


Genetics and population analysis

TeraPCA: a fast and scalable software package to study genetic variation in tera-scale genotypes

Aritra Bose ^{1,†}, Vassilis Kalantzis^{2,†}, Eugenia-Maria Kontopoulou^{1,†}, Mai Elkady¹, Peristera Paschou^{3,*} and Petros Drineas¹

¹Computer Science Department, Purdue University, West Lafayette, IN 47907, USA, ²IBM Research, Thomas J. Watson Research Center, Yorktown Heights, NY 10598, USA and ³Department of Biological Sciences, Purdue University, West Lafayette, IN 47907, USA

*To whom correspondence should be addressed.

[†]The authors wish it to be known that, in their opinion, the first three authors should be regarded as Joint First Authors.

Associate Editor: Russell Schwartz

Received on August 8, 2018; revised on February 26, 2019; editorial decision on February 28, 2019; accepted on April 4, 2019

Abstract

Motivation: Principal Component Analysis is a key tool in the study of population structure in human genetics. As modern datasets become increasingly larger in size, traditional approaches based on loading the entire dataset in the system memory (Random Access Memory) become impractical and out-of-core implementations are the only viable alternative.

Results: We present TeraPCA, a C++ implementation of the Randomized Subspace Iteration method to perform Principal Component Analysis of large-scale datasets. TeraPCA can be applied both in-core and out-of-core and is able to successfully operate even on commodity hardware with a system memory of just a few gigabytes. Moreover, TeraPCA has minimal dependencies on external libraries and only requires a working installation of the BLAS and LAPACK libraries. When applied to a dataset containing a million individuals genotyped on a million markers, TeraPCA requires <5 h (in multi-threaded mode) to accurately compute the 10 leading principal components. An extensive experimental analysis shows that TeraPCA is both fast and accurate and is competitive with current state-of-the-art software for the same task.

Availability and implementation: Source code and documentation are both available at <https://github.com/aritra90/TeraPCA>.

Contact: ppaschou@purdue.edu

Supplementary information: [Supplementary data](#) are available at *Bioinformatics* online.

1 Introduction

Principal Component Analysis (PCA) is perhaps the most fundamental unsupervised linear dimensionality reduction technique. It was invented by Pearson in the early 1900s (Pearson, 1901); and later reinvented and named by Hotelling in the 1930s (Hotelling, 1933, 1936). In statistical parlance, PCA converts a set of observations of possibly correlated variables into a set of linearly uncorrelated (orthogonal) variables called principal components (PCs). The seminal work of Luca Cavalli-Sforza and collaborators in the late 1970s

(Chisholm *et al.*, 1995; Menozzi *et al.*, 1978) pioneered the application of PCA for the study of human genetic variation.

PCA analyses and plots appear in virtually every single paper that analyzes human genetic variation in order to make inferences about population structures. Given m samples genotyped on n genetic loci, it is well-known that applying PCA on the $m \times m$ covariance matrix that emerges by computing any reasonable notion of genotypic distance between every pair of samples using the n genotyped loci results in the observation that the leading PCs mirror

geography (Novembre *et al.*, 2008; Paschou *et al.*, 2014; Wang *et al.*, 2010) for detailed discussions and examples. This observation was leveraged by Price *et al.* (2006, 2010) and Patterson *et al.* (2006) to derive one of the most established methods to account (and correct) for the confounding effects of population stratification in genome-wide association studies (GWAS). The method in Price *et al.* (2006, 2010) and Patterson *et al.* (2006) is essentially equivalent to using a small number of leading PCs as covariates in order to check for associations between genetic loci and affection status in statistical tests, and is implemented in the EIGENSTRAT software package which is routinely used in GWAS analyses to correct for population stratification. Other applications of PCA include the identification of sets of genetic loci that are ancestry-informative or are under selective pressure (Paschou *et al.*, 2007; 2008; Price *et al.*, 2006); and, when combined with other lines of evidence such as social structure and linguistics, the extraction of complex population histories and demographic structures (Bose *et al.*, 2017). We also note that PCA extracts the fundamental features of a dataset without complex computational modeling. Interestingly, even the output of model-based, more complex, methods to detect population structure [such as ADMIXTURE (Alexander *et al.*, 2009)] typically exhibits high correlation with the output of PCA, rendering further support to the significance of PCA in the analysis of human genetics data.

