

## Adaptive testing for association between two random vectors in moderate to high dimensions

Journal:	<i>Genetic Epidemiology</i>
Manuscript ID	GenEpi-17-0023.R1
Wiley - Manuscript type:	Research Article
Date Submitted by the Author:	26-Apr-2017
Complete List of Authors:	Xu, Zhiyuan; University of Minnesota, Xu, Gongjun; University of Michigan Pan, Wei; University of Minnesota, Division of Biostatistics, SPH
Key Words:	dCov test, eQTL, GEE-aSPU test, Multitrait association testing

SCHOLARONE™  
Manuscripts

Review

Author Manuscript

John Wiley & Sons, Inc.

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the Version of Record. Please cite this article as doi: 10.1002/gepi.22059

This article is protected by copyright. All rights reserved.

1  
2  
3 **Adaptive testing for association between two random vectors in**  
4 **moderate to high dimensions**  
5  
6

7 ZHIYUAN XU<sup>1</sup>, GONGJUN XU<sup>2</sup>, WEI PAN<sup>1</sup>, FOR THE ALZHEIMER'S DISEASE NEUROIMAGING  
8  
9 INITIATIVE<sup>3</sup>  
10

11 <sup>1</sup>*Division of Biostatistics, University of Minnesota*

12 <sup>2</sup>*Department of Statistics, University of Michigan*  
13  
14

15 February 27, 2017; revised April 26, 2017  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27

28 Correspondence author: Wei Pan

29 Telephone: (612) 626-2705

30 Fax: (612) 626-0660

31 Email: weip@biostat.umn.edu

32 Address: Division of Biostatistics, MMC 303,

33 School of Public Health, University of Minnesota,

34 Minneapolis, Minnesota 55455-0392, U.S.A.  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48

49 <sup>3</sup>Data used in preparation of this article were obtained from the Alzheimer's Disease Neuroimag-  
50 ing Initiative (ADNI) database ([adni.loni.usc.edu](http://adni.loni.usc.edu)). As such, the investigators within the ADNI  
51 contributed to the design and implementation of ADNI and/or provided data but did not participate  
52 in analysis or writing of this report. A complete listing of ADNI investigators can be found at: [http:](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf)  
53 [http:](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf)  
54 [http:](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf)  
55 [http:](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf)  
56 [http:](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf)  
57 [http:](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf)  
58 [http:](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf)  
59 [http:](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf)  
60

## Abstract

Testing for association between two random vectors is a common and important task in many fields, however, existing tests, such as Escoufier’s RV test, are suitable only for low-dimensional data, not for high-dimensional data. In moderate to high dimensions, it is necessary to consider sparse signals, which are often expected with only a few, but not many, variables associated with each other. We generalize the RV test to moderate to high dimensions. The key idea is to data-adaptively weight each variable pair based on its empirical association. As the consequence, the proposed test is adaptive, alleviating the effects of noise accumulation in high-dimensional data, and thus maintaining the power for both dense and sparse alternative hypotheses. We show the connections between the proposed test with several existing tests, such as a generalized estimating equationsG-based adaptive test, multivariate kernel machine regression, and kernel distance methods. Furthermore, we modify the proposed adaptive test so that it can be powerful for non-linear or non-monotonic associations. We use both real data and simulated data to demonstrate the advantages and usefulness of the proposed new test. The new test is freely available in R package aSPC at <https://github.com/jasonzyx/aSPC>.

*Key words:* aSPC test, dCov test, eQTL, GEE-aSPU test, RV test

## 1 Introduction

To investigate genetic control of gene expression, it is common and useful to conduct association analysis between single nucleotide polymorphisms (SNPs) and gene expression (i.e. mRNA or transcript) levels, also known as eQTL analysis. This often involves massive univariate testing. For example, Colantuoni et al. (2011) examined 30,176 expression probes and 625,439 SNPs, leading to  $1.89 \times 10^{10}$  (19 billion) possible SNP-gene associations. After the conservative Bonferroni adjustment, only 1,628 individual associations surpassed the genome-wide significance level. However, when they conducted a global test for possible association between all SNPs and all transcripts, no association was detected. They noted: “This dramatic lack of association between genetic distance and transcriptome distance across our sample is a surprising result that requires further interrogation. It is possible that no association is found in Fig. 4 because most of the genetic polymorphisms measured do not impact on gene expression.” We agree with Colantuoni et al. (2011) on the possible reason for the lack of a global association in striking contrast to the presence of some individual associations: it is due to the lack of power of a global test for high-dimensional data with only sparse signals. Furthermore, the authors also commented on that, surprisingly, no association was found even for smaller subsets of the SNPs and genes. We note that their used method

1  
2  
3 52 was Mantel's (1967) test, which was originally proposed for low-dimensional data and may have  
4  
5 53 only limited power for moderate- to high-dimensional data as to be confirmed. Nevertheless, this  
6  
7 54 example pinpoints the importance of conducting global association testing with high-dimensional  
8  
9 55 data, given that most of the existing tests were almost exclusively developed for low-dimensional  
10  
11 56 data for historical reasons, as reviewed in Josse and Holmes (2014).

