

## SPARSE LOGISTIC PRINCIPAL COMPONENTS ANALYSIS FOR BINARY DATA

BY SEOKHO LEE, JIANHUA Z. HUANG<sup>1</sup> AND JIANHUA HU<sup>2</sup>

*Harvard School of Public Health, Texas A&M University  
 and University of Texas M. D. Anderson Cancer Center*

We develop a new principal components analysis (PCA) type dimension reduction method for binary data. Different from the standard PCA which is defined on the observed data, the proposed PCA is defined on the logit transform of the success probabilities of the binary observations. Sparsity is introduced to the principal component (PC) loading vectors for enhanced interpretability and more stable extraction of the principal components. Our sparse PCA is formulated as solving an optimization problem with a criterion function motivated from a penalized Bernoulli likelihood. A Majorization–Minimization algorithm is developed to efficiently solve the optimization problem. The effectiveness of the proposed sparse logistic PCA method is illustrated by application to a single nucleotide polymorphism data set and a simulation study.

**1. Introduction.** Principal components analysis (PCA) is a widely used method for dimensionality reduction, feature extraction and visualization of multivariate data. Several sparse PCA methods have recently been introduced to improve the standard PCA [e.g., Jolliffe, Trendafilov and Uddine (2003); Zou, Hastie and Tibshirani (2006); Shen and Huang (2008)]. By requiring the principal component loading vectors to be sparse, sparse PCA methods yield PCs that are more easily interpretable. Sparsity also regularizes the extraction of PCs and thus makes the extraction more stable. Such stability is much desired when the dimension is high, especially in the so-called high-dimension low-sample-size settings. As extensions of the standard PCA, however, these sparse PCA methods are mostly suitable to variables of a continuous type, they are not generally appropriate for other data types such as binary data or counts. Although the basic objective of PCA, or its sparse version, can be achieved regardless of the nature of the original

---

Received January 2009; revised January 2010.

<sup>1</sup>Supported in part by Grants from the National Science Foundation (DMS-06-06580, DMS-09-07170), the National Cancer Institute (CA57030), the Virtual Center for Collaboration between Statisticians in the US and China, and King Abdullah University of Science and Technology (KAUST, Award KUS-CI-016-04).

<sup>2</sup>Supported in part by Grants from the National Science Foundation (DMS-07-06818) and the National Institute of Health (R01-RGM080503A, R21-CA129671).

*Key words and phrases.* Binary data, dimension reduction, MM algorithm, LASSO, PCA, regularization, sparsity.

variable, it is true that variances and covariances have especial relevance for multivariate Gaussian variables, and that linear functions of binary variables are less readily interpretable than linear functions of continuous variables [Jolliffe (2002)]. The goal of this paper is to develop a sparse PCA method for binary data.

There are two commonly used definitions of PCA that give rise to the same result. PCA can be defined by finding the orthogonal projection of the data onto a low dimensional linear subspace such that the variance of the projected data is maximized [Hotelling (1933)]. Alternatively, PCA can also be defined by finding the linear projection that minimizes the mean squared distance between the data points and their projections [Pearson (1901)]. Shen and Huang (2008) developed their sparse PCA method following the viewpoint of Pearson. Suppose  $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^d$  are the  $n$  data points and consider a  $k$ -dimensional ( $k < d$ ) linear manifold spanned by a bases  $\tilde{\mathbf{b}}_1, \dots, \tilde{\mathbf{b}}_k$  with a shift vector  $\boldsymbol{\mu}$ . According to Pearson, the PCA minimizes the following reconstruction error,

$$(1.1) \quad \sum_{i=1}^n \|\mathbf{y}_i - (\boldsymbol{\mu} + a_{i1}\tilde{\mathbf{b}}_1 + \dots + a_{ik}\tilde{\mathbf{b}}_k)\|^2,$$

subject to the constraint that  $\mathbf{A} = (a_{ij})$  has orthonormal columns. Usually the variables presented in  $\mathbf{y}_i$  are scaled so that they have the same order of magnitude. Note that (1.1) is a least squares regression if  $a_{ik}$ 's were known. In light of this connection to regression and borrowing the idea from LASSO [Tibshirani (1996)], Shen and Huang (2008) proposed to add an  $L_1$  penalty  $\|\tilde{\mathbf{b}}_1\|_1 + \dots + \|\tilde{\mathbf{b}}_k\|_1$  to the reconstruction error (1.1) to obtain sparse loading vectors  $\tilde{\mathbf{b}}_1, \dots, \tilde{\mathbf{b}}_k$ . Since the reconstruction error (1.1) can be viewed as the negative log likelihood up to a constant for the Gaussian distributions with mean vectors  $\boldsymbol{\theta}_i = \boldsymbol{\mu} + a_{i1}\tilde{\mathbf{b}}_1 + \dots + a_{ik}\tilde{\mathbf{b}}_k$  for  $i = 1, \dots, n$  and identity covariance, the method of Shen and Huang can be interpreted as a penalized likelihood approach for the sparse PCA. The key idea of the current paper is to replace the Gaussian likelihood by the Bernoulli likelihood where  $\boldsymbol{\theta}_i$  will be the logit transform of the success probabilities. We refer to the proposed PCA method as sparse logistic PCA. The relationship of the proposed sparse logistic PCA to the sparse PCA of Shen and Huang is analogous to the relationship between logistic and linear LASSO regression.

We develop an iterative weighted least squares algorithm to perform the proposed sparse logistic PCA. Since the log Bernoulli likelihood is not quadratic and the  $L_1$  penalty function is nondifferentiable, the optimization problem defining the sparse logistic PCA is not straightforward to solve. Our algorithm applies the general idea of optimization transfer or Majorization–Minimization (MM) algorithm [Lange, Hunter and Yang (2000); Hunter and Lange (2004)]. By iteratively replacing the complex objective function with suitably defined quadratic surrogates, each step of our algorithm solves a weighted least squares problem and has closed form. The algorithm is easy to implement and guaranteed at each iteration to improve the penalized PCA log-likelihood. We show that the same MM algorithm

1 is applicable when there are missing data. We also develop a method for choos- 1  
 2 ing the penalty parameters and for choosing the number of important principal 2  
 3 components. PCA of binary data using Bernoulli likelihood has previously been 3  
 4 studied by Collins, Dasgupta and Schapire (2001), Schein, Saul and Ungar (2003) 4  
 5 and de Leeuw (2006), but none of these works considered sparse loading vectors. 5  
 6 As we demonstrate using simulation and real data, sparsity can enhance interpre- 6  
 7 tation of results and improve the stability and accuracy of the extracted principal 7  
 8 components. 8

9 Other approaches of sparse PCA are not as easily extendible to binary data. 9  
 10 Jolliffe, Trendafilov and Uddine (2003) modified the defining maximum variance 10  
 11 problem of the standard PCA by applying an  $L_1$ -norm constraint on the PC load- 11  
 12 ing vectors to obtain PCA with sparse loadings. Its use of sample variance makes 12  
 13 it unappealing for binary data. Zou, Hastie and Tibshirani (2006) rewrote PCA as 13  
 14 a regression-type optimization problem and then applied the LASSO penalty [Tib- 14  
 15 shirani (1996)] to obtain sparse loadings. However, since the data appear both as 15  
 16 regressors and responses in their regression-type problem, the connection of their 16  
 17 approach to the penalized likelihood is not as natural as Shen and Huang (2008). 17

18 The rest of this article is organized as follows. In Section 2 we introduce the 18  
 19 optimization problem that yields the sparse logistic PCA and provides methods 19  
 20 for tuning parameter selection. Section 3 applies the sparse logistic PCA to a sin- 20  
 21 gle nucleotide polymorphism data set and compares it with the nonsparse version 21  
 22 of logistic PCA. Section 4 presents a Majorization–Minimization algorithm for 22  
 23 efficient computation of the sparse logistic PCA and Section 5 discusses how to 23  
 24 handle missing data. Results of a simulation study are given in Section 6. Sec- 24  
 25 tion 7 concludes the paper with some discussion. An Appendix contains proofs of 25  
 26 theorems. 26

## 27 2. Sparse logistic PCA with penalized likelihood. 27

28  
 29 2.1. *Penalized Bernoulli likelihood.* Consider the  $n \times d$  binary data matrix 29  
 30  $\mathbf{Y} = (y_{ij})$ , each row of which represents a vector of observations from binary 30  
 31 variables. We assume that entries of  $\mathbf{Y}$  are realizations of mutually independent 31  
 32 random variables and that  $y_{ij}$  follows the Bernoulli distribution with success prob- 32  
 33 ability  $\pi_{ij}$ . Let  $\theta_{ij} = \log\{\pi_{ij}/(1 - \pi_{ij})\}$  be the logit transformation of  $\pi_{ij}$ . Define 33  
 34 the inverse logit transformation  $\pi(\theta) = \{1 + \exp(-\theta)\}^{-1}$ . Then the success prob- 34  
 35 abilities can be represented using the canonical parameters as  $\pi_{ij} = \pi(\theta_{ij})$ . The 35  
 36 individual data generating probability becomes 36

$$37 \Pr(Y_{ij} = y_{ij}) = \pi(\theta_{ij})^{y_{ij}} \{1 - \pi(\theta_{ij})\}^{1-y_{ij}} = \pi(q_{ij}\theta_{ij}), 37$$

38 with  $q_{ij} = 2y_{ij} - 1$  since  $\pi(-\theta) = 1 - \pi(\theta)$ . This representation leads to the 38  
 39 compact form of the log likelihood as 39  
 40

$$41 (2.1) \quad \ell = \sum_{i=1}^n \sum_{j=1}^d \log \pi(q_{ij}\theta_{ij}). 41$$

1 Note that the Bernoulli distributions are in the exponential family and  $\theta_{ij}$  are the 1  
2 corresponding canonical parameters. 2

3 To build a probabilistic model for principal components analysis of binary 3  
4 data, the  $d$ -dimensional canonical parameter vectors  $\theta_i = (\theta_{i1}, \dots, \theta_{id})^T$  are 4  
5 constrained to reside in a low dimensional manifold of  $\mathbb{R}^d$  with the dimensionality 5  
6  $k$ . (The choice of  $k$  will be discussed later in Section 2.3.) Specifically, we 6  
7 assume that, for some vectors  $\mu, \tilde{\mathbf{b}}_1, \dots, \tilde{\mathbf{b}}_k \in \mathbb{R}^d$ , the vector of canonical para- 7  
8 meters satisfies  $\theta_i = \mu + a_{i1}\tilde{\mathbf{b}}_1 + \dots + a_{ik}\tilde{\mathbf{b}}_k$  for  $i = 1, \dots, n$ . We call  $\tilde{\mathbf{b}}_1, \dots, \tilde{\mathbf{b}}_k$  8  
9 the principal component loading vectors and the coefficients  $\mathbf{a}_i = (a_{i1}, \dots, a_{ik})^T$  9  
10 the principal component scores (PC scores) for the  $i$ th observation. Geometrically, 10  
11 the vectors of canonical parameters  $\theta_i$  are projected onto the  $k$ -dimensional man- 11  
12 ifold which is the affine subspace spanned by  $k$  PC loading vectors and trans- 12  
13 lated by the intercept vector  $\mu$ . In matrix form, the canonical parameter matrix 13  
14  $\Theta = (\theta_{ij}) = (\theta_1, \dots, \theta_n)^T$  is represented as 14

$$15 \quad (2.2) \quad \Theta = \mathbf{1}_n \otimes \mu^T + \mathbf{A}\mathbf{B}^T, \quad 15$$

17 where  $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_n)^T$  is the  $n \times k$  principal component score matrix and  $\mathbf{B} =$  17  
18  $(\tilde{\mathbf{b}}_1, \dots, \tilde{\mathbf{b}}_k)$  is the  $p \times k$  principal component loading matrix. For identifiability 18  
19 purpose, we require that  $\mathbf{A}$  has orthonormal columns. 19

