

Scalable and flexible probabilistic PCA for large-scale genetic variation data

Sriram Sankararaman

Department of Computer Science
Department of Human Genetics
UCLA

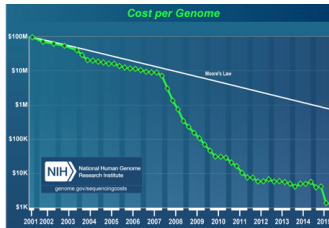
The genomic revolution

Scalable
Probabilistic
PCA

PCA

Probabilistic
PCA

Results



Challenge: scalable methods to analyze and visualize this data

Genetic data

Scalable
Probabilistic
PCA

PCA

Probabilistic
PCA

Results

	Individual 2		
	Individual 1		Individual 3
	.	.	.
	.	.	.
	.	.	.
	.	.	.
	A A	A G	A G
	.	.	.
	.	.	.
	C T	C T	C C
	.	.	.
	.	.	.
	.	.	.

Position along Genome ↓

	Individual 2			Genotype 2		
	Individual 1	Individual 3		Genotype 1	Genotype 3	
	A A	AG	AG	0	1	1
	C T	CT	CC	1	1	2

SNPs ↓

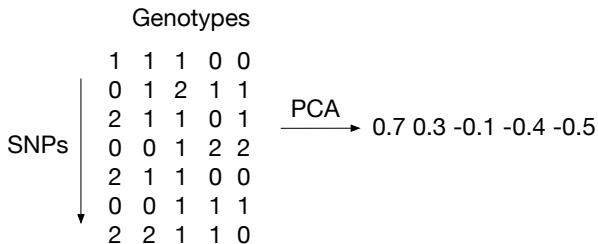
PCA on genetic data

Scalable
Probabilistic
PCA

PCA

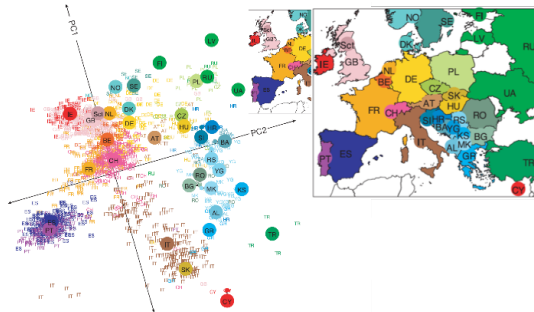
Probabilistic
PCA

Results



PCA on genetic data

Visualize genetic structure



PCA

Scalable
Probabilistic
PCA

PCA

Probabilistic
PCA

Results

Given N genotype vectors over M SNPs
 $\mathbf{x}_n \in \mathbb{R}^M, n \in \{1, \dots, N\}$ and $K \leq M$.

$$\begin{aligned}\mathbf{x}_n &\approx \mathbf{w}_1 z_{n,1} + \dots + \mathbf{w}_K z_{n,K} \\ &= [\mathbf{w}_1 \dots \mathbf{w}_K] \begin{bmatrix} z_{n,1} \\ \vdots \\ z_{n,K} \end{bmatrix} \\ &= \mathbf{W} \mathbf{z}_n\end{aligned}$$

- **PCA Constraint:** columns of \mathbf{W} are orthonormal.
- The PCA solution $\widehat{\mathbf{W}} = \mathbf{U}_K$ where \mathbf{U}_K contains the top K eigenvectors of the sample covariance matrix.

Challenges with PCA

Computational

- Compute all eigenvalue, eigenvectors.
 - Singular-Value Decomposition (SVD):

$$\mathcal{O}(MN \min(M, N)) \approx \mathcal{O}(MN^2)$$

- Infeasible for genetic datasets (large number of SNPs M or individuals N).
- Recent Randomized approximation algorithms

Statistical

- Missing genotypes.
- Correlation among SNPs.