12 57 Some commonly used tests for association between two random vectors include the RV test  
13  
14 58 (Escoufier, 1970), the Mantel test (Mantel, 1967) and the distance covariance (dCov) test (Székely,  
15  
16 59 Rizzo and Bakirov, 2007). The RV test is based on the RV coefficient as a multivariate generalization  
17  
18 60 of Pearson's correlation coefficient. It is perhaps the most popular one in many fields, especially in  
19  
20 61 ecology. The Mantel test aims to detect a possible correlation between two distance matrices among  
21  
22 62 the subjects based on the two random vectors respectively; it is noted that the Mantel test was used  
23  
24 63 by Colantuoni et al. (2011). The dCov test has only become popular recently due to its attracting  
25  
26 64 property of being consistent in detecting any possible associations, including non-linear and non-  
27  
28 65 monotonic relationships. A common problem with the above tests is their treating all the variables  
29  
30 66 in the two random vectors equally a priori, which is perhaps reasonable for low-dimensional data,  
31  
32 67 but not for moderate- to high-dimensional data: as for the SNP-gene expression data of Colantuoni  
33  
34 68 et al. (2011), most of the SNPs do not have regulatory function; even for those regulatory ones,  
35  
36 69 their targets are likely only a few, not most, of the genes. That is, for high-dimensional data, we  
37  
38 70 expect that many or even most (e.g. SNP-gene) pairs are not associated, which is ignored by the  
39  
40 71 above existing tests, leading to their noise accumulations and thus substantial power loss as to be  
41  
42 72 confirmed in later numerical studies. Hence, to boost power, it is important to conduct variable  
43  
44 73 selection or variable weighting. With weak associations, it is difficult for accurate variable selection,  
45  
46 74 so we take a variable weighting approach. In our approach, we use the data to adaptively determine  
47  
48 75 a weight for each pair of the variables: if a pair is more likely to be associated, we assign a higher  
49  
50 76 weight to it. This will effectively down-weight many of those non-associated pairs, alleviating the  
51  
52 77 effects of noise accumulation hindering most existing tests for high-dimensional data. Our adaptive  
53  
54 78 test can be regarded as a generalization of the RV test to high-dimensional data, as to be shown  
55  
56 79 later.

57  
58  
59 80 We note that the above tests aim to tackle the same problem as SNP-set- or gene-based associ-  
60  
61 81 ation testing for multiple traits or longitudinal traits in genetics (e.g., Maity, Sullivan and Tzeng,  
62  
63 82 2012; He et al., 2015; Fan et al., 2016; Wang, Lee, Zhu, Redline and Lin, 2013; Wang et al.,  
64  
65 83 2015; Wang, Xu, Zhang, Wu and Wang, 2017; Kim, Zhang and Pan, 2016 and references therein),

1  
2  
3  
4 84 but the two lines of research seem to be largely non-overlapping; it is also our goal here to bridge  
5 85 the gap between the two lines of research. In particular, our proposed test is related to another  
6 86 adaptive test, called adaptive sum of powered score test based on generalized estimating equations  
7 87 (GEE-aSPU), originally designed in genetics for testing for multi-trait and multi-SNP associations  
8 88 in low to moderate dimensions (Kim et al., 2016), but we will also show some computational ad-  
9 89 vantages of the proposed test over GEE-aSPU. It is also connected with kernel machine regression  
10 90 and kernel distance methods (Hua and Ghosh, 2015). Furthermore, due to the simplicity of our  
11 91 proposed test, it can be also extended to detect non-linear or even non-monotonic associations by  
12 92 borrowing the idea from the dCov test, though our test is much more powerful than the dCov test  
13 93 for sparse signals in moderate- to high-dimensions.

14 94 The rest of the article is organized as follows. In section 2 we will briefly review the RV test,  
15 95 which serves to motivate our proposed aSPC test. We then outline the connections of the aSPC  
16 96 test to some existing tests before presenting its several generalizations. Section 3 applies the new  
17 97 and some existing tests to an SNP-gene expression dataset drawn from the Alzheimer’s Disease  
18 98 Neuroimaging Initiative (ADNI), highlighting some advantages of the new tests over some existing  
19 99 ones. In section 4 more simulation results are shown to support the power and flexibility of the  
20 100 aSPC test. We end with a summary of the main conclusions in section 5.