20 We target a method that can produce a sparse loading matrix, a loading 20  
21 matrix with many zero elements. A sparse loading matrix implies variable selection 21  
22 in principal components analysis, since each principal component only involves 22  
23 those variables corresponding to the nonzero elements of the loading vector. We 23  
24 propose to perform variable selection using the penalized likelihood with a spar- 24  
25 sity inducing penalty. Let  $\mathbf{b}_j^T$  denote the  $j$ th row of  $\mathbf{B}$ . Then (2.2) implies that 25  
26  $\theta_{ij} = \mu_j + \mathbf{a}_i^T \mathbf{b}_j$  where  $\mu_j$  is the  $j$ th element of  $\mu$ . The log likelihood can be 26  
27 written as 27

$$28 \quad (2.3) \quad \ell(\mu, \mathbf{A}, \mathbf{B}) = \sum_{j=1}^d \sum_{i=1}^n \log \pi\{q_{ij}(\mu_j + \mathbf{a}_i^T \mathbf{b}_j)\}. \quad 28$$

30 If  $\mathbf{a}_i$  were observable, (2.3) is the log likelihood for  $d$  logistic regressions 30  
31

$$32 \quad \text{logit } P(Y_{ij} = 1) = \mu_j + \mathbf{a}_i^T \mathbf{b}_j. \quad 32$$

33 This connection with logistic regression suggests use of the  $L_1$  penalty to get a 33  
34 sparse loading matrix, as in the LASSO regression [Tibshirani (1996)]. 34

35 Specifically, consider the penalty 35

$$36 \quad (2.4) \quad P_\lambda(\mathbf{B}) = \sum_{l=1}^k \lambda_l \|\tilde{\mathbf{b}}_l\|_1 = \lambda_1 \sum_{j=1}^d |b_{j1}| + \dots + \lambda_k \sum_{j=1}^d |b_{jk}|, \quad 36$$

37 where  $\lambda_l$  are regularization parameters whose selection will be discussed later. We 37  
38 obtain sparse principal components by maximizing the following penalized log 38  
39  
40  
41  
42  
43

1 likelihood:

$$2 \quad (2.5) \quad f(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B}) = \ell(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B}) - nP_\lambda(\mathbf{B}),$$

3 subject to the constraint that  $\mathbf{A}$  has orthonormal columns. Note that  $\mathbf{B}$  enters the  
4 likelihood together with  $\mathbf{A}$  through  $\mathbf{A}\mathbf{B}^T$  and so  $\mathbf{B}$  can be arbitrarily small by just  
5 increasing the magnitude of  $\mathbf{A}$  and not changing the likelihood. The orthonormal  
6 constraint on  $\mathbf{A}$  prevents elements of  $\mathbf{A}$  becoming arbitrarily large and thus vali-  
7 dates our use of the  $L_1$  penalty on  $\mathbf{B}$ .

8 The sparse principal components can be equivalently formulated as minimizing  
9 the following criterion function:

$$10 \quad (2.6) \quad S(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B}) = -\ell(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B}) + nP_\lambda(\mathbf{B}),$$

11 subject to the constraint that  $\mathbf{A}$  has orthonormal columns. In (2.6) the negative log  
12 likelihood can be interpreted as a loss function and the  $L_1$  penalties increase the  
13 loss for nonzero elements of  $\mathbf{B}$  according to their magnitude. This penalized loss  
14 interpretation is also appealing in the sense that the independent Bernoulli trials as-  
15 sumption for obtaining the likelihood (2.3) need not be a realistic representation of  
16 the actual data generating process but rather a device for generating a suitable loss  
17 function. Since the  $L_1$  penalties regularize the loss minimization, the sparse logis-  
18 tic PCA is sometimes also referred to as the regularized logistic PCA. We shall  
19 focus on the minimization problem (2.6) for the rest of the paper. A computational  
20 algorithm for solving the minimization problem is presented in Section 4.

21 The effectiveness of the proposed sparse logistic PCA is illustrated in Figure 1  
22 using a rank-one model (i.e.,  $k = 1$ ). While the sparse logistic PCA can recover  
23 the original loading vector well, the nonregularized logistic PCA gives more noisy  
24 results. A systematic simulation study is reported in Section 6.

25  
26 *2.2. Choosing the penalty parameters.* Although different penalty parameters  
27 can be used for different PC loading vectors for maximal flexibility of the method-  
28 ology, we consider using only a single penalty parameter  $\lambda$  for all PC loadings.  
29 This simplification substantially reduces the computation time, especially when  $k$   
30 is large. Note that a larger value of  $\lambda$  will lead to a smaller number of nonzeros  
31 in the loading matrix  $\mathbf{B}$  and reduced model complexity, but the reduced model com-  
32 plexity is usually associated with less good fit of the model. To compromise the  
33 goodness of fit and model complexity, for fixed  $k$ , we choose  $\lambda$  by minimizing the  
34 following BIC criterion:

$$35 \quad (2.7) \quad \text{BIC}(\lambda) = -2\ell(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B}) + \log n \times m(\lambda),$$

36 where  $m(\lambda)$  is a measure of the degrees of freedom. Note that Zou, Hastie and  
37 Tibshirani (2007) showed that the number of nonzero coefficients is an unbiased  
38 estimate of the degrees of freedom for the LASSO regression. The degrees of free-  
39 dom  $m(\lambda)$  used in (2.7) is defined as  $m(\lambda) = d + nk + |\mathcal{B}(\lambda)|$ , where  $d$  is the length  
40 of the vector  $\boldsymbol{\mu}$ ,  $nk$  is the total number of elements of  $\mathbf{A}$ , and  $|\mathcal{B}(\lambda)|$  is the cardinal-  
41 ity of the index set  $\mathcal{B}(\lambda)$  of the nonzero loadings in  $\mathbf{B}$  when the penalty parameter  
42 is  $\lambda$ . We use a grid search to find the optimal  $\lambda$  that minimizes the BIC.  
43

6 S. LEE, J. Z. HUANG AND J. HU

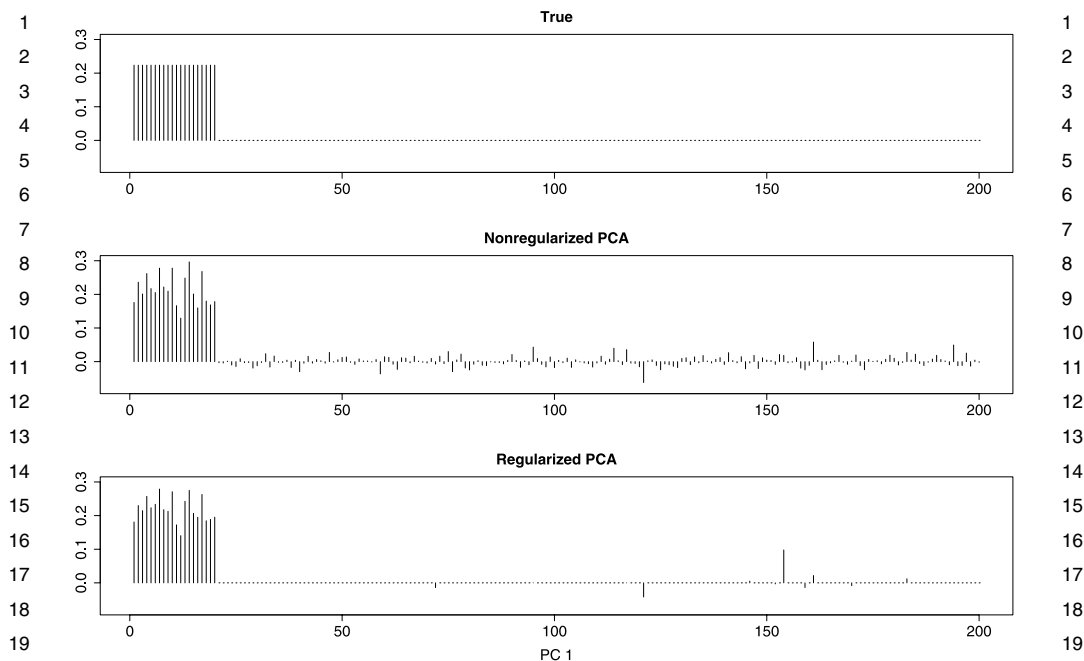


FIG. 1. A simulated data set with  $n = 100$ ,  $d = 200$ , and  $k = 1$ . Top, middle and bottom panels show respectively the true loadings, loadings from the nonregularized logistic PCA and from the regularized logistic PCA. The penalty parameter is selected using the BIC.

2.3. *Determining the dimensionality of the subspace.* The BIC criterion defined in (2.7) can also be used to select a suitable “ $k$ .” A two-dimensional grid search can be used to find the minimizer of the BIC with respect to both  $k$  and  $\lambda$ . To expedite computation, we implement the following strategy: First fix  $k$  at a reasonable large value and select a good  $\lambda$ , then using this  $\lambda$  we refine the choice of  $k$  and, finally, we refine  $\lambda$  with the refined  $k$ . When optimizing with respect to  $\lambda$ , a coarse grid can be used in the first step and a finer grid in the second step. Our simulation study showed that this strategy works reasonably well (see Section 6.3).

REMARK 1. In classical multivariate analysis, the percentage of total variance explained by the principal components provides an intuitive measure that can be used for subjectively choosing the appropriate number of principal components. Zou, Hastie and Tibshirani (2006) and Shen and Huang (2008) extended it to sparse PCA by modifying the definition of variance explained by the PCs. Since there is no clear definition of total variance for the binary data, extension of the notion of “percentage of variance explained” to logistic PCA is an interesting but unsolved problem.

1 **3. Application to single nucleotide polymorphism data.** Association stud- 1  
2 ies based on high-throughput single nucleotide polymorphism (SNP) data 2  
3 [Brookes (1999); Kwok et al. (1996)] have become a popular way to detect ge- 3  
4 nomic regions associated with human complex diseases. A SNP is a single base 4  
5 pair position in genomic DNA at which the sequence (alleles) variation occurs 5  
6 between members of a species, wherein the least frequent allele has an abundance 6  
7 of 1% or greater. A crucial issue in association studies is population stratifica- 7  
8 tion detection [Hao et al. (2004)], which is to determine whether a population is 8  
9 homogeneous or has hidden structures within it. With the presence of population 9  
10 stratification, the naive case-control approach not accounting for this factor would 10  
11 yield biased results [Ewens and Spielman (1995)] and, therefore, draw inaccurate 11  
12 scientific conclusions. See Liang and Kelemen (2008) for an extensive discussion 12  
13 of statistical methods and difficulties for SNP data analysis. 13

14 The proposed sparse logistic PCA method can be used for population strat- 14  
15 ification detection. For the purpose of demonstration, we use the SNP data set 15  
16 available in the International HapMap project [The International HapMap Consor- 16  
17 tium (2005)]. It consists of 3 different ethnic populations of 90 Caucasians (Utah 17  
18 residents with ancestry from northern and western Europe; CEO), 90 Africans 18  
19 (Yoruba in Ibadan, Nigeria; YRI) and 90 Asians (45 Han Chinese in Beijing, 19  
20 China; CHB and 45 Japanese in Tokyo, Japan; JPT). Our task is to detect this 20  
21 three-subpopulation structure using the SNP data on the 270 subjects. At many 21  
22 SNP locations, heterozygosity distribution and allele frequency are known to be 22  
23 different among populations and could confound the effect of the risk of disease. 23  
24 To account for this factor, Serre et al. (2008) selected 1536 SNPs with similar het- 24  
25 erozygosity distribution and allele frequency. The locations of these SNPs cover 25  
26 all the chromosomes except for the sex-determining chromosome. Among these 26  
27 1536 SNPs, 1392 are shared by three ethnic groups, which are used in our analy- 27  
28 sis. We coded 0 for the most prevalent homogeneous base pair (wild-type) and 1 28  
29 for others (mutant), resulting in a  $270 \times 1392$  binary matrix. This data matrix 29  
30 has 2.37% missing entries. 30