From a computational viewpoint, PCA essentially amounts to computing eigenvectors of the $m \times m$ (normalized) covariance matrix associated with the dataset at hand. When m does not exceed a few thousands, all eigenvectors can be computed by appropriate dense linear algebra routines in LAPACK, a Fortran 90 matrix factorization-based library which is widely used for solving systems of linear equations, least-squares problems, eigenvalue problems and singular value problems (Anderson *et al.*, 1999). Matrix factorization-based dense eigenvalue solvers return all m eigenvectors with a time complexity in the order of $O(m^3)$, which becomes impractical as m , the number of samples, increases. Practical applications of PCA in population genetics only require the computation of those PCs determined by the eigenvectors associated with only a few (say 10–20) of the largest eigenvalues. Computing a few of the leading eigenvalues and associated eigenvectors of large (sparse or dense) matrices is typically achieved by first projecting the original eigenvalue problem onto a low-dimensional subspace which includes an invariant subspace associated with the relevant eigenvectors. This low-dimensional subspace can be formed in many different ways, e.g. by means of subspace iteration or Krylov projection schemes and much work in the Numerical Analysis community has been devoted in understanding the theoretical properties of such approaches (Parlett, 1998; Saad, 2011). In particular, a variant of the family of Krylov projection schemes, the so-called Implicitly Restarted Arnoldi method (IRA), is the projection scheme of choice in FlashPCA2 (Abraham *et al.*, 2017), a software package which has been shown to outperform other PCA software packages, both in terms of memory usage and wall-clock time. On the other hand, recent advances in the design and analysis of Randomized Numerical Linear Algebra (Drineas and Mahoney, 2016) algorithms have yielded novel insights as well as fast and efficient alternatives to approximate the leading PCs of large matrices (Drineas *et al.*, 2018; Drineas and Mahoney, 2018; Halko *et al.*, 2011; Musco and Musco, 2015). Indeed, FastPCA (Galinsky *et al.*, 2016) applied such randomized algorithms to perform PCA analyses in population genetics data.

This paper presents TeraPCA, a C++ software package to perform PCA of tera-scale genotypic datasets that cannot fully reside in

the system memory. TeraPCA is essentially an out-of-core implementation of the Randomized Subspace Iteration method (Halko *et al.*, 2011; Rokhlin *et al.*, 2010) and features minimal dependencies to external (in contrast to FlashPCA2 which relies on the IRA implementation on the Spectra C++ library, TeraPCA comes with an in-house implementation of the Randomized Subspace Iteration algorithm.) libraries. As the amount of time spent on input/output typically dominates the wall-clock time in out-of-core scenarios, TeraPCA builds a high-dimensional initial approximation subspace by loading the dataset from secondary storage exactly once. The dimension of this initial approximation subspace can be controlled directly by the user. Each subsequent iteration of Randomized Subspace Iteration ‘corrects’ the initial subspace so that an invariant subspace associated with the leading target eigenvectors is computed. The dataset needs to be accessed twice in each iteration, but, fortunately, a few steps of Randomized Subspace Iteration are typically sufficient in practice in order to get highly accurate approximations to the leading eigenvectors. Note here that the above idea is somewhat orthogonal to the ideas underlying IRA, which builds the approximation subspace in a vector-by-vector manner, thus necessitating a large number of dataset fetches from secondary storage to even form an approximation subspace whose dimension is equal to or slightly larger than the number of PCs that we seek to approximate.

TeraPCA was tested extensively on both real [Human Genome Diversity Panel (HGDP), 1000 Genomes, etc.] and synthetic datasets. Our synthetic datasets were generated via the Pritchard–Stephens–Donnelly (PSD) model (Gopalan *et al.*, 2016; Pritchard *et al.*, 2000). Our results suggest that TeraPCA is both fast and accurate and in most cases outperforms other out-of-core PCA libraries such as FlashPCA2. Specific highlights include the computation of the 10 leading PCs of a dataset of 1 million samples genotyped on 1 million genetic markers (this dataset exceeds 3.5 TBs in uncompressed format) in about 13 h (using a single thread) and in <4.5 h (using 12 threads).