## 2 Methods

21 102 Our goal is to test for association between two random vectors  $\mathbf{x}_{p \times 1}$  and  $\mathbf{y}_{q \times 1}$  in  $p$  and  $q$  dimensions  
22 103 respectively. We have  $n$  iid observations on  $\mathbf{x}$ - $\mathbf{y}$  pair as stored in two matrices  $X_{n \times p}$  and  $Y_{n \times q}$ ,  
23 104 respectively; each row of the two matrices corresponds to an observed  $\mathbf{x}$ - $\mathbf{y}$  pair. Denote  $X_l$  as the  
24 105  $l$ th ( $l = 1, \dots, p$ ) column of matrix  $X$  and  $Y_m$  as the  $m$ th ( $m = 1 \dots q$ ) column of  $Y$ . It is assumed  
25 106 throughout that each column of the two matrices is centered at mean 0 with a unit variance. We  
26 107 will use  $X$  and  $Y$  to test for association between  $\mathbf{x}$  and  $\mathbf{y}$ ; with some abuse of notation, we also  
27 108 call it association between  $X$  and  $Y$ .

### 2.1 Review: the RV test

28 110 For the purpose of comparison, we first briefly review the RV test, largely following Josse and  
29 111 Holmes (2014). The two cross-product matrices of  $X$  and  $Y$  are  $W_X = XX^T$  and  $W_Y = YY^T$ ,  
30 112 both of which are of size  $n \times n$ . To measure their proximity, the Hilbert-Schmidt inner product

113 between matrices  $W_X$  and  $W_Y$  can be used:

$$\langle W_X, W_Y \rangle = \text{tr}(XX^TYY^T) = (n-1)^2 \sum_{l=1}^p \sum_{m=1}^q \text{Cov}_n^2(X_{.l}, Y_{.m}), \quad (1)$$

114 where  $\text{Cov}_n(X_{.l}, Y_{.m})$  is the sample covariance between columns  $X_{.l}$  and  $Y_{.m}$ . The RV coefficient,  
115 a correlation coefficient proposed by Escoufier (1973) for two random vectors, is computed by  
116 normalizing the Hilbert-Schmidt inner product by the matrix norms:

$$\text{RV}(X, Y) = \frac{\langle W_X, W_Y \rangle}{\|W_X\| \|W_Y\|} = \frac{\text{tr}(XX^TYY^T)}{\sqrt{\text{tr}(XX^T)^2 \text{tr}(YY^T)^2}}, \quad (2)$$

117 which accounts for possibly different scales of  $\mathbf{x}$  and  $\mathbf{y}$ . The population RV coefficient is  $\rho(\mathbf{x}, \mathbf{y}) =$   
118  $\text{tr}(\Sigma_{\mathbf{x}\mathbf{y}}\Sigma_{\mathbf{y}\mathbf{x}})/\sqrt{\text{tr}(\Sigma_{\mathbf{x}\mathbf{x}}^2)\text{tr}(\Sigma_{\mathbf{y}\mathbf{y}}^2)}$ , where  $\Sigma_{\mathbf{x}\mathbf{y}}$  is the population covariance between  $\mathbf{x}$  and  $\mathbf{y}$ . Our goal  
119 is to test  $H_0 : \rho(\mathbf{x}, \mathbf{y}) = 0$ .

120 If each column of  $X$  and of  $Y$  is standardized to have a zero mean and a unit variance, as always  
121 assumed here, the RV coefficient can be simplified as:

$$\text{RV}(X, Y) = \frac{\text{tr}(XX^TYY^T)}{(n-1)^2 pq} = \frac{\sum_{l=1}^p \sum_{m=1}^q \text{corr}_n^2(X_{.l}, Y_{.m})}{pq} \propto \sum_{l=1}^p \sum_{m=1}^q \text{corr}_n^2(X_{.l}, Y_{.m}), \quad (3)$$

122 where  $\text{corr}_n(X_{.l}, Y_{.m})$  is the sample Pearson correlation coefficient between columns  $X_{.l}$  and  $Y_{.m}$ .

123 A permutation method can be used to calculate the  $P$ -value. Specifically, for each permutation  
124  $b = 1, \dots, B$ , we permute the rows of matrix  $X$  (or  $Y$ ), then calculate the corresponding RV  
125 coefficient  $\text{RV}^{(b)}$ ; the  $P$ -value is calculated as the sample proportion  $[\sum_{n=1}^B I(\text{RV} \leq \text{RV}^{(b)}) +$   
126  $1]/(B+1)$ .

## 127 2.2 New method: an adaptive sum of powered correlation (aSPC) test

128 To generalize the RV coefficient as reformulated in equation (3), we propose a family of so-called  
129 sum of powered correlation (SPC) tests:

$$\text{SPC}(\gamma) = \sum_{l=1}^p \sum_{m=1}^q \text{corr}_n^\gamma(X_{.l}, Y_{.m}) \quad (4)$$

130 for a set of integers  $\gamma \geq 1$ . Each term  $\text{corr}_n^\gamma(X_{.l}, Y_{.m})$  in equation (4) can be re-written as  
131  $\text{corr}_n^\gamma(X_{.l}, Y_{.m}) = w_{lm} \text{corr}_n(X_{.l}, Y_{.m})$ , where  $w_{lm} = \text{corr}_n^{\gamma-1}(X_{.l}, Y_{.m})$  is regarded as a weight for