31 We applied the sparse logistic PCA to this SNP data set to explore variabil- 31  
32 ity among high dimensional SNP variables, using the computation algorithm 32  
33 given in Sections 4 and 5 below. The method described in Section 2.3 was 33  
34 used for model selection. Specifically, we initially fixed the reduced dimen- 34  
35 sion to  $k = 30$  and chose the penalty parameter  $\lambda$  among the rough grid of 35  
36  $0, 1.5^{-18}, 1.5^{-17}, \dots, 1.5^{-10}$  using the BIC criterion defined in Section 2.3. Given 36  
37 the selected  $\lambda = 1.5^{-16}$ , the dimension  $k$  was refined by minimizing the BIC, 37  
38 giving  $k = 10$ . Finally, with  $k = 10$ , we refined  $\lambda$  by searching over the grid 38  
39  $0, 0.0005, 0.0010, 0.0015, \dots, 0.0100$ , resulting in  $\lambda = 0.0015$ . As a comparison, 39  
40 we also applied the nonregularized logistic PCA to the data, which corresponds to 40  
41  $\lambda = 0$  in our general formulation of regularized logistic PCA. 41

42 To examine which principal components represent the variability associated 42  
43 with three racial groups, we used a  $F$ -test where scores for each fixed PC is 43

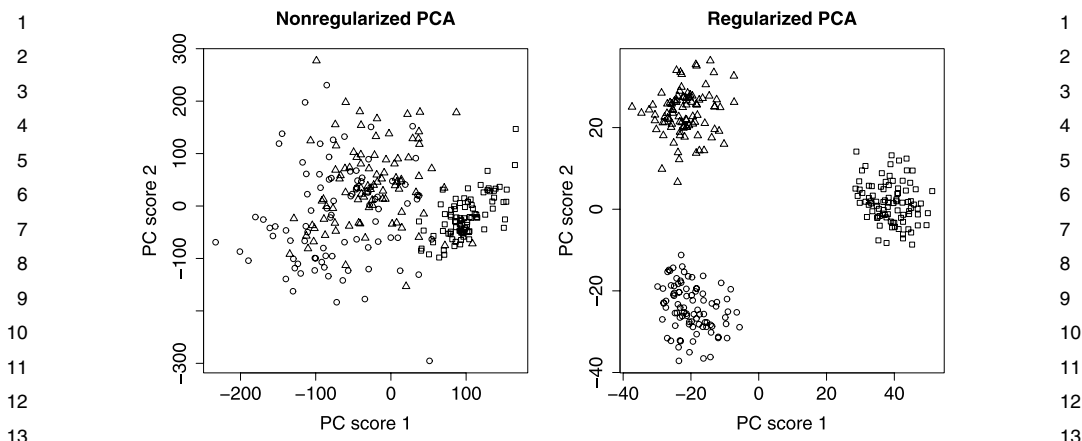


FIG. 2. The scatterplots of the first two PC scores from the nonregularized (left) and regularized logistic PCA. Circle, rectangle and triangle represent Caucasian, African and Asian population respectively.

regressed on the group dummy variables. For the sparse logistic PCA, only the first two PCs were highly significant with both  $p$ -values less than 0.0001 and the remaining eight PCs were not significant with large  $p$ -values (0.7681, 0.9109, 0.4764, 0.5523, 0.3376, 0.5415, 0.4480, 0.6441 for the third to the tenth PCs respectively). This result suggests that the sparse logistic PCA can effectively compress the racial group information into two leading PCs. Similar compression was not achieved by the nonregularized logistic PCA; the  $F$ -test was significant for all the first ten PCs with  $p$ -values  $<0.0001$ ,  $<0.0001$ , 0.0002, 0.0001,  $<0.0001$ ,  $<0.0001$ ,  $<0.0001$ , 0.0028,  $<0.0001$  and 0.0299 respectively.

Pairwise scatterplots were used to check clustering of subjects using the PC scores. Figure 2 shows the scatterplots of first 2 PC scores with and without regularization. The three ethnic groups are clearly separated by the regularized PCA but not by the nonregularized PCA. To verify that the group separation obtained is not because of luck, we permuted observations for each SNP and applied the sparse logistic PCA to the permuted data set; no clear clustering showed up in the PC scores.

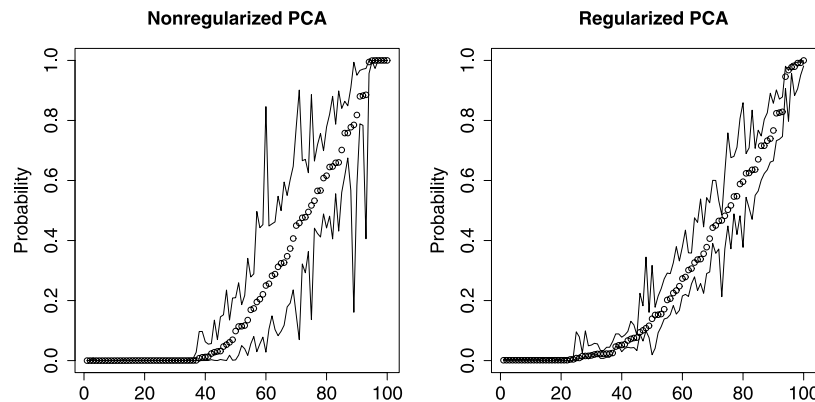
The proposed sparse PCA method allows directly identifying the SNPs that contribute to the group separation. The selected model has 790 and 658 nonzero loadings (representing the SNPs) respectively for the first 2 PCs, among which 509 SNPs are shared. Therefore, 939 SNPs involved in the first 2 PC directions are claimed to be associated with the ethnic group effect. Our result suggests that the population stratification factor should be taken into consideration at these 939 SNP locations in the subsequent study of the association between SNPs and the disease phenotype to avoid biased conclusion. Although in light of our simulation results, some selected SNPs could be false positives, we believe that a large proportion of the selected SNPs are relevant in differentiation among the three racial groups,



1 because the studied SNPs were delicately selected to represent the most genetic 1  
 2 diversity of the whole genome [Serre et al. (2008)] and the genetic differentia- 2  
 3 tion is the greatest when defined on a continental basis, which is the case for our 3  
 4 comparison between Caucasian, Asian and African [Risch et al. (2002)]. 4

5 We further compared the regularized and nonregularized logistic PCA by as- 5  
 6 sessing the variability of the probability estimates using the parametric bootstrap. 6  
 7 For each method, we generated 100 bootstrapped data sets of binary matrices; each 7  
 8 binary matrix has entries that are independently drawn from the Bernoulli distrib- 8  
 9 ution with success probability  $\hat{\pi}_{ij}$  for the  $(i, j)$ th entry, where  $\hat{\pi}_{ij}$  is the estimated 9  
 10 probability. We then applied the method to these bootstrapped data sets to obtain 10  
 11 100 bootstrapped probabilities for each  $(i, j)$  combination and to construct a 90% 11  
 12 variability interval using the 5% and 95% quantiles of the bootstrapped probabil- 12  
 13 ities. These 90% variability intervals were plotted against the ordered  $\hat{\pi}_{ij}$  to form 13  
 14 a variability envelop. The variability envelop for the regularized PCA is narrower 14  
 15 than that for the nonregularized PCA, indicating that regularization indeed reduces 15  
 16 the variability of the probability estimates (Figure 3). 16

17 Our working model for the logistic PCA specified by (2.1) and (2.2) assumes 17  
 18 that, conditional on the principal component scores, the observations are independ- 18  
 19 ent. Since there exists spatial dependency among SNPs, one may have concerns 19  
 20 about the validity of our analysis results if the dependence is strong. In our data 20  
 21 set, the 1536 SNPs were selected from the whole genome to capture most of the 21  
 22 genetic diversity in population considering factors of physical distances, allele fre- 22  
 23 quencies and linkage disequilibrium patterns. The selected SNPs are sufficiently 23  
 24 well separated within each chromosome so that they can be representative of the 24  
 25 whole genome [Serre et al. (?Serre07)]. Therefore, we expect that the spatial depen- 25  
 26 dency in this data set should not be too serious to invalidate our results. To 26  
 27



28  
 29  
 30  
 31  
 32  
 33  
 34  
 35  
 36  
 37  
 38  
 39  
 40  
 41 FIG. 3. The SNP data: 90% bootstrap variability envelope (showed as lines) of the probability 41  
 42 estimates, using 100 randomly selected SNPs. Circles are the estimated probabilities  $\hat{\pi}_{ij}$  from the 42  
 43 SNP data. Results are based on 100 bootstrap samples. 43

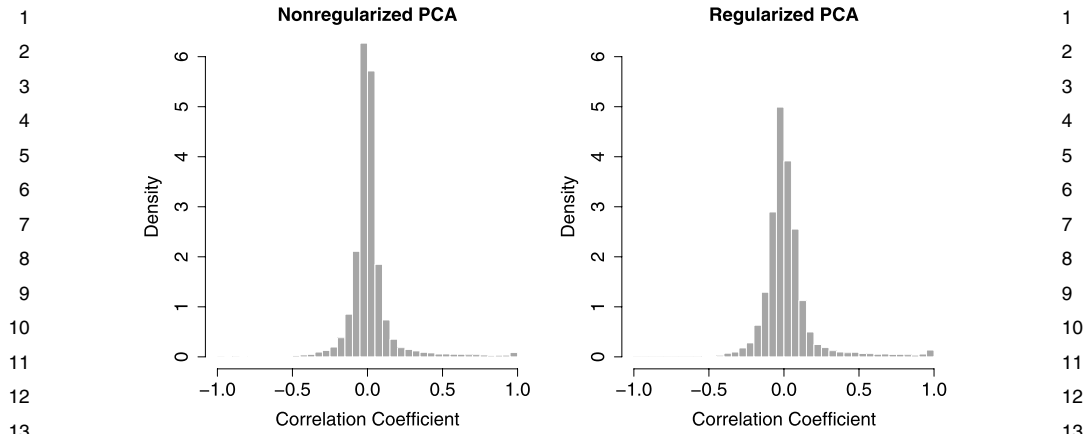


FIG. 4. Histograms of pairwise correlations of Pearson's residuals from nonregularized (left) and regularized (right) logistic PCA.

address this issue empirically, we first computed Pearson's residuals after fitting the models for the nonregularized and regularized logistic PCA, then calculated pairwise correlations of these Pearson's residuals for all SNP pairs for each chromosome. Figure 4 shows the histogram of the pairwise correlations for each model. For both models most pairwise correlations are close to zero, indicating that the SNPs are weakly correlated. We noticed that there exists a very small proportion of SNP pairs that are highly correlated. Examination of the physical locations revealed that those highly correlated SNP pairs consist of SNPs in close vicinity, indicating the imperfection of the initial SNP selection process.

**4. Computational algorithm.** We develop a Majorization–Minimization (MM) algorithm for minimizing (2.6), which iteratively minimizes a suitably defined quadratic upper bound of (2.6). Instead of directly dealing with the non-quadratic log likelihood and the nondifferentiable sparsity inducing  $L_1$  penalty, the MM algorithm sequentially optimizes a quadratic surrogate objective function. A function  $g(x|y)$  is said to majorize a function  $f(x)$  at  $y$  if

$$g(x|y) \geq f(x) \quad \text{for all } x \quad \text{and} \quad g(y|y) = f(y).$$

In the geometrical view the function surface  $g(x|y)$  lies above the function  $f(x)$  and is tangent to it at the point  $y$  so  $g(x|y)$  becomes an upper bound of  $f(x)$ . To minimize  $f(x)$ , the MM algorithm starts from an initial guess  $x^{(0)}$  of  $x$ , and iteratively minimizes  $g(x|x^{(m)})$  until convergence, where  $x^{(m)}$  is the estimate of  $x$  at the  $m$ th iteration. The MM algorithm decreases the objective function in each step and is guaranteed to converge to a local minimum of  $f(x)$ . When applying the MM algorithm, the majorizing function  $g(x|y)$  is chosen such that it is easier to minimize than the original objective function  $f(x)$ . See Hunter and Lange (2004) for an introductory description of the MM algorithm.

