

Genetics and population analysis

# FlashPCA: fast sparse canonical correlation analysis of genomic data

Gad Abraham<sup>1,2\*</sup> and Michael Inouye<sup>1,2</sup>

<sup>1</sup> Centre for Systems Genomics, School of BioSciences, University of Melbourne, Parkville 3010, VIC, Australia.

<sup>2</sup> Department of Pathology, Faculty of Medicine, Dentistry, and Health Sciences, University of Melbourne, Parkville 3010, VIC, Australia.

\*To whom correspondence should be addressed.

Associate Editor: XXXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

## Abstract

**Summary:** Sparse canonical correlation analysis (SCCA) is a useful approach for correlating one set of measurements, such as single nucleotide polymorphisms (SNPs), with another set of measurements, such as gene expression levels. We present a fast implementation of SCCA, enabling rapid analysis of hundreds of thousands of SNPs together with thousands of phenotypes. Our approach is implemented both as an R package `flashpcaR` and within the standalone commandline tool `flashpca`.

**Availability and implementation:** <https://github.com/gabraham/flashpca>

**Contact:** gad.abraham@unimelb.edu.au

**Supplementary information:** Supplementary data are available at *Bioinformatics* online.

## 1 Introduction

Canonical correlation analysis (CCA) is a well-known statistical approach for multivariate analysis of two datasets (Hotelling, 1936). In the context of large-scale genomic and multi-omic analyses, CCA can prove useful in identifying relationships amongst complex data, for example single nucleotide polymorphisms (SNPs) and gene expression levels. Approaches that consider one SNP at a time together with multiple phenotypes have been shown to increase power to detect QTLs over the simpler but commonly utilised single-SNP/single-phenotype approach (Ferreira *et al.*, 2009; Inouye *et al.*, 2012), particularly when analysing correlated phenotypes that are modulated by the same genetic variants.

Analysis of multiple SNPs simultaneously is an attractive extension of the single-SNP multiple-phenotype approach, however, standard CCA is not well-defined when the number of samples is lower than the number of SNPs or phenotypes ( $n < \min\{p, m\}$ ). One solution is Sparse CCA (SCCA) (Witten *et al.*, 2009a,b; Parkhomenko *et al.*, 2009), an  $L_1$ -penalised variant of CCA which allows for tuning the number of variables that effectively contribute to the canonical correlation, thus making the problem well-defined. Owing to the induced sparsity, SCCA can be useful for identifying a small subset of SNPs and a small subset of the phenotypes exhibiting strong correlations. However, the rapidly increasing size and coverage of genotyping arrays (exacerbated by genotype imputation), together with the availability of large phenotypic datasets (transcriptomic,

metabolomic, and others; e.g., Bartel *et al.* (2015); The GTEx Consortium (2015)), makes it challenging to perform such analyses using existing tools.

We have developed an efficient implementation of SCCA that is capable of analysing genome-wide SNP datasets (1 million SNPs or more) together with thousands of phenotypes, as part of the tool `flashpca` (Abraham *et al.*, 2014). The tool is implemented in C++ using the Eigen 3 numerical library (Guennebaud *et al.*, 2010), as well as an R interface (package `flashpcaR`) based on RcppEigen (Bates *et al.*, 2013).

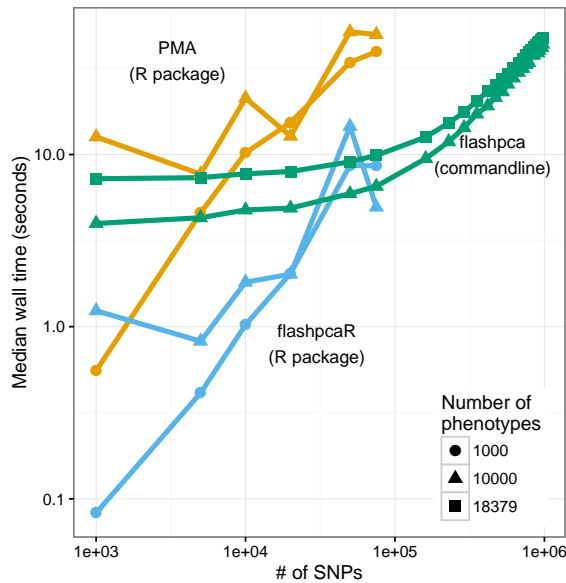
Here, we compare the SCCA implementation in `flashpcaR` and `flashpca` with a widely-used implementation (PMA, by Witten *et al.* (2013)), and demonstrate the substantial improvements in speed of our tool, allowing for large analyses to be performed rapidly.