2 Materials and methods

2.1 Simulated datasets

The first group of the datasets used for our experiments was generated using the PSD model of simulating genotypes. In particular, a recent study (Gopalan *et al.*, 2016) simulated genotypic data by obtaining individual ancestry proportions from the PSD model to fit the 1000 Genomes dataset and then modeling the per-population allele frequencies using Wright’s F_{ST} and the Weir and Cockerham estimate (Weir and Cockerham, 1984). We developed a multi-threaded C++ package which is essentially an efficient implementation of the R code developed in Tera-Structure (Gopalan *et al.*, 2016). We generated various datasets in order to evaluate TeraPCA’s performance, with the number of markers ranging from 100 000 to 1 000 000 and the number of samples ranging from 5000 to 1 000 000.

2.2 Real datasets

The HGDP dataset consists of 1043 individuals genotyped at 660 734 single nucleotide polymorphisms (SNPs), across 51 populations across Africa, Europe, Middle East, South and Central Asia, East Asia, Oceania and the Americas (Cann *et al.*, 2002). We ran Quality Control (QC) on the data by filtering SNPs with minor allele frequency below 0.01 and subsequently pruning for LD using a window size of 1000 kb. Moreover, we set the variance inflation

factor to 50 and set $r^2 > 0.2$, thus retaining 1 544 71 variants. We applied the same parameters for LD pruning on the 1000 Genomes dataset which has 2504 individuals sampled from 26 different populations across all continents genotyped at 39 million SNPs. After QC, we retained ~808 704 SNPs and ran our experiments on the pruned dataset.

We also tested the performance of TeraPCA on case-control data, which are ubiquitous in population genetics. We used the Wellcome Trust Case Control Consortium's Type 2 Diabetes (T2D) and Parkinson's (PRK) datasets. The T2D dataset had 6371 individuals (1816 cases and 4555 controls) genotyped on 313 654 SNPs and the PRK dataset had 5000 individuals (2000 cases and 3000 controls) genotyped on 500 000 SNPs. We removed related samples from these datasets and pruned them using the aforementioned QC parameters resulting in datasets with 6370 individuals genotyped on 72 457 SNPs for T2D and 4706 individuals genotyped on 111 831 SNPs for PRK.

2.3 TeraPCA

TeraPCA first normalizes the genotypes using the same procedure that was used by both FlashPCA (Abraham and Inouye, 2014) and FastPCA (Galinsky *et al.*, 2016) (see our [Supplementary Material](#) for details) and then applies Randomized Subspace Iteration in an out-of-core fashion.

The main parameters of TeraPCA are as follows (see our [Supplementary Material](#) for more details and our code release for full documentation):

1. Number of PCs to be computed (denoted by k). Default value is set to $k := 10$.
2. Number of contiguous rows of the SNP-major input matrix fetched from the secondary storage at each time unit (denoted by β). This can be user-defined or automatically determined based on the available system memory.
3. Dimension of the initial approximation subspace (denoted by s). Default value is set to $s := 2k$.
4. Convergence tolerance (denoted by tol). Default value is set to $\text{tol} := 1e - 3$.

The wall-clock time of TeraPCA is affected by all of the above parameters. Clearly, reducing tol or increasing k results in an increase of the wall-clock time. Using a higher-dimensional approximation subspace, i.e. increasing s , might reduce the corresponding wall-clock time as it typically enhances convergence toward the k -leading eigenvectors. On the other hand, increasing the value of s also increases the amount of floating-point operations performed. Finally, since only a part of the dataset can fit in the system memory at any time unit, the choice of β is typically determined automatically by TeraPCA based on the size of the system memory. The total amount of time spent on I/O is largely independent of the value of β but we have observed that the value of β has an effect on the wall-clock time of the LAPACK routines.

3 Results and discussion

The performance of TeraPCA was tested on both simulated and real-world genotypic datasets. All our experiments were performed at Purdue's Brown cluster on a dedicated node which features an Intel Xeon Gold 6126 processor running at 2.6 GHz with 96 GB of Random Access Memory and a 64-bit CentOS Linux 7 operating system. [Table 1](#) lists the number of samples, number of SNPs, and size of each dataset. Datasets S_1 through S_7 are synthetic datasets

Table 1. Our data sets (simulated and real)