Probabilistic PCA

Scalable
Probabilistic
PCA

Model

$$\begin{aligned} \mathbf{z}_n &\stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_K) \\ p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{W}, \sigma^2) &= \mathcal{N}(\mathbf{W}\mathbf{z}_n, \sigma^2 \mathbf{I}_M) \end{aligned}$$

Log likelihood

$$\mathcal{LL}(\mathbf{W}, \sigma^2) \equiv \log P(\mathbf{X} | \mathbf{W}, \sigma^2) = \log \prod_{n=1}^N p(\mathbf{x}_n, \mathbf{z}_n | \mathbf{W}, \sigma^2)$$

The maximum likelihood estimator is equivalent to PCA.

PCA

Probabilistic
PCA

Results

Probabilistic PCA

Scalable
Probabilistic
PCA

PCA

Probabilistic
PCA

Results

EM algorithm

- E-step:

$$\mathbf{Z} = (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T \mathbf{X}$$

- M-step:

$$\mathbf{W} = \mathbf{X} \mathbf{Z}^T (\mathbf{Z} \mathbf{Z}^T)^{-1}$$

- Assume: $\sigma^2 \rightarrow 0$

Probabilistic PCA

Scalable
Probabilistic
PCA

EM algorithm: computational complexity

- E-step:

$$\mathbf{Z} = \underbrace{(\mathbf{W}^T \mathbf{W})^{-1}}_{K \times K} \underbrace{\mathbf{W}^T}_{K \times M} \underbrace{\mathbf{X}}_{M \times N}$$

$$\mathcal{O}(NMK)$$

- M-step:

$$\mathbf{W} = \underbrace{\mathbf{X}}_{M \times N} \underbrace{\mathbf{Z}^T}_{N \times K} \underbrace{(\mathbf{Z}\mathbf{Z}^T)^{-1}}_{K \times K}$$

$$\mathcal{O}(NMK)$$

PCA

Probabilistic
PCA

Results

Probabilistic PCA

Scalable
Probabilistic
PCA

PCA

Probabilistic
PCA

Results

EM algorithm: computational complexity

- Run for I iterations with each iteration costing $\mathcal{O}(NMK)$.
- For small K , leads to a linear-time algorithm.

Probabilistic PCA

Scalable
Probabilistic
PCA

PCA

Probabilistic
PCA

Results

EM algorithm: computational complexity

- Run for I iterations with each iteration costing $\mathcal{O}(NMK)$.
- For small K , leads to a linear-time algorithm.
- Ignores the special structure of the genotype matrix \mathbf{X} .

Probabilistic PCA

Scalable
Probabilistic
PCA

EM algorithm: computational complexity

Each E and M step, perform the following operations K times:

$$\mathbf{c} = \mathbf{X}\mathbf{b}$$

- \mathbf{X} is a fixed $M \times N$ matrix of genotypes.
- \mathbf{b} is a real-valued vector that could potentially change each iteration.
- Naive multiplication takes $\mathcal{O}(NM)$.

PCA

Probabilistic
PCA

Results

Probabilistic PCA

Scalable
Probabilistic
PCA

EM algorithm: computational complexity

Each E and M step, perform the following operations K times:

$$\mathbf{c} = \mathbf{X}\mathbf{b}$$

- \mathbf{X} is a fixed $M \times N$ matrix of genotypes.
- \mathbf{b} is a real-valued vector that could potentially change each iteration.
- Naive multiplication takes $\mathcal{O}(NM)$.
- For a genotype matrix, can we do some pre-processing so that $\mathbf{X}\mathbf{b}$ can be computed more efficiently ?

PCA

Probabilistic
PCA

Results

Probabilistic PCA

EM algorithm: computational complexity

Each E and M step, perform the following operations K times:

$$\mathbf{c} = \mathbf{X}\mathbf{b}$$

- \mathbf{X} is a fixed $M \times N$ matrix of genotypes.
- \mathbf{b} is a real-valued vector that could potentially change each iteration.
- Naive multiplication takes $\mathcal{O}(NM)$.
- For a genotype matrix, can we do some pre-processing so that $\mathbf{X}\mathbf{b}$ can be computed more efficiently ?
- Yes! For a matrix with binary entries: $\mathcal{O}\left(\frac{MN}{\log_2(N)}\right)$.