## 2 Methods

In standard CCA, we assume that we have two matrices  $\mathbf{X}$  ( $n \times p$ ) and  $\mathbf{Y}$  ( $n \times m$ ), measured for the same  $n$  samples. We further assume that both  $\mathbf{X}$  and  $\mathbf{Y}$  have been column-wise standardised (zero mean, unit variance). For a single pair of canonical variables  $a$  and  $b$ , CCA involves solving the problem

$$\arg \max_{a,b} \frac{a^T \Sigma_{XY} b}{\sqrt{a^T \Sigma_{XX} a b^T \Sigma_{YY} b}}, \quad (1)$$

where  $\Sigma_{XX}$  and  $\Sigma_{YY}$  are the covariance matrices of  $\mathbf{X}$  and  $\mathbf{Y}$ , respectively, and  $\Sigma_{XY}$  is the covariance matrix of  $\mathbf{X}$  and  $\mathbf{Y}$ . The solution is given by the singular value decomposition (SVD) of



**Fig. 1.** Timing (median of 30 runs) of SCCA implemented in (i) the `flashpcaR` (R package) and (ii) `flashpca` (stand-alone commandline tool), compared with PMA, using subsets of the HapMap3 dataset with gene expression levels as phenotypes. The stand-alone `flashpca` timing includes data loading into memory.

$\Sigma_{XX}^{-1/2} \Sigma_{XY} \Sigma_{YY}^{-1/2}$ , with  $a = \Sigma_{XX}^{-1/2} u_1$  and  $b = \Sigma_{YY}^{-1/2} v_1$ , where  $u_1$  and  $v_1$  are the first left and right singular vectors, respectively.

SCCA is typically used for high-dimensional data, where a useful assumption is that the columns of  $\mathbf{X}$  and  $\mathbf{Y}$  are uncorrelated, i.e.,  $\Sigma_{XX} = \Sigma_{YY} = \mathbf{I}$  (Parkhomenko *et al.*, 2009), hence,  $a = u$  and  $b = v$ . Thus, SCCA involves solving another form of CCA,

$$\begin{aligned} \arg \max_{u,v} \quad & u^T \Sigma_{XY} v \\ \text{s.t.} \quad & \|u\|_2^2 = 1, \|v\|_2^2 = 1, \|u\|_1 \leq s_u, \|v\|_1 \leq s_v, \end{aligned} \quad (2)$$

where  $u$  and  $v$  are the left and right canonical vectors, respectively, and  $s_u$  and  $s_v$  are constraints on the  $L_1$  norms of these canonical vectors.

The problem can be converted into the penalised (Lagrangian) form, making it solvable using iterative soft-thresholding (Parkhomenko *et al.*, 2009), which we employ here. Unlike standard CCA, SCCA is well-defined even when  $n < \min\{p, m\}$ , and induces sparse canonical vectors, depending on the choice of  $L_1$  penalties (higher penalties lead to higher sparsity). The optimal set of penalties can be found via cross-validation: the data (both  $\mathbf{X}$  and  $\mathbf{Y}$ ) are split into training and test sets, SCCA is run on the training set ( $\mathbf{X}_{\text{train}}, \mathbf{Y}_{\text{train}}$ ) using a 2D grid of penalties, and the pair of penalties that produce the highest correlations in the test set,  $\text{Cor}(\mathbf{X}_{\text{test}} u, \mathbf{Y}_{\text{test}} v)$ , are selected. Optionally, a new model may be fitted to the entire data using these penalties.

### 3 Results

We utilised the HapMap3 phase III genotypes (International HapMap 3 Consortium, 2010), together with gene expression data of 709 individuals (Stranger *et al.*, 2012), with 1.4M SNPs and 47,000 gene expression probes (out of which 21,800 probes were analysed by Stranger *et al.* (2012) and used in our analysis). After quality control (see Supplementary Material) and taking the intersection of SNPs across the populations, the data consisted of 709 individuals, 973,983 SNPs, and 18,379 gene expression probes.