1  
2  
3 132  $\text{corr}_n(X_{.l}, Y_{.m})$ . Therefore, a larger  $|\text{corr}_n(X_{.l}, Y_{.m})|$  will yield higher weight  $|w_{lm}|$ , which will help  
4  
5 133 improve power with sparse alternatives that are common for moderate- to high-dimensional data.  
6  
7 134 Specifically, when  $\gamma = 1$ , all  $\text{corr}_n(X_{.l}, Y_{.m})$ 's will be assigned an equal weight 1, which will be  
8  
9 135 beneficial for dense alternatives (i.e. if all or most of the columns of the two matrices  $X$  and  $Y$  are  
10  
11 136 associated) with the same association direction; however, when  $\gamma \geq 2$ , the larger the  $\gamma$ , the higher  
12  
13 137 weights would be assigned to those larger  $\text{corr}_n(X_{.l}, Y_{.m})$ 's, more and more favoring sparse alterna-  
14  
15 138 tives (i.e. when only few of the columns of  $X$  and  $Y$ , as indicated by those larger  $\text{corr}_n(X_{.l}, Y_{.m})$ 's,  
16  
17 139 are truly associated with each other); an even integer  $\gamma$  would give a test robust to varying asso-  
18  
19 140 ciation directions while an odd  $\gamma$  would not. In the extreme case of a sparse alternative with only  
20  
21 141 one or few associated column-pairs between  $X$  and  $Y$ , for an even integer  $\gamma \rightarrow \infty$ , we have

$$22 \quad \text{SPC}(\gamma) \propto \left( \sum_{l=1}^p \sum_{m=1}^q \text{corr}_n^\gamma(X_{.l}, Y_{.m}) \right)^{1/\gamma} \rightarrow \max_j |\text{corr}_n(X_{.l}, Y_{.m})| = \text{SPC}(\infty), \quad (5)$$

23  
24  
25 142 which we can see largely eliminates the effects of non-associated pairs and thus is expected to be  
26  
27 143 more powerful for more sparse alternatives. We emphasize that, with large  $p$  and  $q$  in moderate  
28  
29 144 to high dimensions, noise accumulation is a severe problem for sparse alternatives, which explains  
30  
31 145 power loss of many non-adaptive tests like the RV test, as to be shown later.

32 146 In summary, depending on the type of a true alternative hypothesis to be tested, i.e. dense or  
33  
34 147 sparse, a small or a large  $\gamma$  would yield higher power for the SPU( $\gamma$ ) test. In practice, because it is  
35  
36 148 unknown what is the true alternative and thus which  $\gamma$  value would yield high power, we develop  
37  
38 149 an adaptive SPC (aSPC) test to combine the evidence across the SPC tests:

$$40 \quad \text{aSPC} = \min_{\gamma \in \Gamma} P_{\text{SPC}(\gamma)} \quad (6)$$

41  
42  
43 150 where  $P_{\text{SPC}(\gamma)}$  is the  $P$ -value of the SPC( $\gamma$ ) test, and  $\Gamma$  contains a set of candidate values for  $\gamma$ . In  
44  
45 151 general,  $\Gamma = \{1, 2, \dots, \gamma_u, \infty\}$  with  $1 < \gamma_u < \infty$  can be used; larger  $p$  and  $q$  require a larger  $\gamma_u$ ; a  
46  
47 152 practical guideline on the choice of  $\gamma_u$  is that SPC( $\gamma_u$ ) gives results similar to SPC( $\infty$ ). We used  
48  
49 153  $\Gamma = \{1, \dots, 8, \infty\}$  throughout this paper for its good performance based on our limited experience.

50 154 A permutation method can be used to obtain the  $P$ -values of all the SPC and aSPC tests in  
51  
52 155 a *single loop* (or layer) of permutations. Briefly,  $B$  copies of the null statistic  $\text{SPC}(\gamma)^{(b)}$  for each  
53  
54 156  $\gamma \in \Gamma$  and  $b = 1, \dots, B$  can be calculated by permuting the rows of matrices  $X$  (or  $Y$ )  $B$  times. The  
55  
56 157  $P$ -value of each SPC( $\gamma$ ) is calculated as  $P_{\text{SPC}(\gamma)} = [\sum_{b=1}^B I(|\text{SPC}(\gamma)^{(b)}| \geq |\text{SPC}(\gamma)|) + 1]/(B + 1)$ .



1  
2  
3 Furthermore, based on the same  $B$  copies of the null statistics, we calculate the  $P$ -value for the  
4  
5 aSPU test as  $P_{\text{aSPC}} = [\sum_{b=1}^B I(\text{aSPC}^{(b)} \leq \text{aSPC}) + 1]/(B + 1)$  with  $\text{aSPC}^{(b)} = \min_{\gamma \in \Gamma} p_{\gamma}^{(b)}$  and  
6  
7  $p_{\gamma}^{(b_1)} = [\sum_{b \neq b_1} I(|\text{SPC}(\gamma)^{(b)}| \geq |\text{SPC}(\gamma)^{(b_1)}|) + 1]/B$ .