1 To find a suitable majorizing function of (2.6), we treat the log likelihood term 1  
 2 and the penalty term separately. For the log likelihood term, note that, for a given 2  
 3 point  $y$ , 3

$$4 \quad (4.1) \quad -\log \pi(x) \leq -\log \pi(y) - \{1 - \pi(y)\}(x - y) + \frac{2\pi(y) - 1}{4y}(x - y)^2 \quad 4$$

$$5 \quad (4.2) \quad \leq -\log \pi(y) - \{1 - \pi(y)\}(x - y) + \frac{1}{8}(x - y)^2, \quad 5$$

6 and the equalities hold when  $x = y$  [Jaakkola and Jordan (2000); de Leeuw 6  
 7 (2006)]. These inequalities provide quadratic upper bounds for the negative log 7  
 8 inverse logit function at the tangent point  $y$ . We refer to the former bound as the 8  
 9 tight bound, and the latter bound as the uniform bound since its curvature does 9  
 10 not change with  $y$ . We pursue here the MM algorithm by using the uniform bound 10  
 11 and leave the discussion of using the tight bound to the supplemental article [Lee, 11  
 12 Huang and Hu (2010)]. Use of the tight bound usually leads to a smaller number of 12  
 13 iterations of the algorithm but longer computation time because of the complexity 13  
 14 involved in computing the bound. For the penalty term, the inequality 14  
 15 16  
 17

$$18 \quad (4.3) \quad |x| \leq \frac{x^2 + y^2}{2|y|}, \quad y \neq 0, \quad 18$$

19 gives an upper bound for  $|x|$  and the equality holds when  $x = y$  [Hunter and Li 19  
 20 (2005)]. Application of (4.2) and (4.3) yields a suitable majorizing function of (2.6) 20  
 21 and thus an MM algorithm. 21

22 Now we present details of the MM algorithm via the uniform bound. Let  $\Theta^{(m)}$  22  
 23 be the estimate of  $\Theta$  obtained in the  $m$ th step of the algorithm, with the entries 23  
 24  $\theta_{ij}^{(m)} = \mu_j^{(m)} + \mathbf{a}_i^{(m)T} \mathbf{b}_j^{(m)}$ . By completing the square, the uniform bound (4.2) can 24  
 25 be rewritten as 25  
 26

$$27 \quad (4.4) \quad -\log \pi(x) \leq -\log \pi(y) + \frac{1}{8}[x - y - 4\{1 - \pi(y)\}]^2. \quad 27$$

28 Substituting  $x$  and  $y$  with  $q_{ij}\theta_{ij}$  and  $q_{ij}\theta_{ij}^{(m)}$  respectively in (4.4) and noticing that 28  
 29  $q_{ij} = \pm 1$ , we obtain 29

$$30 \quad (4.5) \quad -\log \pi(q_{ij}\theta_{ij}) \leq -\log \pi(q_{ij}\theta_{ij}^{(m)}) + w_{ij}^{(m)}(\theta_{ij} - x_{ij}^{(m)})^2, \quad 30$$

31 where  $w_{ij}^{(m)} = 1/8$  and 31

$$32 \quad (4.6) \quad x_{ij}^{(m)} = \theta_{ij}^{(m)} + 4q_{ij}\{1 - \pi(q_{ij}\theta_{ij}^{(m)})\}. \quad 32$$

33 The superscript  $m$  of  $w_{ij}^{(m)}$  and  $x_{ij}^{(m)}$  indicates the dependence on  $\Theta^{(m)}$ . Summing 33  
 34 over all  $i, j$  of (4.5) and ignoring a constant term that does not depend on unknown 34  
 35 parameters, we obtain the following quadratic upper bound of the negative log- 35  
 36 likelihood: 36

$$37 \quad (4.7) \quad \sum_{i=1}^n \sum_{j=1}^d w_{ij}^{(m)} (\theta_{ij} - x_{ij}^{(m)})^2 = \sum_{i=1}^n \sum_{j=1}^d w_{ij}^{(m)} \{x_{ij}^{(m)} - (\mu_j + \mathbf{a}_i^T \mathbf{b}_j)\}^2. \quad 37$$

43

1 On the other hand, (4.3) implies that the penalty  $P_\lambda(\mathbf{B})$  has the following quadratic  
2 upper bound:

$$(4.8) \quad P_\lambda(\mathbf{B}) \leq \lambda_1 \sum_{j=1}^d \frac{b_{j1}^2 + b_{j1}^{(m)2}}{2|b_{j1}^{(m)}|} + \cdots + \lambda_k \sum_{j=1}^d \frac{b_{jk}^2 + b_{jk}^{(m)2}}{2|b_{jk}^{(m)}|}.$$

3  
4  
5  
6  
7 Combining (4.7) and (4.8) yields the following quadratic upper bound (up to a  
8 constant) of the criterion function  $S(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B})$  defined in (2.6):

$$(4.9) \quad \begin{aligned} & g(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B} | \boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)}) \\ &= \sum_{i=1}^n \sum_{j=1}^d [w_{ij}^{(m)} \{x_{ij}^{(m)} - (\mu_j + \mathbf{a}_i^T \mathbf{b}_j)\}^2 + \mathbf{b}_j^T \mathbf{D}_{\lambda, j}^{(m)} \mathbf{b}_j], \end{aligned}$$

9  
10  
11  
12  
13 where  $\mathbf{D}_{\lambda, j}^{(m)}$  is a diagonal matrix with diagonal elements  $\lambda_l / \{2|b_{jl}^{(m)}|\}$  for  $l =$   
14  $1, \dots, k$ .

15  
16  
17 **THEOREM 4.1.** (i) *Up to a constant that depends on  $\boldsymbol{\mu}^{(m)}$ ,  $\mathbf{A}^{(m)}$  and  $\mathbf{B}^{(m)}$*   
18 *but not on  $\boldsymbol{\mu}$ ,  $\mathbf{A}$  and  $\mathbf{B}$ , the function  $g(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B} | \boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)})$  defined in (4.9)*  
19 *majorizes  $S(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B})$  at  $(\boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)})$ .*

20 (ii) *Let  $(\boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)})$ ,  $m = 1, 2, \dots$ , be a sequence obtained by iteratively*  
21 *minimizing the majorizing function. Then  $S(\boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)})$  decreases as  $m$  gets*  
22 *larger and it converges to a local minimum of  $S(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B})$  as  $m$  goes to infinity.*

23  
24 The majorizing function given in (4.9) is quadratic in each of  $\boldsymbol{\mu}$ ,  $\mathbf{A}$  and  $\mathbf{B}$  when  
25 the other two are fixed and, thus, alternating minimization of (4.9) with respect  
26 to  $\boldsymbol{\mu}$ ,  $\mathbf{A}$  and  $\mathbf{B}$  has closed-form solutions, which are given below. We now drop  
27 the superscript in  $x_{ij}^{(m)}$  for notational convenience. Recall that  $w_{ij}^{(m)} = 1/8$  is a  
28 constant. For fixed  $\mathbf{A}$  and  $\mathbf{B}$ , set  $x_{ij}^\dagger = x_{ij} - \mathbf{a}_i^T \mathbf{b}_j$ , the optimal  $\hat{\mu}_j$  is given by

$$(4.10) \quad \hat{\mu}_j = \arg \min_{\mu_j} \sum_{i=1}^n (x_{ij}^\dagger - \mu_j)^2 = \frac{1}{n} \sum_{i=1}^n x_{ij}^\dagger, \quad j = 1, \dots, d.$$

29  
30  
31  
32 This leads to a simple matrix formula  $\hat{\boldsymbol{\mu}} = \frac{1}{n} \mathbf{X}^{\dagger T} \mathbf{1}_n$ , which is obtained by taking  
33 the column means of  $\mathbf{X}^\dagger = (x_{ij}^\dagger)$ .

34  
35 To update  $\mathbf{A}$  and  $\mathbf{B}$  for fixed  $\boldsymbol{\mu}$ , set  $x_{ij}^* = x_{ij} - \mu_j$  or in matrix form,  $\mathbf{X}^* =$   
36  $(x_{ij}^*) = \mathbf{X} - \mathbf{1}_n \otimes \boldsymbol{\mu}^T$ . Denote the  $i$ th row vector of  $\mathbf{X}^*$  as  $\mathbf{x}_i^{*T}$ . For fixed  $\boldsymbol{\mu}$  and  $\mathbf{B}$ ,  
37 the  $i$ th row of  $\mathbf{A}$  is updated by minimizing with respect to  $\mathbf{a}_i$  the sum of squares  
38  $\sum_{j=1}^d (x_{ij}^* - \mathbf{a}_i^T \mathbf{b}_j)^2 = (\mathbf{x}_i^* - \mathbf{B} \mathbf{a}_i)^T (\mathbf{x}_i^* - \mathbf{B} \mathbf{a}_i)$ , which has a closed form solution

$$(4.11) \quad \hat{\mathbf{a}}_i = (\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T \mathbf{x}_i^*, \quad i = 1, \dots, n,$$

39  
40  
41 or  $\hat{\mathbf{A}} = \mathbf{X}^* \mathbf{B} (\mathbf{B}^T \mathbf{B})^{-1}$  in matrix form. The columns of updated  $\mathbf{A}$  can be made or-  
42 thonormal by using the QR decomposition. Denote the  $j$ th column vector of  $\mathbf{X}^*$   
43

1 as  $\tilde{\mathbf{x}}_j^*$ . For fixed  $\boldsymbol{\mu}$  and  $\mathbf{A}$ , the  $j$ th row of  $\mathbf{B}$  is updated by solving the ridge regres- 1  
 2 sion problem that minimizes with respect to  $\mathbf{b}_j$  the penalized sum of squares 2

$$\begin{aligned} & \frac{1}{8} \sum_{i=1}^n (x_{ij}^* - \mathbf{a}_i^T \mathbf{b}_j)^2 + n \sum_{l=1}^k \lambda_l \frac{b_{jl}^2}{2|b_{jl}^{(m)}|} \\ & = \frac{1}{8} (\tilde{\mathbf{x}}_j^* - \mathbf{A} \mathbf{b}_j)^T (\tilde{\mathbf{x}}_j^* - \mathbf{A} \mathbf{b}_j) + n \mathbf{b}_j^T \mathbf{D}_{\lambda, j} \mathbf{b}_j, \end{aligned}$$

3 which has a closed form solution 3  
 4  
 5  
 6  
 7  
 8

$$(4.12) \quad \hat{\mathbf{b}}_j = (\mathbf{A}^T \mathbf{A} + 8n \mathbf{D}_{\lambda, j})^{-1} \mathbf{A}^T \tilde{\mathbf{x}}_j^*, \quad j = 1, \dots, d.$$

9 Since, during the iteration,  $\mathbf{A}$  is made orthonormal,  $\mathbf{A}^T \mathbf{A}$  becomes the identity 9  
 10 matrix of size  $k$ . Therefore, since the matrices to be inverted are diagonal matrices, 10  
 11  $\hat{\mathbf{b}}_j$  can be obtained by component-wise shrinkage 11  
 12  
 13  
 14

$$\hat{b}_{jl} = \frac{|b_{jl}^{(m)}|}{|b_{jl}^{(m)}| + 4n\lambda_l} \tilde{\mathbf{a}}_l^T \tilde{\mathbf{x}}_j^*, \quad l = 1, \dots, k, j = 1, \dots, d,$$