EM with the Mailman algorithm

Scalable
Probabilistic
PCA

PCA

Probabilistic
PCA

Results

- Entries in genotype matrix take one of three values: $\{0, 1, 2\}$.
- Using the Mailman algorithm, per-iteration time complexity of EM for genotype matrix: $\mathcal{O}\left(\frac{MNK}{\log_3(N)}\right)$.

Sub-linear time algorithm for computing PCA

Simulations

Accuracy

50,000 SNPs, 10,000 individuals

F_{st}	MEV	
	$K = 5$	$K = 10$
0.001	0.987	1.000
0.002	0.999	1.000
0.003	0.999	1.000
0.004	0.999	1.000
0.005	1.000	1.000
0.006	1.000	1.000
0.007	1.000	1.000
0.008	1.000	1.000
0.009	1.000	1.000
0.010	1.000	1.000

Scalable
Probabilistic
PCA

PCA

Probabilistic
PCA

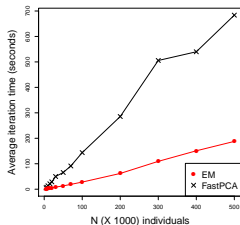
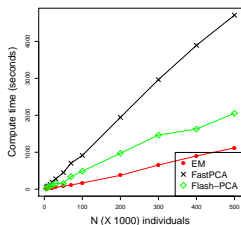
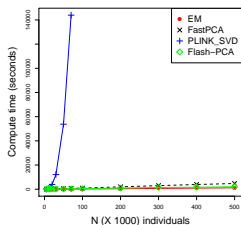
Results

Simulations

Scalable
Probabilistic
PCA

Efficiency

$M = 100,000$ SNPs, $K = 5$, $F_{ST} = 0.01$



PCA

Probabilistic
PCA

Results

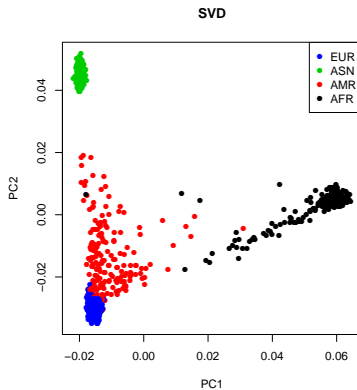
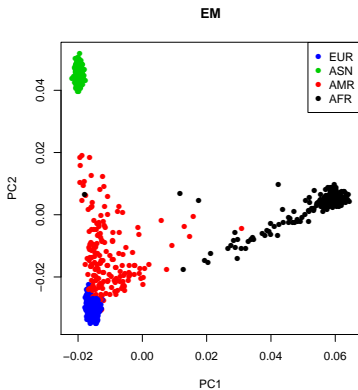
Application to 1000 Genomes data

Scalable
Probabilistic
PCA

PCA

Probabilistic
PCA

Results



Other advantages of Probabilistic PCA

Scalable
Probabilistic
PCA

PCA

Probabilistic
PCA

Results

Can naturally handle missing data

- E-step involves inferring hidden variables \mathbf{Z} as well as hidden (missing observations).
- Can handle missing data efficiently.

Can use model selection to infer K .

- Choose K to maximize the marginal likelihood $P(\mathbf{X}|K)$.
- Use cross-validation and pick K that maximizes likelihood on held out data.

Open questions

Scalable
Probabilistic
PCA

PCA

Probabilistic
PCA

Results

- Modeling correlations .
- Beyond Gaussian outputs

Baran et al. 2013, Wen and Stephens 2012
Collins et al. 2002

Summary

Scalable
Probabilistic
PCA

PCA

Probabilistic
PCA

Results

- PCA can be interpreted as a latent variable model with continuous latent variable.
 - Probabilistic interpretation useful to generalize PCA.
 - Leads to efficient inference.
- Fast matrix-vector multiplication for genotype data more generally applicable and can lead to sub-linear time algorithms.

Acknowledgments

Scalable
Probabilistic
PCA

PCA

Probabilistic
PCA

Results

- Aman Agrawal
- Minh Le
- Eran Halperin

Email sriram@cs.ucla.edu