### 161 2.3 Connections with some existing tests

162 We start by establishing a relationship between the aSPC test (with the Pearson correlation co-  
163 efficient) and an existing test called GEE-aSPU, which was proposed by Kim et al. (2016) for  
164 multiple trait-multiple SNP associations. We first review the GEE-aSPU test before pointing out  
165 its connection to the aSPC test.

166 First we need some notations. Denote  $X_i = (x_{i1}, \dots, x_{ip})$  and  $Y_i = (y_{i1}, \dots, y_{iq})^T$  as  $i$ th row  
167 in matrices  $X$  and transpose of  $i$ th row in  $Y$  for  $i = 1, \dots, n$ , respectively; denote  $X_i = I \otimes X_{i\cdot}$ ,  
168 where  $I$  is a  $q \times q$  identity matrix, and  $\otimes$  represents the Kronecker product.

169 Suppose we treat each column  $Y_m$  for  $m = 1, \dots, q$  in  $Y_{n \times q}$  as a response, each column  $X_l$   
170 for  $l = 1, \dots, p$  in  $X_{n \times p}$  as a covariate or predictor of interest; recall that  $Y_m$  and  $X_l$  has been  
171 standardized to have zero mean and unit variance. We can then test if there is any association  
172 between the columns of  $X$  and those of  $Y$  with a marginal generalized linear model

$$32 \quad g(E(Y_i | X_i)) = X_i \beta, \quad (7)$$

35 where  $g(\cdot)$  is a canonical link function, and  $\beta$  is a  $pq$ -dimensional vector of unknown parameters of  
36 interest. We aim to test the null hypothesis  $H_0 : \beta = 0$ . Denote  $\bar{Y}$  as the mean vector of columns of  
37  $Y$ , which is a zero vector of length  $q$ . With a canonical link function and a working independence  
38 model in GEE (Liang and Zeger, 1986), the generalized score vector for  $\beta$  is  
39  
40

$$43 \quad U = \frac{1}{n-1} \sum_{i=1}^n X_i^T (Y_i - \bar{Y}) = \frac{1}{n-1} \sum_{i=1}^n X_i^T Y_{i\cdot}. \quad (8)$$

47 It is easy to verify  $U = (U_{11}, \dots, U_{p1}, \dots, U_{1q}, \dots, U_{pq})^T$  with  $U_{lm} = X_l^T Y_m / (n-1) = \text{corr}_n(X_l, Y_m)$ .  
48 That is, each element  $U_{lm}$  measures the association between columns  $X_l$  and  $Y_m$ . The GEE-SPU  
49 test statistic is defined by  
50

$$53 \quad \text{SPU}(\gamma_1, \gamma_2) = \sum_{m=1}^q \left[ \left( \sum_{l=1}^p U_{lm}^{\gamma_1} \right)^{\frac{1}{\gamma_1}} \right]^{\gamma_2} = \sum_{m=1}^q \left[ \left( \sum_{l=1}^p \text{corr}_n^{\gamma_1}(X_l, Y_m) \right)^{\frac{1}{\gamma_1}} \right]^{\gamma_2}. \quad (9)$$



180 Denote  $\Gamma_1$  and  $\Gamma_2$  are two sets of positive integers. The GEE-aSPU test statistic is then defined as  
 181 the minimum p-value of SPU( $\gamma_1, \gamma_2$ )' tests for all  $\gamma_1 \in \Gamma_1$  and  $\gamma_2 \in \Gamma_2$ :

$$\text{aSPU} = \min_{\gamma_1, \gamma_2} p_{\gamma_1, \gamma_2} \quad (10)$$

182 Here we observe a close connection between the SPC test and the GEE-SPU test: if  $\gamma_1 = \gamma_2 = \gamma$ ,  
 183 we have SPU( $\gamma, \gamma$ ) = SPC( $\gamma$ ). The difference between the aSPC and aSPU tests is that the latter  
 184 searches for two optimal  $(\gamma_1, \gamma_2)$  in a two-dimensional space (i.e. over  $\Gamma_1 \times \Gamma_2$ ), while aSPC  
 185 searches over only a one-dimensional space (i.e.  $\Gamma$ ); the GEE-aSPU test reduces to aSPC if we  
 186 impose  $\gamma_1 = \gamma_2 = \gamma$ .

187 Due to the currently inefficient implementation of the GEE-aSPU test (in its general regression  
 188 framework) in R package **GEE-aSPU**, it cannot be applied to high-dimensional data: it requires  
 189 a large memory space for its inefficient storage of the design matrix with dimension  $np \times pq$  (or  
 190  $nq \times pq$ ) if  $Y$  (or  $X$ ) is treated as the response. As an example, the GEE-aSPU test will need  
 191 about a 40GB memory space if  $p = q = 300$  and sample size  $n = 200$ , not yet available on many  
 192 computers. In contrast, due to its simplicity, the aSPC test is applicable to high-dimensional data.