Dataset	Size (.PED file)	Size (.BED file)	# Samples	# SNPs
S_1 (simulated)	19 GB	120 MB	5000	1 000 000
S_2 (simulated)	38 GB	239 MB	10 000	1 000 000
S_3 (simulated)	373 GB	24 GB	100 000	1 000 000
S_4 (simulated)	1.9 TB	117 GB	500 000	1 000 000
S_5 (simulated)	3.7 TB	233 GB	1 000 000	1 000 000
S_6 (simulated)	38 GB	2.4 GB	100 000	100 000
S_7 (simulated)	150 GB	9.4 GB	2000	20 000 000
HGDP	615 MB	39 MB	1043	154 417
1000 Genomes	8.4 GB	483 MB	2504	808 704
PRK	2 GB	126 MB	4706	111 831
T2D	1.8 GB	111 MB	6370	72 457

and the remaining ones are real-world datasets. This section provides comparisons between TeraPCA and FlashPCA2. The latter has already been shown to be faster than previous methods such as FlashPCA (Abraham and Inouye, 2014), FastPCA (Galinsky *et al.*, 2016), etc. The results reported throughout the remainder of this section were obtained by setting the amount of system memory made available to TeraPCA (as well as FlashPCA2) to 2 GBs. This is precisely the amount of memory allowed to FlashPCA2 in prior work.

3.1 Synthetic datasets

Datasets S_1 through S_5 in [Table 1](#) have a fixed number of SNPs (=1 million) and a varying number of samples (from 5 000 to 1 million). On the other hand, dataset S_6 was used to fine-tune prior state-of-the-art methods and contains 100 000 samples genotyped on 100 000 SNPs. S_7 was used to test the performance of TeraPCA on extremely rectangular matrices, where the number of SNPs heavily outnumbers the number of individuals.

We first consider the plots of the three leading PCs returned by both TeraPCA and FlashPCA2 for dataset S_6 (see [Supplementary Fig. S1](#)). TeraPCA and FlashPCA2 show a complete visual agreement with each other and both libraries agree with the expected outcome of the PSD model. For this particular example, TeraPCA terminated in just under 40 min, while FlashPCA2 required 141 min [to be fair in our comparisons between TeraPCA and FlashPCA2, we performed multiple runs of FlashPCA2 on dataset S_6 in order to explore and understand its properties. In particular, we varied the convergence criterion in FlashPCA2 and recorded the resulting trade-off between wall-clock time and digits of accuracy for the top 10 computed eigenvalues. Fixing the convergence tolerance in FlashPCA2 to three digits of accuracy and the maximum number of iterations of FlashPCA2 to 100 was the best choice in terms of the trade-off between running time and accuracy (see [Supplementary Material](#) for more details)].

[Table 2](#) lists the wall-clock times achieved by TeraPCA when applied on datasets S_1 through S_7 . For datasets S_4 and S_5 , which were the largest ones in our collection, TeraPCA terminated after 7.3 and 13.2 h, respectively. On the other hand, FlashPCA2 did not terminate within the 50 h limit that we imposed. TeraPCA outperformed FlashPCA2 on all synthetic datasets, with a speed-up that ranged between 1.3 and 4.5, at least for those datasets where FlashPCA2 terminated within our 50 h limit. We note that for all synthetic datasets the leading PCs returned by TeraPCA and FlashPCA2 showed perfect correlation as measured by the Pearson correlation coefficient (=1 in all cases). To further test TeraPCA's performance on datasets where the number of SNPs heavily

Table 2. Wall-clock running times comparisons for the datasets of Table 1 using a single thread and 2 GB of system memory (* indicates no convergence after 50 h)

Dataset	TeraPCA	FlashPCA2	Speed-up
S_1	26.2 min	33.3 min	1.27
S_2	39.3 min	87.5 min	2.22
S_3	7.9 h	35.6 h	4.50
S_4	7.3 h	n/a*	∞
S_5	13.2 h	n/a*	∞
S_6	39.5 min	141.1 min	3.57
S_7	37.3 min	106.5 min	2.86
HGDP	6.5 sec	7.7 s	1.22
1000 Genomes	4.3 min	3.5 min	0.81
T2D	96 s	119 s	1.24
PRK	76 s	73 s	0.96

outnumbers the number of individuals, we applied it to S_7 and observed that even in a heavily under-determined system, TeraPCA outperformed FlashPCA2 by a factor of 2.9, with similar accuracy guarantees.

3.2 Real datasets

We first considered the HGDP dataset (Cann et al., 2002). TeraPCA was marginally faster than FlashPCA2 and both libraries required about 7 s. A plot of the projection of the HGDP dataset along the two leading PCs computed by TeraPCA is shown in Supplementary Figure S2. Given the relatively small size of this dataset, we were able to compute the exact 10 leading eigenvectors using LAPACK. Figure 1 reports the entry-wise error of the 10 leading eigenvectors returned by TeraPCA. As expected, eigenvectors associated with the largest eigenvalues are captured more accurately since they converge faster.