We first confirmed that `flashpcaR::scca` produced models comparable with `PMA::CCA`, by comparing the results in cross-validation on HapMap3 genotypes together with simulated gene expression

levels (Supplementary Material). Next, to further assess the relative speed improvement using real-world data, we used subsets of the HapMap3 genotypes and real gene expression levels (Stranger *et al.*, 2012) and compared the runtime of: (i) `PMA::CCA` (R package), (ii) `flashpcaR::scca` (R package), and (iii) `flashpca` (commandline tool). Whereas both `PMA::CCA` and `flashpcaR::flashpca` are bound by the memory limitations of R, the commandline tool `flashpca` allows much larger analyses; hence we also ran larger analyses of increasing size: chromosomes 1–2, 1–3, ..., 1–22, up to all 973,983 SNPs. Figure 1 shows that `flashpcaR::scca` was 3–12× faster than `PMA::CCA`, with an analysis of 75,000 SNPs and 10,000 gene expression levels completing in 5s and 50s, respectively. The commandline `flashpca` was faster than `PMA::CCA` as well, and completed an analysis of 709 individuals, 973,983 SNPs and 18,379 gene expression levels in median wall time of ~47s (including all overheads), using ~10GiB of RAM. Note that runtime for the commandline `flashpca` includes all steps such as loading data into RAM, unlike the R version where the data is pre-loaded into R. Performing cross-validation over a grid of penalties will increase these times considerably, in which case we recommend parallelisation over several cores (Supplementary Material).

### 4 Conclusion

`flashpca` provides a fast implementation of sparse canonical correlation analysis, making it possible to rapidly analyse high dimensional datasets. SCCA is available in the R package `flashpcaR`, which enables analysis of metabolomic, transcriptomic, or any other quantitative set of measurements. The commandline version is targeted at SNP/phenotype data, enabling large QTL analyses of >1 million SNPs and thousands of phenotypes, that would otherwise be too large to fit within R.

### Funding

This work has been supported by the NHMRC (grant no. 1062227). GA was supported by an NHMRC Early Career Fellowship (no. 1090462). MI was supported by a Career Development Fellowship co-funded by the NHMRC and Heart Foundation (no. 1061435).

### References

- Abraham, G. *et al.* (2014). Fast principal component analysis of large-scale genome-wide data. *PLoS ONE*, **9**(4), e93766.
- Bartel, J. *et al.* (2015). The human blood metabolome-transcriptome interface. *PLoS Genet*, **11**(6), e1005274.
- Bates, D. *et al.* (2013). Fast and elegant numerical linear algebra using the RcppEigen package. *J Stat Soft*, **52**(5), 1–24.
- Ferreira, M. *et al.* (2009). A multivariate test of association. *Bioinformatics*, **25**, 132–133.
- Guennebaud, G. *et al.* (2010). Eigen v3. <http://eigen.tuxfamily.org>.
- Hotelling, H. (1936). Relations between two sets of variates. *Biometrika*, **28**, 321–377.
- Inouye, M. *et al.* (2012). Novel Loci for Metabolic Networks and Multi-Tissue Expression Studies Reveal Genes for Atherosclerosis. *PLoS Genet*, **8**(8), e1002907.
- International HapMap 3 Consortium (2010). Integrating common and rare genetic variation in diverse human populations. *Nature*, **467**, 52–58.
- Parkhomenko, E. *et al.* (2009). Sparse canonical correlation analysis with application to genomic data integration. *Stat Appl Genet Mol Biol*, **8**(1), Article 1.
- Stranger, B. *et al.* (2012). Patterns of Cis Regulatory Variation in Diverse Human Populations. *PLoS Genet*, **8**(4), e1002639.
- The GTEx Consortium (2015). The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue gene regulation in humans. *Science*, **348**, 648–660.
- Witten, D. *et al.* (2009a). A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics*, **10**, 515–34.
- Witten, D. *et al.* (2009b). Extensions of Sparse Canonical Correlation Analysis with Applications to Genomic Data. *Stat Appl Genet Mol Biol*, **8**(1), Article 29.
- Witten, D. *et al.* (2013). *PMA: Penalized Multivariate Analysis*. R package version 1.0.9.