193 Finally, we comment on that the SPC(2) test is also closely related to several other tests, further  
 194 illustrating the potential power of the aSPC test. First, since the dCov test and the Hilbert-  
 195 Schmidt independence criterion (HSIC) test are equivalent (Sejdinovic, Sriperumbudur, Gretton  
 196 and Fukumizu, 2013), Hua and Ghosh (2015) called them kernel distance covariance method (KDC);  
 197 they further established the equivalence of KDC and multivariate kernel machine regression (KMR)  
 198 test (Maity et al., 2012) (if the same kernels are used in the two). On the other hand, Kim et al.  
 199 (2016) pointed out that GEE-SPU(2,2) is similar to multivariate KMR with a linear kernel; the  
 200 two are exactly the same if the true correlation matrix is used as the working correlation structure  
 201 in GEE for the former, which in general does not hold (unless the columns of  $Y$  are independent),  
 202 because the working independence model is used in GEE-SPU tests. Now, by the equivalence  
 203 between SPC(2) and GEE-aSPU(2,2) and by the above results, we see the close similarity between  
 204 SPC(2) and other tests. Using the weighting argument motivating the development of other SPC( $\gamma$ )  
 205 tests with  $\gamma > 2$ , we expect that the other tests (i.e. dCov, HSIC and KMR with linear kernels)  
 206 may lose power with sparse association patterns, which will be confirmed in our later simulations.

## 207 2.4 Extensions

208 So far we define the SPC test with the Pearson correlation coefficients between the columns of  
 209 the two matrices. Here we generalize the SPC and thus aSPC tests with several other dependence  
 210 measures and with covariates.

### 211 2.4.1 Fisher's transformation

212 We may take Fisher's z-transformation on the sample Pearson correlation coefficient  $r_{lm} = \text{corr}_n(X_{.l}, Y_{.m})$   
 213 before plugging into equation (4). The reason is to account for heterogeneous variances of the  
 214 sample correlations for an alternative hypothesis; as to be shown next, the variance of a sample  
 215 correlation increases monotonically as the absolute value of the true correlation decreases (un-  
 216 der the normality assumption). Specifically, the sample correlation  $r_{lm} = \text{corr}_n(X_{.l}, Y_{.m})$  is re-  
 217 placed by  $z_{lm} = \frac{1}{2} \ln((1 + r_{lm})/(1 - r_{lm}))$  in equation (4). Under the normality assumption (on  
 218 each pair of the columns of  $X$  and  $Y$ ),  $z_{lm}$  is approximately normally distributed with mean  
 219  $\frac{1}{2} \ln((1 + \rho_{lm})/(1 - \rho_{lm}))$  and a constant variance  $1/(n - 3)$ , where  $\rho_{lm}$  is the population Pearson  
 220 correlation coefficient.

221 Given that  $z_{lm} \sim N(\frac{1}{2} \ln((1 + \rho_{lm})/(1 - \rho_{lm})), 1/(n - 3))$ , it is not hard to find the approximate  
 222 distribution of the sample Pearson correlation coefficient is  $r_{lm} \sim N(\rho_{lm}, (1 - \rho_{lm}^2)^2/(n - 3))$ ; the  
 223 variance  $(1 - \rho_{lm}^2)^2/(n - 3)$  is obtained by the delta method and clearly confirms the monotonic-  
 224 ity mentioned above. In particular, since the variance is largest for no correlations, not taking  
 225 Fisher's transformation or not stabilizing the variance may lead to loss of power, especially for  
 226 high-dimensional data, for which sparse alternatives are expected with many non-associated pairs.

227 Whenever needed, to distinguish using Fisher's z-transformed Pearson correlation coefficients  
 228 from using other dependence measures for the SPC and aSPC tests, we will use SPC.P and aSPC.P  
 229 to refer to the former:

$$\text{SPC.P}(\gamma) = \sum_{l=1}^p \sum_{m=1}^q z_{lm}^\gamma, \quad (11)$$

230 and the aSPC.P test is similarly defined as before.

### 231 2.4.2 The aSPC test with Spearman's correlation

232 More generally, the sample Pearson correlation coefficient term  $r_{lm} = \text{corr}_n(X_{.l}, Y_{.m})$  in equation  
 233 (4) can be replaced by a different dependence measure. For example, we can use Spearman's  
 234 (1904) rank correlation coefficient, which is effective for monotonic relationships, in contrast to only