15 where  $\tilde{\mathbf{a}}_l$  is the  $l$ th column of  $\mathbf{A}$ . 15  
 16  
 17  
 18  
 19

20 The MM algorithm will alternate between (4.10), (4.11) and (4.12) until conver- 20  
 21 gence. The details are summarized in Algorithm 1. In this algorithm,  $k$ , the number 21  
 22 of columns of  $\mathbf{A}$  and  $\mathbf{B}$ , should be specified in advance. Different from the sequen- 22  
 23 tial extraction approach of Shen and Huang (2008), the matrices  $\mathbf{A}$  and  $\mathbf{B}$  obtained 23  
 24 after applying Algorithm 1 depend on the value of  $k$ , but the results are reasonably 24  
 25 stable when  $k$  is large enough. See Section 2.3 for discussion on choice of  $k$ . We 25  
 26 use random initial values for  $\boldsymbol{\mu}$ ,  $\mathbf{A}$  and  $\mathbf{B}$ . As with any nonlinear optimization al- 26  
 27 gorithms, our algorithm is not guaranteed to converge to a global minimum. We 27  
 28 can follow the common practice to random start the algorithm several times and 28  
 29 find the best solution. Our experience is that the algorithm with different initial 29  
 30 values usually converges to the same solution (within the precision specified by 30  
 31 the convergence criterion). 31  
 32  
 33

34 ALGORITHM 1 (Sparse logistic PCA algorithm I). 34

- 35 1. Initialize with  $\boldsymbol{\mu}^{(1)} = (\mu_1^{(1)}, \dots, \mu_d^{(1)})^T$ ,  $\mathbf{A}^{(1)} = (\mathbf{a}_1^{(1)}, \dots, \mathbf{a}_n^{(1)})^T$  and  $\mathbf{B}^{(1)} =$  35  
 36  $(\mathbf{b}_1^{(1)}, \dots, \mathbf{b}_d^{(1)})^T$ . Set  $m = 1$ . 36
- 37 2. Compute  $x_{ij}^{(m)}$  using (4.6) and set  $\mathbf{X}^{(m)} = (x_{ij}^{(m)})$ . 37
- 38 3. Set  $\mathbf{X}^{(m)\dagger} = (x_{ij}^{(m)\dagger})$  with  $x_{ij}^{(m)\dagger} = x_{ij}^{(m)} - \mathbf{a}_i^{(m)T} \mathbf{b}_j^{(m)}$ . Update  $\boldsymbol{\mu}$  using  $\boldsymbol{\mu}^{(m+1)} =$  38  
 39  $\frac{1}{n} \mathbf{X}^{(m)\dagger T} \mathbf{1}_n$ . 39
- 40 4. Set  $\mathbf{X}^{(m+1)*} = \mathbf{X}^{(m)} - \mathbf{1}_n \otimes \boldsymbol{\mu}^{(m+1)T}$ . 40
- 41 5. Update  $\mathbf{A}$  by  $\mathbf{A}^{(m+1)} = \mathbf{X}^{(m+1)*} \mathbf{B}^{(m)} (\mathbf{B}^{(m)T} \mathbf{B}^{(m)})^{-1}$ . Compute the QR decom- 41  
 42 position  $\mathbf{A}^{(m+1)} = \mathbf{Q} \mathbf{R}$  and then replace  $\mathbf{A}^{(m+1)}$  by  $\mathbf{Q}$ . 42  
 43 43

14

S. LEE, J. Z. HUANG AND J. HU

1 6. Set  $\mathbf{C}^{(m+1)} = (c_{jl}^{(m+1)}) = \mathbf{X}^{(m+1)*T} \mathbf{A}^{(m+1)}$ . Update  $\mathbf{B}$  by  $\mathbf{B}^{(m+1)} = (b_{jl}^{(m+1)})$  1  
 2 where 2

$$3 \quad b_{jl}^{(m+1)} = \frac{|b_{jl}^{(m)}|}{|b_{jl}^{(m)}| + 4n\lambda_l} c_{jl}^{(m+1)}, \quad l = 1, \dots, k, j = 1, \dots, d. \quad 3$$

4 7. Repeat steps 2 through 6 with  $m$  replaced by  $m + 1$  until convergence. 4  
 5 5  
 6 6  
 7 7  
 8 8

9 **REMARK 2.** The orthogonalization in step 5 of Algorithm 1 does not 9  
 10 change the descent property of the MM algorithm. Let  $A^{(m+1)}$  be the opti- 10  
 11 mizer before orthogonalization. Then  $S(A^{(m+1)}, B^{(m)}) \leq S(A^{(m)}, B^{(m)})$ , where, 11  
 12 for simplicity,  $\mu$  is omitted from the objective function  $S$ . Let  $A^{(m+1)} =$  12  
 13  $\tilde{A}^{(m+1)} R$  be the QR decomposition of  $A^{(m+1)}$  and let  $\tilde{B}^{(m)} = B^{(m)} R^T$ . Then 13  
 14  $\tilde{A}^{(m+1)} \tilde{B}^{(m)T} = A^{(m+1)} B^{(m)T}$  and so  $S(\tilde{A}^{(m+1)}, \tilde{B}^{(m)}) = S(A^{(m+1)}, B^{(m)})$ . Con- 14  
 15 sequently,  $S(\tilde{A}^{(m+1)}, \tilde{B}^{(m)}) \leq S(A^{(m)}, B^{(m)})$ . 15  
 16 16  
 17 17

18 **5. Handling missing data.** Missing data are commonly encountered in real 18  
 19 applications. In this section we extend our sparse logistic PCA method to cases 19  
 20 when missing data are present. 20

21 Let  $\mathcal{N} = \{(i, j) | y_{ij} \text{ is not observed}\}$  denote the index set for missing values. The 21  
 22 sparse logistic PCA minimizes the following criterion function: 22

$$23 \quad (5.1) \quad T(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B}) = -\ell_{obs}(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B}) + nP_\lambda(\mathbf{B}), \quad 23$$

24 where 24  
 25 25

$$26 \quad (5.2) \quad \ell_{obs}(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B}) = \sum_{(i,j) \notin \mathcal{N}} \sum \log \pi \{q_{ij}(\mu_j + \mathbf{a}_i^T \mathbf{b}_j)\} \quad 26$$

27 can be interpreted as the observed data log likelihood for model (2.2). Similar 27  
 28 to the nonmissing data case, direct minimization of (5.1) is not straightforward 28  
 29 because the log likelihood term is not quadratic and the penalty term is nondif- 29  
 30 ferentiable. Direct minimization of (5.1) is also complicated by the fact that the 30  
 31 summation in the definition of the observed data log likelihood is not over a rec- 31  
 32 tangular region. Again, we develop an iterative MM algorithm to solve the opti- 32  
 33 mization problem. The strategy is to fill in the missing data with the fitted values 33  
 34 based on the current parameter estimates, then proceed with the algorithm that 34  
 35 assumes complete data, and iterate until convergence. 35  
 36 36  
 37 37  
 38 38

39 Define the working variables 39

$$40 \quad (5.3) \quad z_{ij}^{(m)} = \begin{cases} x_{ij}^{(m)}, & (i, j) \notin \mathcal{N}, \\ \theta_{ij}^{(m)} = \mu_j^{(m)} + \mathbf{a}_i^{(m)T} \mathbf{b}_j^{(m)}, & (i, j) \in \mathcal{N}, \end{cases} \quad 40$$

43

1 where  $x_{ij}^{(m)}$  is defined in (4.6). Let

$$(5.4) \quad h(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B} | \boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)}) \\ = \sum_{i=1}^n \sum_{j=1}^d [w_{ij}^{(m)} \{z_{ij}^{(m)} - (\mu_j + \mathbf{a}_i^T \mathbf{b}_j)\}^2 + \mathbf{b}_j^T \mathbf{D}_{\lambda, j}^{(m)} \mathbf{b}_j],$$

2 where  $\mathbf{D}_{\lambda, j}^{(m)}$  are diagonal matrices with diagonal elements  $\lambda_l / \{2|b_{jl}^{(m)}|\}$  for  $l =$   
 3  $1, \dots, k$ . The following result extends Theorem 4.1 to the missing data case. The  
 4 proof is given in the [Appendix](#).

5 **THEOREM 5.1.** (i) *Up to a constant that depends on  $\boldsymbol{\mu}^{(m)}$ ,  $\mathbf{A}^{(m)}$ , and  $\mathbf{B}^{(m)}$*   
 6 *but not on  $\boldsymbol{\mu}$ ,  $\mathbf{A}$ , and  $\mathbf{B}$ , the function  $h(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B} | \boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)})$  defined in (5.4)*  
 7 *majorizes  $T(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B})$  at  $(\boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)})$ .*

8 (ii) *Let  $(\boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)})$ ,  $m = 1, 2, \dots$ , be a sequence obtained by iteratively*  
 9 *minimizing the majorizing function. Then  $T(\boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)})$  decreases as  $m$  gets*  
 10 *larger and it converges to a local minimum of  $T(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B})$  as  $m$  goes to infinity.*

11 Note that the majorizing functions given in (5.4) have the same form as those  
 12 given in (4.9) except that  $x_{ij}^{(m)}$  in (4.9) is changed to  $z_{ij}^{(m)}$  in (5.4). Thus, the com-  
 13 putation algorithm developed in Section 4 is readily applicable in the missing data  
 14 case with a simple replacement of  $x_{ij}^{(m)}$  by  $z_{ij}^{(m)}$ . The working variable  $z_{ij}^{(m)}$  in (5.4)  
 15 is easily understood: It is the same as the nonmissing data case if  $y_{ij}$  is observable;  
 16 otherwise, it is an imputed  $\theta_{ij}$  value based on the reduced rank model (2.2) and the  
 17 current guess of  $\boldsymbol{\mu}$ ,  $\mathbf{A}$  and  $\mathbf{B}$ .

18 **6. Simulation study.** In this section we demonstrate our sparse logistic PCA  
 19 method using a simulation study. The method worked well in various settings that  
 20 we tested, but here we only report results in a challenging case that the number of  
 21 variables  $d$  is bigger than the sample size  $n$ .

22 **6.1. The signal-to-noise ratio.** To facilitate setting up simulation studies, we  
 23 introduce a notion of signal-to-noise ratio for logistic PCA. In our logistic PCA  
 24 model, the entries of the  $n \times d$  data matrix are independent Bernoulli random  
 25 variables with success probability  $\pi_{ij} = \{1 + \exp(-\theta_{ij})\}^{-1}$  for the  $(i, j)$ th cell.  
 26 The matrix of canonical parameters  $\boldsymbol{\Theta} = (\theta_{ij})$  has a reduced rank representation  
 27  $\boldsymbol{\Theta} = \mathbf{1} \otimes \boldsymbol{\mu}^T + \mathbf{A}\mathbf{B}^T$ , where  $\mathbf{A}$  is a  $n \times k$  matrix of PC scores and  $\mathbf{B}$  is a sparse  
 28  $d \times k$  PC loading matrix. In our simulation study, elements of the  $l$ th column of  $\mathbf{A}$   
 29 are independent draws from a zero-mean Gaussian distribution with variance  $\sigma_{al}^2$ ,  
 30  $1 \leq l \leq k$ . The variance  $\sigma_{al}^2$  measures the signal level of the  $l$ th PC. We set up the  
 31 PC variances relative to a suitably defined baseline noise level.

