 2010, pp. 1–6.

© 2019 by the authors. Submitted to *Intelligent Data Analysis* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (http://creativecommons.org/licenses/by/4.0/).