In addition, Supplementary Table S1 reports the relative and absolute errors of the 10 leading eigenvalues returned by TeraPCA and FlashPCA2. For TeraPCA, the (much) higher accuracy in the approximation of the 3–4 leading eigenvalues is due to the fact that these approximate eigenvalues kept improving as Randomized Subspace Iteration kept iterating to approximate the trailing eigenvalues and eigenvectors. On the other hand, the accuracy in the approximation of the eigenvalues returned by FlashPCA2 was somewhat uniform for all eigenvalues.

TeraPCA and FlashPCA2 showed similar qualitative and computational performance on the pruned 1000 Genomes dataset (see Fig. 2), with FlashPCA2 terminating slightly faster than TeraPCA. Notice that this dataset is also the one in which the number of SNPs outnumbered the number of individuals by the largest factor.

PCA is an essential tool to detect population stratification in GWAS. In order to evaluate TeraPCA's performance on real-world case-control studies, we applied it on Wellcome Trust Case Control Consortium's T2D and PRK datasets. Like other real-world datasets, both FlashPCA2 and TeraPCA performed similarly, needing roughly the same wall-clock time. Execution of TeraPCA on these datasets can also be done in-core, as they fit in the system memory, leading to comparatively faster computation.

3.3 Multi-threading

The wall-clock times of TeraPCA and FlashPCA2 can significantly improve by executing the associated linear algebra computations using more than one threads. This is indeed the most obvious way to speed-up software such as ours. To test the performance of TeraPCA as a function of the number of threads, we focused on datasets

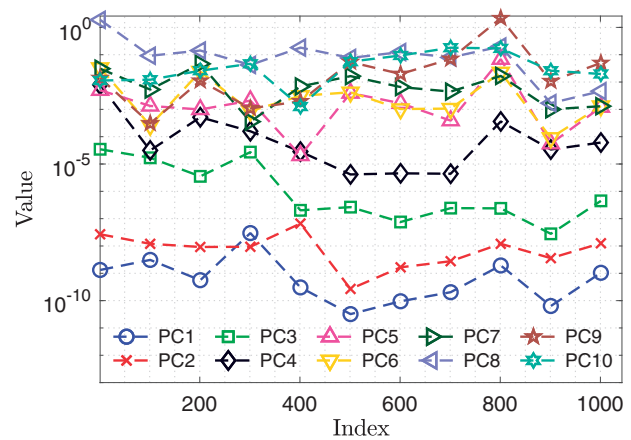


Fig. 1. Entry-wise relative error of the top 10 leading eigenvectors returned by TeraPCA for the HGDP dataset, compared to the eigenvectors returned by LAPACK. The y-axis shows the relative error; recall that each eigenvector has 1043 entries. We observe that the relative error is roughly the same for each entry of a specific eigenvector

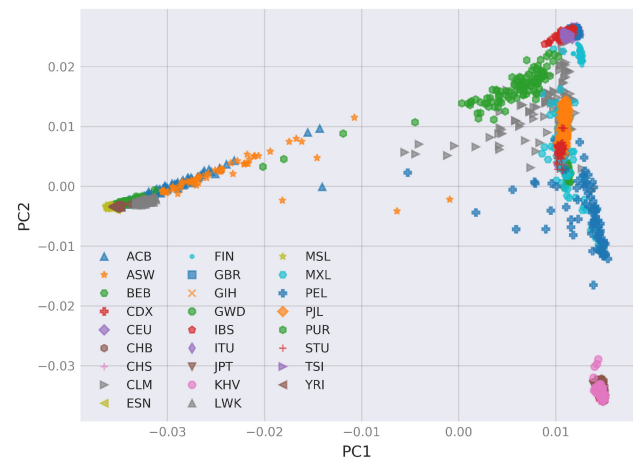


Fig. 2. Projection of the samples of the 1000 Genomes dataset on the top two left singular vectors (PC1 and PC2), as computed by TeraPCA

S_1 , S_2 , S_4 , S_6 , S_7 , and the 1000 Genomes dataset. The number of threads was set to 4, 8 and 12 and the speedups reported in Figure 3 are against the single-thread execution of TeraPCA. Generally speaking, we observed a $1.6\times$ to $2.8\times$ speed-up, which is somewhat sub-optimal. The reason underlying this non-optimality is that we used multi-threading only for the linear algebra operations. However, much of the wall-clock time is spent on I/O operations in order to load the dataset from secondary memory, a procedure that cannot be multi-threaded. We emphasize that FlashPCA2 did not demonstrate comparable improvements when multi-threading was enabled. In particular, when applied to the dataset S_6 , the wall-clock time of FlashPCA2 reduced only by 2 min, i.e. from 141 to 139 min.