32 We define a baseline noise level for fixed  $n$ ,  $d$  and  $k$  as follows. First we cre-  
 33 ate a binary data matrix by generating  $n \times d$  independent binary variables from  
 34  
 35  
 36  
 37  
 38  
 39  
 40  
 41  
 42  
 43

1 Bernoulli distribution with the success probability  $1/2$ . These binary variables are 1  
 2 understood to come from the pure noise since they are generated without having 2  
 3 any structure on the success probabilities. Then, we conduct a  $k$ -component logis- 3  
 4 tic PCA without regularization and compute the average of the sample varian- 4  
 5 ces of the obtained  $k$  PC scores, which is denoted as  $\sigma_b^2$ . We repeat the above process 5  
 6 of generating “pure noise” binary data matrices a large number of times (e.g., 100) 6  
 7 and take the mean of  $\sigma_b^2$  computed from these matrices as the baseline noise level. 7

8 With the notion of baseline noise level, we define the signal-to-noise ratio (SNR) 8  
 9 for a PC as 9

$$10 \quad (6.1) \quad \text{SNR} = \frac{\text{variance of PC scores}}{\text{baseline noise level}}. \quad 10$$

11 In our simulation study we first compute the baseline noise level for a given com- 11  
 12 bination of  $n$ ,  $d$  and  $k$ , then use the above formula to specify the variances of PC 12  
 13 scores based on the fixed values of SNR. 13  
 14  
 15  
 16

17 **6.2. Simulation setup.** We set the intrinsic dimension to be  $k = 2$  and the num- 17  
 18 ber of rows of the data matrix to be  $n = 100$ . We varied the number of variables  $d$  18  
 19 and the signal-to-noise ratio SNR. We considered three choices of  $d$ :  $d = 200$ , 19  
 20  $d = 500$  and  $d = 1000$ . The scores of the  $l$ th PC were randomly drawn from the 20  
 21  $N(0, \sigma_{al}^2)$  distribution with  $\sigma_{al}^2 = \text{SNR}_l \cdot (\text{baseline noise level})$ , where  $\text{SNR}_l$  is the 21  
 22 SNR for the  $l$ th PC. We considered two settings of SNR:  $(3, 2)$  and  $(5, 3)$ . For ex- 22  
 23 ample, when the SNR is  $(3, 2)$ , the variance of the first PC is 3 times the baseline 23  
 24 noise level and the variance of the second PC is 2 times the baseline noise level. 24  
 25 We construct two sparse PC loading vectors as follows: Let  $b_{j1}$  and  $b_{j2}$  denote 25  
 26 correspondingly the components of the first and the second PC loading vectors. 26  
 27 We let  $b_{j1} = 1$  for  $j = 1, \dots, 20$ ,  $b_{j2} = 1$  for  $j = 21, \dots, 40$ , and the rest of  $b_{jl}$  27  
 28 are all taken to be 0. The mean vector  $\mu$  was set to be a vector of zeros. 28  
 29

30 **6.3. Simulation results.** Logistic PCA with and without sparsity inducing reg- 30  
 31 ularization was conducted on 100 simulated data sets for each setting. When ap- 31  
 32 plying the sparse logistic PCA algorithm, three choice of  $k$  were considered:  $k$  is 32  
 33 fixed at the true value ( $k = 2$ ), at a moderately large value ( $k = 30$ ), and selected 33  
 34 using the BIC. The penalty parameter was selected using the method described in 34  
 35 Section 2.2. 35

36 To measure the closeness of the estimated PC loading matrix  $\hat{\mathbf{B}}$  and the true 36  
 37 loading matrix  $\mathbf{B}$ , we use the principal angle between spaces spanned by  $\hat{\mathbf{B}}$  and  $\mathbf{B}$ . 37  
 38 The principal angle measures the maximum angle between any two vectors on 38  
 39 the spaces generated by the columns of  $\hat{\mathbf{B}}$  and  $\mathbf{B}$ . More precisely, it is defined 39  
 40 by  $\cos^{-1}(\rho) \times 180/\pi$ , where  $\rho$  is the minimum eigenvalue of the matrix  $\mathbf{Q}_{\hat{\mathbf{B}}}^T \mathbf{Q}_{\mathbf{B}}$ , 40  
 41 where  $\mathbf{Q}_{\hat{\mathbf{B}}}$  and  $\mathbf{Q}_{\mathbf{B}}$  are orthogonal basis matrices obtained by the QR decomposi- 41  
 42 tion of matrices  $\hat{\mathbf{B}}$  and  $\mathbf{B}$ , respectively [Golub and van Loan (1996)]. 42  
 43



## SPARSE LOGISTIC PCA

17

TABLE 1

*The results of logistic PCA with and without sparsity inducing regularization, based on 100 simulated data sets for each setting. The reported values are the mean (standard error) of the principal angle ( $^{\circ}$ ) between the estimated and the true PC loading matrices*

$d$	SNR	$k = 2$	$k = 30$	Selected $k$
200	SNR = (3, 2)			
	Nonregularized	12.532 (0.115)	35.725 (0.177)	–
	Regularized	5.860 (0.123)	10.125 (0.324)	5.816 (0.125)
	SNR = (5, 3)			
	Nonregularized	11.913 (0.122)	36.350 (0.189)	–
	Regularized	5.803 (0.128)	9.843 (0.321)	5.769 (0.127)
500	SNR = (3, 2)			
	Nonregularized	10.890 (0.095)	31.884 (0.188)	–
	Regularized	4.731 (0.115)	9.413 (0.282)	4.690 (0.101)
	SNR = (5, 3)			
	Nonregularized	10.166 (0.095)	31.941 (0.193)	–
	Regularized	4.729 (0.121)	9.242 (0.252)	4.544 (0.119)
1000	SNR = (3, 2)			
	Nonregularized	12.018 (0.167)	36.040 (0.181)	–
	Regularized	7.015 (0.486)	11.807 (0.433)	4.534 (0.141)
	SNR = (5, 3)			
	Nonregularized	11.370 (0.156)	36.144 (0.180)	–
	Regularized	6.767 (0.474)	10.825 (0.475)	4.196 (0.127)

The mean and standard deviation of principal angles for logistic PCA with and without regularization are presented in Table 1. Since smaller principal angles indicate better estimates of the PC loading matrix, the sparsity inducing regularization has a clear benefit—it can substantially reduce the principal angles. The benefit is even more profound when the number of PCs used in the program ( $k = 30$ ) is larger than the true number that was used to generate the data ( $k = 2$ ). The performance of sparse logistic PCA with selected  $k$  is similar to that when  $k$  is fixed at the true value. Frequencies of the selected  $k$  from 100 simulation data sets in each settings of Table 1 are shown in Table 2. When  $d = 200$ , the BIC finds well the true  $k = 2$  but, as  $d$  gets larger, there is a trend that a slightly larger  $k$  is selected. The performance of using BIC to select  $k$  is considered as quite good, given that the sample size is only 100.

A useful feature of the sparse logistic PCA is its ability to select relevant variables when estimating the PC loading vectors. A zero loading of a variable on a PC means that the corresponding variable is not used when forming that PC, and a nonzero loading indicates a useful variable. Our experience with simulated data shows that nonzero loadings can almost always be identified by the method, but some identified nonzero loadings may correspond to irrelevant variables, cases of

TABLE 2  
Frequencies of the selected  $k$  using the BIC

$d$	SNR	Selected $k$						
		1	2	3	4	5	6	7
200	(3, 2)	0	95	5	0	0	0	0
	(5, 3)	0	96	4	0	0	0	0
500	(3, 2)	1	58	37	4	0	0	0
	(5, 3)	0	60	36	3	1	0	0
1000	(3, 2)	3	34	36	15	10	1	1
	(5, 3)	2	31	47	15	4	1	0

false positives. Table 3 presents the percentages of false positives for various settings reported in Table 1. When  $d$  is 500 or 1000, the percentages of false positives are low, all below 20%. But when  $d$  is 200, the percentages of false positives are between 40% and 50%, suggesting big room for improvement in variable selection.

**7. Discussion and extension.** In this paper we propose a sparse PCA method for analyzing binary data by maximizing a penalized Bernoulli likelihood. The sparsity inducing  $L_1$  penalty is used to acquire simple principal components for the sake of easy interpretation and stable estimation. The MM algorithm developed for implementation of our method provides a unified solution for dealing with (i) the nonquadratic likelihood, (ii) the nondifferentiable penalty function, and (iii) presence of missing data. Although the theoretical derivation is not straightforward, the steps of the algorithm are very simple—they are (weighted) penalized least squares with closed-form expressions.

TABLE 3  
The results of logistic PCA with sparsity inducing regularization, based on 100 simulated data sets for each setting in Table 1. The reported values are the mean (standard error) of the percentages of false positives. The description of results is in the text

$d$	SNR	$k = 2$	$k = 30$	Selected $k$
200	(3, 2)	45.05(1.54)	41.51 (1.39)	44.94 (1.51)
	(5, 3)	48.16 (1.63)	40.53 (1.36)	48.26 (1.63)
500	(3, 2)	14.83 (0.74)	18.91 (0.51)	16.70 (0.72)
	(5, 3)	16.06 (0.68)	18.78 (0.42)	16.93 (0.68)
1000	(3, 2)	10.87 (0.75)	12.80 (0.73)	10.13 (0.60)
	(5, 3)	10.89 (0.70)	12.86 (0.73)	9.26 (0.50)

1 We have focused on the logit link so far, but other link functions can also be 1  
 2 used. In particular, a slight modification of the proposed method can handle the 2  
 3 probit link, where the success probabilities  $\theta_{ij} = \Phi^{-1}(\pi_{ij})$  with  $\Phi(\cdot)$  being the 3  
 4 c.d.f. of the standard Gaussian distribution. The log likelihood function (2.3) of 4  
 5 the reduced rank model is changed to 5

$$6 \quad (7.1) \quad \ell(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B}) = \sum_{j=1}^d \sum_{i=1}^n \log \Phi\{q_{ij}(\mu_j + \mathbf{a}_i^T \mathbf{b}_j)\}. \quad 6$$

7 Instead of using the majorization in (4.2), we apply the following upper bound to 7  
 8 majorize the negative log likelihood: 8

$$9 \quad (7.2) \quad -\log \Phi(x) \leq -\log \Phi(y) - \frac{\phi(y)}{\Phi(y)}(x - y) + \frac{1}{2}(x - y)^2, \quad 9$$

10 where  $\phi(\cdot)$  is the Gaussian density [Böhning (1999); de Leeuw (2006)]. Algo- 10  
 11 rithm 1 still applies with appropriate changes to the definitions of the weights  $w_{ij}^{(m)}$  11  
 12 and the working variables  $x_{ij}^{(m)}$ . 12

13 Our method can also be extended in a straightforward way to handle compos- 13  
 14 ite data which consists of both binary and continuous variables. While the binary 14  
 15 variables are modeled with Bernoulli distributions, the continuous variables can 15  
 16 be modeled with Gaussian distributions. Including some continuous variables cor- 16  
 17 responds to adding some negative Gaussian log likelihood terms to the log like- 17  
 18 lihood expression (2.3). Since the Gaussian log likelihood is quadratic, it blends 18  
 19 in easily with the quadratic majorization used for the logistic PCA. Specifically, 19  
 20 if the  $j$ th variable is of a continuous type, we assume  $y_{ij} \sim N(\theta_{ij}, \sigma^2)$  with  $\theta_{ij}$  20  
 21 satisfying (2.2), and simply let  $x_{ij}^{(m)} = y_{ij}$  and  $w_{ij}^{(m)} = 1/\sigma^2$  when forming the 21  
 22 majorizing function (4.9). The residual variance  $\sigma^2$  of fitting the continuous 22  
 23 variables can be estimated using the residual sum of squares. Taking into account the fact 23  
 24 that different weighting schemes are used for the binary variables and the con- 24  
 25 tinuous variables in the majorizing function, a slight modification of Algorithm 2 25  
 26 presented in the supplemental article [Lee, Huang and Hu (2010)] can be used for 26  
 27 computation. 27  
 28  
 29  
 30  
 31  
 32