1  
2  
3 235 linear relationships by Pearson's coefficient. The Spearman correlation coefficient is defined as the  
4  
5 236 Pearson correlation coefficient between the ranked variables. Specifically,  $X_l$  and  $Y_m$  ( $l = 1, \dots, p$   
6  
7 237 and  $m = 1, \dots, q$ ) are converted to the rank score vectors  $\text{rank}(X_l)$  and  $\text{rank}(Y_m)$  (e.g. rank score  
8  
9 238 = 1 for the smallest value in  $X_l$  (or  $Y_m$ ) and rank score =  $n$  for the largest value in  $X_l$  or ( $Y_m$ )).  
10 239 The sample Spearman correlation coefficient is calculated as

$$r_{lm}(\text{Spearman}) = \frac{\text{Cov}_n(\text{rank}(X_l), \text{rank}(Y_m))}{\sqrt{\text{Cov}_n(\text{rank}(X_l), \text{rank}(X_l))\text{Cov}_n(\text{rank}(Y_m), \text{rank}(Y_m))}}, \quad (12)$$

16 240 where  $\text{Cov}_n(u, v)$  is a sample covariance between vectors  $u_{n \times 1}$  and  $v_{n \times 1}$ . Then the SPC statistic  
17  
18 241 with Spearman's rank correlation coefficient is defined as:

$$T_{\text{SPC.Sp}(\gamma)} = \sum_{l=1}^p \sum_{m=1}^q r_{lm}^{\gamma}(\text{Spearman}), \quad (13)$$

24 242 and aSPC.Sp is defined similarly as before.

### 27 243 2.4.3 The aSPC test with the distance correlation

30 244 Another extension is to replace each sample Pearson correlation coefficient in equation (4) by  
31  
32 245 a corresponding distance correlation coefficient (dCor), which is derived based on the distance  
33  
34 246 covariance (dCov) (Szykely et al., 2007) and is consistent in detecting any dependency, not only  
35  
36 247 the linear ones (detectable by Pearson's) or monotonic ones (by Spearman's); for example, in  
37  
38 248 the presence of non-linear (and non-monotonic) dependency, use of dCor is expected to be more  
39  
40 249 powerful, as to be confirmed in our later simulations. We first review the usual dCov test and then  
41  
42 250 modify the SPC test with the distance correlations.

43 251 The standard dCov test utilizes all columns in  $X$  and  $Y$  to calculate the pairwise distance before  
44  
45 252 computing the sample distance covariance:

$$a_{ij} = \|X_i - X_j\|^t, b_{ij} = \|Y_i - Y_j\|^t, \quad (14)$$

49 253 where  $\|\cdot\|$  denotes the Euclidean distance/norm;  $X_i$  and  $Y_i$  denote the  $i$ th row of  $X$  and  $Y$   
50  
51 254 respectively ( $i = 1, \dots, n$ );  $t \in (0, 2]$  and  $t = 1$  corresponds to the Euclidean norm, which was  
52  
53 255 used in our data analysis throughout unless specified otherwise. The pairwise distances are doubly  
54  
55 256 centered:

$$A_{ij} = a_{ij} - \bar{a}_i - \bar{a}_j + \bar{a}_{..}, B_{ij} = b_{ij} - \bar{b}_i - \bar{b}_j + \bar{b}_{..}, \quad (15)$$

where  $\bar{a}_{i\cdot}$ ,  $\bar{a}_{\cdot j}$  and  $\bar{a}_{\cdot\cdot}$  are the  $i$ th row mean, the  $j$ th column mean and the grand mean of matrix  $[a_{ij}]$ ;  $\bar{b}_{i\cdot}$ ,  $\bar{b}_{\cdot j}$  and  $\bar{b}_{\cdot\cdot}$  are similar defined for matrix  $[b_{ij}]$ . Then the squared sample distance covariance of  $X$  and  $Y$  is defined as:

$$\text{dCov}_n^2(X, Y) = \frac{1}{n^2} \sum_{i,j=1}^n A_{ij} B_{ij}. \quad (16)$$

A permutation method can be used to calculate the  $P$ -value. The null statistics  $T_{\text{dCov}}^{(b)} = \frac{1}{n^2} \sum_{i,j=1}^n A_{ij}^{(b)} B_{ij}^{(b)}$  can be calculated based on each permuted sample  $X^{(b)}$  and  $Y^{(b)}$ , where  $X^{(b)}$  (or  $Y^{(b)}$ ) is generated by permuting the rows of  $X$  (or  $Y$ ). The  $P$ -value is calculated as  $P_{\text{dCov}} = (\sum_b^B I(\text{dCov}^{(b)} \geq \text{dCov}) + 1)/(B + 1)$  based on  $B$  permutations.