In all of the above experiments we set $s := 2k$ and $k := 10$. Finally, Supplementary Figure S6 reports the amount of time required to multiply the (normalized) covariance matrix by a set of s vectors using the DGEMM BLAS routine of MKL and a varying number of threads for different values of s and β for datasets S_6 and HGDP. It is worth noting that while an exhaustive analysis lies outside the goals of this paper, it is easy to verify that doubling the value of s does not double the amount of time required to perform the multiplication, while larger values of s also lead to higher

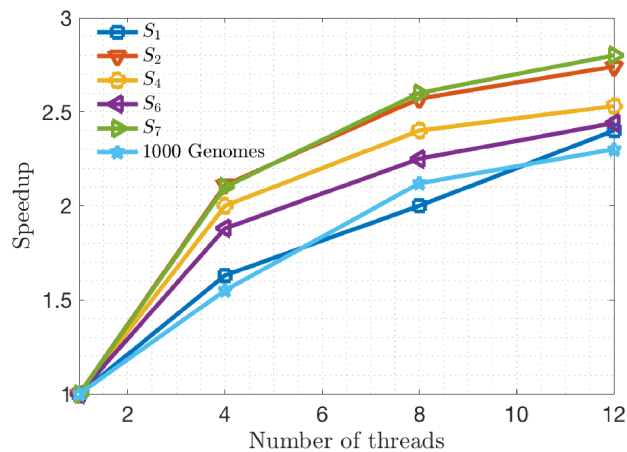


Fig. 3. Speed-up of TeraPCA over single-threaded execution

speedups when multiple threads are used. Similarly, very small values of β are likely to penalize the performance of DGEMM due to non-optimal cache utilization.

4 Summary and future work

In this paper we presented TeraPCA, a C++ library to perform out-of-core PCA analysis of massive genomic datasets. It is based on Randomized Subspace Iteration, building upon principled and theoretically sound methods to approximate the top PCs of massive covariance matrices. TeraPCA returns highly accurate approximations to the top PCs, while taking advantage of modern computer architectures that support multi-threading and it has minimal dependencies to external libraries. TeraPCA can be applied both in-core and out-of-core and is able to successfully operate even on personal workstations with a system memory of just a few gigabytes. Numerical experiments performed on synthetic and real datasets demonstrate that TeraPCA performs similarly or better when compared to state-of-the-art software packages such as FlashPCA2, on a single thread and significantly better with multi-threading.

Future work will focus on implementing a distributed memory version of TeraPCA using the Message Passing Interface standard. Another interesting research direction would be to combine TeraPCA with block Krylov subspace techniques.

Acknowledgements

The authors are grateful to P. Gopalan and W. Hao for sharing their R script to generate the simulated datasets as well as for their valuable comments. The authors are also thankful to P. Anappindi for contributing to the code in the pilot phase.

Author contributions

A.B., V.K., E.K. and P.D. conceived and designed the work. A.B., V.K. and E.K. developed the TeraPCA C++ package. E.K., M.E. and A.B. developed the C++ package to generate the simulated datasets from the PSD model. P.D. and P.P. participated in and discussed analyses. A.B. and V.K. ran the experiments. A.B., V.K., E.K., P.D. and P.P. wrote and revised the manuscript.

Funding

This work was partially supported by the National Science Foundation [IIS-1661760, IIS-1661756, IIS-1715202 to P.P. and P.D.].

Conflict of Interest: none declared.