## 33 APPENDIX 33

34 **A.1. Proof of Theorem 4.1.** We prove the results for both the tight and the 34  
 35 uniform bound case. Applications of (4.1) and (4.2) yield the following majorizing 35  
 36 functions of the negative log likelihood  $-\ell(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B})$ : 36

$$37 \quad \sum_{i=1}^n \sum_{j=1}^d \left[ -\log \pi(q_{ij}\theta_{ij}^{(m)}) - q_{ij}\{1 - \pi(q_{ij}\theta_{ij}^{(m)})\}(\theta - \theta_{ij}^{(m)}) \right. \\ 38 \quad \left. + \frac{2\pi(q_{ij}\theta_{ij}^{(m)}) - 1}{4q_{ij}\theta_{ij}^{(m)}}(\theta - \theta_{ij}^{(m)})^2 \right] \quad 38$$

43

1 for the tight bound, and

$$2 \sum_{i=1}^n \sum_{j=1}^d \left[ -\log \pi(q_{ij}\theta_{ij}^{(m)}) - q_{ij} \{1 - \pi(q_{ij}\theta_{ij}^{(m)})\} (\theta - \theta_{ij}^{(m)}) + \frac{1}{8} (\theta - \theta_{ij}^{(m)})^2 \right]$$

3 for the uniform bound. Note that

$$4 \{2\pi(q_{ij}\theta_{ij}^{(m)}) - 1\} / \{4q_{ij}\theta_{ij}^{(m)}\} = \{2\pi(\theta_{ij}^{(m)}) - 1\} / \{4\theta_{ij}^{(m)}\}$$

5 for  $q_{ij} = \pm 1$ . By completing the squares and using the definitions of  $x_{ij}^{(m)}$  and  $w_{ij}^{(m)}$ , these majorizing functions can be rewritten as

$$6 \begin{aligned} & -\tilde{\ell}(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B} | \boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)}) \\ & = -\ell(\boldsymbol{\Theta}^{(m)}) - 2 \sum_{i=1}^n \sum_{j=1}^d \{1 - \pi(q_{ij}\theta_{ij}^{(m)})\}^2 + \sum_{i=1}^n \sum_{j=1}^d w_{ij}^{(m)} (\theta_{ij} - x_{ij}^{(m)})^2. \end{aligned}$$

7 On the other hand, application of (4.3) yields the following majorizing function of  $P_\lambda(\mathbf{B})$ :

$$8 \begin{aligned} \tilde{P}_\lambda(\mathbf{B} | \mathbf{B}^{(m)}) & = \lambda_1 \sum_{j=1}^d \frac{b_{j1}^2 + b_{j1}^{(m)2}}{2|b_{j1}^{(m)}|} + \dots + \lambda_k \sum_{j=1}^d \frac{b_{jk}^2 + b_{jk}^{(m)2}}{2|b_{jk}^{(m)}|} \\ & = \sum_{j=1}^d \mathbf{b}_j^{(m)T} \mathbf{D}_{\lambda,j}^{(m)} \mathbf{b}_j^{(m)} + \sum_{j=1}^d \mathbf{b}_j^T \mathbf{D}_{\lambda,j}^{(m)} \mathbf{b}_j. \end{aligned}$$

9 Since the majorization relation between functions is closed under the formation of sums,  $-\tilde{\ell} + n\tilde{P}_\lambda(\mathbf{B} | \mathbf{B}^{(m)})$  majorizes  $S(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B})$  at  $(\boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)})$ . Noticing that  $-\tilde{\ell} + n\tilde{P}_\lambda(\mathbf{B} | \mathbf{B}^{(m)})$  equals  $g(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B} | \boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)})$  up to a constant independent of  $(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B})$ , we complete the proof of part (i). Part (ii) of the theorem follows from the general property of the MM algorithm [Hunter and Lange (2004)].  $\square$

10 **A.2. Proof of Theorem 5.1.** Note that the objective function to be minimized is the summation of two terms—the log likelihood term and the penalty term. Because the majorization property is closed under function summation, we deal with the two terms separately. We can find a majorization function of the penalty term as in Theorem 4.1. To find a majorization function of the log likelihood term, we apply the argument in the standard EM algorithm for handling missing data [Dempster, Laird and Rubin (1977)]. The complete data log likelihood is

$$11 \ell_{com}(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B}) = \sum_{(i,j) \notin \mathcal{N}} \log \pi(q_{ij}\theta_{ij}) + \sum_{(i,j) \in \mathcal{N}} \log \pi(q_{ij}\theta_{ij}).$$

1 Its conditional expectation given the observed data and the current guess of the  
2 parameter values is

$$\begin{aligned}
 & Q(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B} | \boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)}) \\
 (A.1) \quad & = \sum_{(i,j) \notin \mathcal{N}} \log \pi(q_{ij} \theta_{ij}) \\
 & + \sum_{(i,j) \in \mathcal{N}} E[\log \pi(q_{ij} \theta_{ij}) | \mathbf{Y}_o, \boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)}],
 \end{aligned}$$

10 where  $\mathbf{Y}_o$  denotes the observed data. By the standard EM theory,

$$\begin{aligned}
 (A.2) \quad & -\tilde{\ell}_{obs}(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B}) \triangleq -Q(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B} | \boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)}) - \ell_{obs}(\boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)}) \\
 & + Q(\boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)} | \boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)})
 \end{aligned}$$

14 majorizes  $-\ell_{obs}(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B})$  at  $(\boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)})$ , that is,  $-\tilde{\ell}_{obs}(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B}) \geq$   
15  $-\ell_{obs}(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B})$ , and the equality holds when  $(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B}) = (\boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)})$ .

16 Now we find a quadratic majorizing function of  $-\tilde{\ell}_{obs}(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B})$ , which in turn  
17 majorizes  $-\ell_{obs}(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B})$  because of the transitivity of the majorization relation.  
18 We need only to find a quadratic majorization function of  $-Q(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B} | \boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)},$   
19  $\mathbf{B}^{(m)})$ , since it is the only term in the definition (A.2) of  $-\tilde{\ell}_{obs}(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B})$  that de-  
20 pends on the unknown parameters. According to (A.1),  $-Q(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B} | \boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)},$   
21  $\mathbf{B}^{(m)})$  can be decomposed into two terms, one corresponding to observed data,  
22 the other corresponding to the missing data. The former term can be treated  
23 as in the proof of Theorem 4.1. When  $(i, j) \notin \mathcal{N}$ ,  $-\log \pi(q_{ij} \theta_{ij})$  is majorized  
24 by  $w_{ij}^{(m)} (\theta_{ij} - x_{ij}^{(m)})^2$ , up to a constant. To treat the latter term, note that, when  
25  $(i, j) \in \mathcal{N}$ ,

$$\begin{aligned}
 & E[\log \pi(q_{ij} \theta_{ij}) | \mathbf{Y}_o, \boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)}] \\
 & = \pi(\theta_{ij}^{(m)}) \log \pi(\theta_{ij}) + \{1 - \pi(\theta_{ij}^{(m)})\} \log \{1 - \pi(\theta_{ij})\} \\
 & = \sum_{q_{ij} = \pm 1} \pi(q_{ij} \theta_{ij}^{(m)}) \log \pi(q_{ij} \theta_{ij}),
 \end{aligned}$$

32 using the fact that the missing data are independent of the observed data, and that  
33  $1 - \pi(\theta) = \pi(-\theta)$ . Then, by applying the inequalities (4.1) and (4.2) and using  
34 the definition of  $w_{ij}^{(m)}$ , we obtain that

$$\begin{aligned}
 & -E[\log \pi(q_{ij} \theta_{ij}) | \mathbf{Y}_o, \boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)}] \\
 & \leq \sum_{q_{ij} = \pm 1} \pi(q_{ij} \theta_{ij}^{(m)}) [-\log \pi(\theta_{ij}^{(m)}) - \{1 - \pi(q_{ij} \theta_{ij}^{(m)})\} \{q_{ij} (\theta_{ij} - \theta_{ij}^{(m)})\}] \\
 & \quad + w_{ij}^{(m)} \{(\theta_{ij} - \theta_{ij}^{(m)})\}^2 \\
 & \leq C_m + w_{ij}^{(m)} \{(\theta_{ij} - \theta_{ij}^{(m)})\}^2,
 \end{aligned}$$

1 where  $C_m$  is a constant independent of  $\boldsymbol{\mu}$ ,  $\mathbf{A}$  and  $\mathbf{B}$ . Combining the above re- 1  
 2 sults, we see that  $-Q(\boldsymbol{\mu}, \mathbf{A}, \mathbf{B} | \boldsymbol{\mu}^{(m)}, \mathbf{A}^{(m)}, \mathbf{B}^{(m)})$  is up to a constant majorized 2  
 3 by  $\sum_{ij} w_{ij}^{(m)} \{(\theta_{ij} - z_{ij}^{(m)})\}^2$ , where  $z_{ij}^{(m)}$  equals  $x_{ij}^{(m)}$  if  $(i, j) \notin \mathcal{N}$ , and  $\theta_{ij}^{(m)}$  if 3  
 4  $(i, j) \in \mathcal{N}$ . The proof of part (i) is thus complete. Part (ii) of the theorem follows 4  
 5 from the general result of the MM algorithm.  $\square$  5  
 6

7 **Acknowledgments.** We would like to thank Editor Michael Stein, an Asso- 7  
 8 ciate Editor and two referees for helpful comments. We would also like to thank 8  
 9 Lan Zhou for help in improving the writing of the paper. 9  
 10

## 11 SUPPLEMENTARY MATERIAL 11

12 **The MM algorithm for sparse logistic PCA using the tight bound** (DOI: 12  
 13 [10.1214/10-AOAS327SUPP](https://doi.org/10.1214/10-AOAS327SUPP); .pdf). We develop the MM algorithm for sparse lo- 13  
 14 gistic PCA using the tight majorizing bound. Comparison of the developed algo- 14  
 15 rithm with the MM algorithm using the uniform bound in terms of computing time 15  
 16 is also presented. 16  
 17

## 18 REFERENCES 18

- 19 BÖHNING, D. (1999). The lower bound method in probit regression. *Comput. Statist. Data Anal.* **30** 19  
 20 13–17. [MR1681451](#) 20  
 21 BROOKES, A. J. (1999). Review: The essence of SNPs. *Gene* **234** 177–186. 21  
 22 COLLINS, M., DASGUPTA, S. and SCHAPIRE, R. E. (2001). A generalization of principal compo- 22  
 23 nent analysis to the exponential family. In *Advanced in Neural Information Processing System* 23  
 24 **14**. 24  
 25 DE LEEUW, J. (2006). Principal component analysis of binary data by iterated singular value decom- 25  
 26 position. *Comput. Statist. Data Anal.* **50** 21–39. [MR2196220](#) 26  
 27 DEMPSTER, A. P., LAIRD, N. M. and RUBIN, D. B. (1977). Maximum likelihood from incomplete 27  
 28 data via the EM algorithm (with discussion). *J. Roy. Statist. Soc. Ser. B* **39** 1–38. [MR0501537](#) 28  
 29 EWENS, W. J. and SPIELMAN, R. S. (1995). The transmission/disequilibrium test: History, subdivi- 29  
 30 sion, and admixture. *The American Journal of Human Genetics* **57** 455–464. 30  
 31 GOLUB, G. and VAN LOAN, C. (1996). *Matrix Computations*, 3rd ed. The Johns Hopkins Univ. 31  
 32 Press. [MR1417720](#) 32  
 33 HAO, K., LI, C., ROSENOW, C. and WONG, W. H. (2004). Detect and adjust for population strat- 33  
 34 ification in population-based association study using genomic control markers: An application 34  
 35 of Affymetrix Genechip<sup>®</sup> Human Mapping 10K array. *European Journal of Human Genetics* **12** 35  
 36 1001–1006. 36  
 37 HOTELLING, H. (1933). Analysis of a complex of statistical variables into principal components. 37  
 38 *Journal of Educational Psychology* **24** 417–441. 37  
 39 HUNTER, D. R. and LANGE, K. (2004). A tutorial on MM algorithms. *Amer. Statist.* **58** 30–37. 38  
 40 [MR2055509](#) 38  
 41 HUNTER, D. R. and LI, R. (2005). Variable selection using MM algorithms. *Ann. Statist.* **33** 1617– 39  
 42 1642. [MR2166557](#) 40  
 43 JAAKKOLA, T. S. and JORDAN, M. I. (2000). Bayesian parameter estimation via variational meth- 41  
 42 ods. *Statist. Comput.* **10** 25–37. 42  
 43 JOLLIFFE, I. T. (2002). *Principal Component Analysis*. Springer. [MR2036084](#) 43

- 1 JOLLIFFE, I. T., TREDAFILOV, M. and UDDINE, M. (2003). A modified principal component 1  
 2 technique based on the LASSO. *J. Comput. Graph. Statist.* **12** 531–547. MR2002634 2  
 3 KWOK, P. Y., DENG, Q., ZAKERI, H., TAYLOR, S. L. and NICKERSON, D. A. (1996). Increasing 3  
 4 the information content of STS-based genome maps: Identifying polymorphisms in mapped 4  
 5 STSs. *Genomics* **31** 123–126. 5  
 6 LANGE, K., HUNTER, D. R. and YANG, I. (2000). Optimization transfer using surrogate objective 6  
 7 functions (with discussion). *J. Comput. Graph. Statist.* **9** 1–20. MR1819865 6  
 8 LEE, S., HUANG, J. Z. and HU, J. (2010). The MM algorithm for sparse logistic PCA using the 7  
 9 tight bound: A supplementary note to “Sparse logistic principal components analysis for binary 8  
 10 data.” DOI: [10.1214/10-AOAS327SUPP](https://doi.org/10.1214/10-AOAS327SUPP). 9  
 11 LIANG, Y. and KELEMEN, A. (2008). Statistical advances and challenges for analyzing correlated 10  
 12 high dimensional SNP data in genomic study for complex diseases. *Statistics Surveys* **2** 43–60. 11  
 13 MR2520980 12  
 14 PEARSON, K. (1901). On lines and planes of closest fit to systems of points in space. *The London,* 13  
 15 *Edinburgh and Dublin Philosophical Magazine and Journal of Science, Sixth Series* **2** 559–572. 14  
 16 RISCH, N., BURCHARD, E., ZIV, E. and TANG, H. (2002). Categorization of humans in biomedical 15  
 17 research: Genes, race and disease. *Genome Biology* **3** comment 2007.1–2007.12. 16  
 18 SCHEIN, A. I., SAUL, L. K. and UNGAR, L. H. (2003). A generalized linear model for princi- 17  
 19 pal component analysis of binary data. In *Proceedings of the Ninth International Workshop on* 18  
 20 *Artificial Intelligence and Statistics* 14–21. 17  
 21 SERRE, D., MONTPETIT, A., PARÉ, G., ENGERT, J. G., YUSUF, S., KEAVNEY, B., HUDSON, 18  
 22 K. J. and ANAND, S. (2008). Correction of population stratification in large multi-ethnic associ- 19  
 23 ation studies. *PLoS ONE* **2** e1382. 20  
 24 SHEN, H. and HUANG, J. Z. (2008). Sparse principal component analysis via regularized low rank 21  
 25 matrix approximation. *J. Multivariate Anal.* **99** 1015–1034. MR2419336 22  
 26 THE INTERNATIONAL HAPMAP CONSORTIUM (2005). A haplotype map of the human genome. 23  
 27 *Nature* **437** 1299–1320. 24  
 28 TIBSHIRANI, R. J. (1996). Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc.* 25  
 29 *Ser. B* **58** 267–288. MR1379242 26  
 30 ZOU, H., HASTIE, T. J. and TIBSHIRANI, R. J. (2006). Sparse principal component analysis. 27  
 31 *J. Comput. Graph. Statist.* **15** 265–286. MR2252527 28  
 32 ZOU, H., HASTIE, T. J. and TIBSHIRANI, R. J. (2007). On the “Degrees of Freedom” of the 29  
 33 LASSO. *Ann. Statist.* **35** 2173–2192. MR2363967 30
- 31 S. LEE J. Z. HUANG 31  
 32 DEPARTMENT OF BIostatISTICS DEPARTMENT OF STATISTICS 32  
 33 HARVARD SCHOOL OF PUBLIC HEALTH TEXAS A&M UNIVERSITY 33  
 34 BOSTON, MASSACHUSETTS 02115 COLLEGE STATION, TEXAS 77843-3143 34  
 35 USA USA 35  
 36 E-MAIL: [seokhol@hsph.harvard.edu](mailto:seokhol@hsph.harvard.edu) E-MAIL: [jianhua@stat.tamu.edu](mailto:jianhua@stat.tamu.edu) 36
- 37 J. HU 37  
 38 DEPARTMENT OF BIostatISTICS 38  
 39 DIVISION OF QUANTITATIVE SCIENCES 39  
 40 UNIVERSITY OF TEXAS M. D. ANDERSON CANCER CENTER 40  
 41 HOUSTON, TEXAS 77030-4009 41  
 42 USA 42  
 43 E-MAIL: [jhu@mdanderson.org](mailto:jhu@mdanderson.org) 43

## 1 THE LIST OF SOURCE ENTRIES RETRIEVED FROM MATHSCINET 1

2 The list of entries below corresponds to the Reference section of your article and was retrieved  
3 from MathSciNet applying an automated procedure. Please check the list and cross out those  
4 entries which lead to mistaken sources. Please update your references entries with the data  
5 from the corresponding sources, when applicable. More information can be found in the sup-  
6 port page:

7 <http://www.e-publications.org/ims/support/mrhelp.html>. 7

8 BÖHNING, D. (1999). The lower bound method in probit regression. *Comput. Statist. Data Anal.* **30**  
9 13–17. MR1681451 (2000a:62149) 9

10 Not Found! 10

11 Not Found! 11

12 DE LEEUW, J. (2006). Principal component analysis of binary data by iterated singular value de-  
13 composition. *Comput. Statist. Data Anal.* **50** 21–39. MR2196220 13

14 DEMPSTER, A. P., LAIRD, N. M., AND RUBIN, D. B. (1977). Maximum likelihood from incom-  
15 plete data via the EM algorithm. *J. Roy. Statist. Soc. Ser. B* **39** 1–38. With discussion. MR0501537  
16 (58 #18858) 16

17 Not Found! 17

18 GOLUB, G. H. AND VAN LOAN, C. F. (1996). *Matrix computations*, Third ed. Johns Hopkins Stud-  
19 ies in the Mathematical Sciences. Johns Hopkins University Press, Baltimore, MD. MR1417720  
20 (97g:65006) 19

21 Not Found! 21

22 HUNTER, D. R. AND LANGE, K. (2004). A tutorial on MM algorithms. *Amer. Statist.* **58** 30–37.  
23 MR2055509 22

24 HUNTER, D. R. AND LI, R. (2005). Variable selection using MM algorithms. *Ann. Statist.* **33**  
25 1617–1642. MR2166557 24

26 Not Found! 25

27 JOLLIFFE, I. T. (2002). *Principal component analysis*, Second ed. Springer Series in Statistics.  
28 Springer-Verlag, New York. MR2036084 (2004k:62010) 26

29 JOLLIFFE, I. T., TREDAFILOV, N. T., AND UDDIN, M. (2003). A modified principal component  
30 technique based on the LASSO. *J. Comput. Graph. Statist.* **12** 531–547. MR2002634 27

31 Not Found! 28

32 LANGE, K., HUNTER, D. R., AND YANG, I. (2000). Optimization transfer using surrogate objective  
33 functions. *J. Comput. Graph. Statist.* **9** 1–59. With discussion, and a rejoinder by Hunter and  
34 Lange. MR1819865 (2001k:62029) 31

35 Not Found! 32

36 LIANG, Y. AND KELEMEN, A. (2008). Statistical advances and challenges for analyzing corre-  
37 lated high dimensional SNP data in genomic study for complex diseases. *Stat. Surv.* **2** 43–60.  
38 MR2520980 33

39 Not Found! 34

40 Not Found! 35

41 Not Found! 36

42 Not Found! 37

43 SHEN, H. AND HUANG, J. Z. (2008). Sparse principal component analysis via regularized low rank  
matrix approximation. *J. Multivariate Anal.* **99** 1015–1034. MR2419336 (2009m:62184) 39

44 Not Found! 40

45 TIBSHIRANI, R. (1996). Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser.*  
46 *B* **58** 267–288. MR1379242 (96j:62134) 42

47 Not Found! 43



1	ZOU, H., HASTIE, T., AND TIBSHIRANI, R. (2006). Sparse principal component analysis. <i>J.</i>	1
2	<i>Comput. Graph. Statist.</i> <b>15</b> 265–286. MR2252527	2
3	ZOU, H., HASTIE, T., AND TIBSHIRANI, R. (2007). On the “degrees of freedom” of the lasso. <i>Ann.</i>	3
4	<i>Statist.</i> <b>35</b> 2173–2192. MR2363967 (2009d:62096)	4
5		5
6		6
7		7
8		8
9		9
10		10
11		11
12		12
13		13
14		14
15		15
16		16
17		17
18		18
19		19
20		20
21		21
22		22
23		23
24		24
25		25
26		26
27		27
28		28
29		29
30		30
31		31
32		32
33		33
34		34
35		35
36		36
37		37
38		38
39		39
40		40
41		41
42		42
43		43

1 META DATA IN THE PDF FILE 1

2 Following information will be included as pdf file Document Properties: 2

3  
 4 **Title** : Sparse logistic principal components analysis for binary 4  
 data 4  
 5 **Author** : Seokho Lee, Jianhua Z. Huang, Jianhua Hu 5  
 6 **Subject** : The Annals of Applied Statistics , 2010, Vol.0, No.00, 1- 6  
 7 26 7  
 8 **Keywords**: Binary data, dimension reduction, MM algorithm, LASSO, PCA, 8  
 9 regularization, sparsity. 9

10  
 11 THE LIST OF URI ADDRESSES 11

12  
 13  
 14 Listed below are all uri addresses found in your paper. The non-active uri addresses, if any, are 14  
 15 indicated as ERROR. Please check and update the list where necessary. The e-mail addresses 15  
 16 are not checked – they are listed just for your information. More information can be found in 16  
 17 the support page: 17

18 <http://www.e-publications.org/ims/support/urihelp.html>. 18

19 200 <http://www.imstat.org/aoas/> [2:pp.1,1] OK 19  
 20 200 <http://dx.doi.org/10.1214/10-AOAS327> [2:pp.1,1] OK 20  
 21 200 <http://www.imstat.org> [2:pp.1,1] OK 21  
 22 200 <http://lib.stat.cmu.edu/aoas/327/supplement.pdf> [2:pp.22,22] OK 22  
 23 200 <http://dx.doi.org/10.1214/10-AOAS327SUPP> [2:pp.23,23] OK 23  
 24 --- <mailto:seokhol@hsph.harvard.edu> [2:pp.23,23] Check skip 24  
 25 --- <mailto:jianhua@stat.tamu.edu> [2:pp.23,23] Check skip 25  
 26 --- <mailto:jhu@mdanderson.org> [2:pp.23,23] Check skip 26

27 27  
 28 28  
 29 29  
 30 30  
 31 31  
 32 32  
 33 33  
 34 34  
 35 35  
 36 36  
 37 37  
 38 38  
 39 39  
 40 40  
 41 41  
 42 42  
 43 43