References

- Abraham, G. and Inouye, M. (2014) Fast principal component analysis of large-scale genome-wide data. *PLoS One*, **9**, 1–5.
- Abraham, G. *et al.* (2017) FlashPCA2: principal component analysis of Biobank-scale genotype datasets. *Bioinformatics*, **33**, 2776–2778.
- Alexander, D.H. *et al.* (2009) Fast model-based estimation of ancestry in unrelated individuals. *Genome Res.*, **19**, 1655–1664.
- Anderson, E. *et al.* (1999) *LAPACK Users' Guide*. 3rd edn. Society for Industrial and Applied Mathematics, Philadelphia, PA.
- Bose, A. *et al.* (2017) Dissecting Population Substructure in India via Correlation Optimization of Genetics and Geodemographics. *bioRxiv*.
- Cann, H.M. *et al.* (2002) A human genome diversity cell line panel. *Science*, **296**, 261–262.
- Chisholm, B. *et al.* (1995) The history and geography of human genes. *J. Asian Stud.*, **54**, 490.
- Drineas, P. and Mahoney, M.W. (2016) RandNLA: randomized numerical linear algebra. *Commun. ACM*, **59**, 80–90.
- Drineas, P. and Mahoney, M.W. (2018) Lectures on randomized numerical linear algebra. In: *The Mathematics of Data*, IAS/Park City Mathematics Series. Vol. 25. American Mathematical Society, Providence, RI, pp. 1–45. <https://dblp.org/rec/bib/journals/corr/abs-1712-08880>.
- Drineas, P. *et al.* (2018) Structural convergence results for low-rank approximations from block Krylov spaces. *SIAM J. Matrix Anal. Appl.*, **39**, 567–586.
- Galinsky, K.J. *et al.* (2016) Fast principal-component analysis reveals convergent evolution of ADH1B in Europe and East Asia. *Am. J. Hum. Genet.*, **98**, 456–472.
- Gopalan, P. *et al.* (2016) Scaling probabilistic models of genetic variation to millions of humans. *Nat. Genet.*, **48**, 1587–1590.
- Halko, N. *et al.* (2011) Finding structure with randomness: probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Rev.*, **53**, 217–288.
- Hotelling, H. (1933) Analysis of a complex of statistical variables into principal components. *J. Educ. Psychol.*, **24**, 417–441.
- Hotelling, H. (1936) Relations between two sets of variates. *Biometrika*, **28**, 321–377.
- Menozi, P. *et al.* (1978) Synthetic maps of human gene frequencies in Europeans. *Science*, **201**, 786–792.
- Musco, C. and Musco, C. (2015) Randomized block Krylov methods for stronger and faster approximate singular value decomposition. In: Cortes, C. *et al.* (eds) *Advances in Neural Information Processing Systems 28*. Curran Associates, Inc., Montreal, Canada, pp. 1396–1404.
- Novembre, J. *et al.* (2008) Genes mirror geography within Europe. *Nature*, **456**, 98.
- Parlett, B. (1998) *The Symmetric Eigenvalue Problem*. Society for Industrial and Applied Mathematics, Philadelphia, PA.
- Paschou, P. *et al.* (2007) PCA-correlated SNPs for structure identification in worldwide human populations. *PLoS Genet.*, **3**, 1–15.
- Paschou, P. *et al.* (2008) Tracing sub-structure in the European American population with PCA-informative markers. *PLoS Genet.*, **4**, 1–13.
- Paschou, P. *et al.* (2014) Maritime route of colonization of Europe. *Proc. Natl. Acad. Sci. USA*, **111**, 9211–9216.
- Patterson, N. *et al.* (2006) Population structure and eigenanalysis. *PLoS Genet.*, **2**, 1–20.
- Pearson, K. (1901) On lines and planes of closest fit to systems of points in space. *Lond. Edimb. Dubl. Phil. Mag.*, **2**, 559–572.
- Price, A.L. *et al.* (2006) Principal components analysis corrects for stratification in genome-wide association studies. *Nat. Genet.*, **38**, 904.
- Price, A.L. *et al.* (2010) New approaches to population stratification in genome-wide association studies. *Nat. Rev. Genet.*, **11**, 459–463.
- Pritchard, J.K. *et al.* (2000) Inference of population structure using multilocus genotype data. *Genetics*, **155**, 945–959. [pmid].
- Rokhlin, V. *et al.* (2010) A randomized algorithm for principal component analysis. *SIAM J. Matrix Anal. Appl.*, **31**, 1100–1124.
- Saad, Y. (2011) *Numerical Methods for Large Eigenvalue Problems*. Society for Industrial and Applied Mathematics, Philadelphia, PA.
- Wang, C. *et al.* (2010) Comparing spatial maps of human population-genetic variation using procrustes analysis. *Stat. Appl. Genet. Mol. Biol.*, **9**, 13.
- Weir, B.S. and Cockerham, C.C. (1984) Estimating f-statistics for the analysis of population structure. *Evolution*, **38**, 1358–1370.