In the standard dCov test, all columns of  $X$  and  $Y$  are used to calculate the pairwise distances; that is, each variable (or dimension) is treated equally a priori, which may not be a good idea for high-dimensional data for the abundance of sparse alternatives. In contrast, in our SPC test, each column/variable of  $X$  and  $Y$  is treated differently according to the magnitudes of their estimated pairwise associations. Specifically, similar to the standard dCov test, first we define all pairwise distances among the observations based on the  $i$ th and  $j$ th elements of  $X_{\cdot l}$  and  $Y_{\cdot m}$  as

$$a_{ij(l)} = \|X_{il} - X_{jl}\|^t, b_{ij(m)} = \|Y_{im} - Y_{jm}\|^t, \quad (17)$$

which computes the  $n \times n$  distance matrices  $(a_{ij(l)})$  and  $(b_{ij(m)})$  for  $i = 1, \dots, n$ ,  $j = 1, \dots, n$ ,  $l = 1, \dots, p$  and  $m = 1, \dots, q$ . Denote  $\bar{a}_{i\cdot(l)}$ ,  $\bar{a}_{\cdot j(l)}$  and  $\bar{a}_{\cdot\cdot(l)}$  as the  $i$ th row mean, the  $j$ th column mean and the grand mean of  $[a_{ij(l)}]$ ; similarly, denote  $\bar{b}_{i\cdot(m)}$ ,  $\bar{b}_{\cdot j(m)}$  and  $\bar{b}_{\cdot\cdot(m)}$  for  $[b_{ij(m)}]$ . The elements  $a_{ij(l)}$  and  $b_{ij(m)}$  are then doubly centered as:

$$A_{ij(l)} = a_{ij(l)} - \bar{a}_{i\cdot(l)} - \bar{a}_{\cdot j(l)} + \bar{a}_{\cdot\cdot(l)}, B_{ij(m)} = b_{ij(m)} - \bar{b}_{i\cdot(m)} - \bar{b}_{\cdot j(m)} + \bar{b}_{\cdot\cdot(m)}, \quad (18)$$

then the squared sample distance covariance is defined as:

$$\text{dCov}_n^2(X_{\cdot l}, Y_{\cdot m}) = \frac{1}{n^2} \sum_{i,j=1}^n A_{ij(l)} B_{ij(m)}. \quad (19)$$

The sample distance correlation (dCor) between  $X_{\cdot l}$  and  $Y_{\cdot m}$  is then defined as

$$\text{dCor}_n(X_{\cdot l}, Y_{\cdot m}) = \frac{\text{dCov}_n(X_{\cdot l}, Y_{\cdot m})}{\sqrt{\text{dCov}_n(X_{\cdot l}, X_{\cdot l}) \text{dCov}_n(Y_{\cdot m}, Y_{\cdot m})}}. \quad (20)$$

1  
2  
3 276 The SPC.dCor test statistic is defined as:  
4

$$5 \text{ SPC.dCor}(\gamma) = \sum_{l=1}^p \sum_{m=1}^q \text{dCor}_n^{\gamma}(X_{.l}, Y_{.m}) \quad (21)$$

6  
7  
8  
9 277 and the aSPC.dCor is similarly defined as before.  
10

11 278 As to be shown later in simulations, the aSPC.dCor test was much more powerful than the  
12  
13 279 standard distance covariance (dCov) test for sparse alternatives in even only moderate dimensions,  
14  
15 280 presumably because the former's weighting on the pairwise dCor's alleviates the harmful effects of  
16  
17 281 noise accumulations in the latter.  
18

#### 19 282 **2.4.4 The aSPC test with covariates**

20  
21 283 The aSPC test can be applied to situations with covariates. We only need to first regress  $X$  and/or  
22  
23 284  $Y$  on the covariates, then use the residuals to construct the SPC tests. We will illustrate such an  
24  
25 285 application in the example section.  
26

### 27 286 **2.5 Software**

28  
29  
30 287 The asymptotic- and permutation-based RV tests are available as functions `coeffRV()` and `RV.rtest()`  
31  
32 288 in R packages `FactoMineR` and `ade4`, respectively. The permutation-based Mantel test, dCov test  
33  
34 289 and GEE-aSPU test are in functions `mantel()`, `dcov.test()`, `GEEaSPUset()` in R packages `vegan`,  
35  
36 290 `energy` and `GEEaSPU`, respectively. We implemented various versions of the new SPC and aSPC  
37  
38 291 tests in an R package `aSPC`, which is available on github (and CRAN).  
39

## 40 292 **3 Simulations**

### 41 293 **3.1 Simulation I: linear associations**

42  
43  
44 294 To further investigate the operating characteristics of the proposed tests, we compare their power  
45  
46 295 performance with several existing tests. We first consider an ideal situation with a linear association  
47  
48 296 between two sets of normal variates.  
49

50  
51 297 To generate a simulated dataset, two matrices  $X_{n \times p}$  and  $Y_{n \times p}$  were simulated with  $n = 500$ .  
52  
53 298 First, for each  $X$  and  $Y$ ,  $p$  ( $= 25, 45$  or  $65$ ) independent columns were simulated from a standard  
54  
55 299 multivariate normal distribution. Second, a matrix  $Z_{n \times 10}$  with ten columns were simulated from a  
56  
57 300 multivariate normal distribution with mean 0 and a compound symmetry covariance matrix (with  
58  
59  
60

all diagonal elements equal to 1 and all off-diagonal elements equal to 0.1); for power comparisons, we added the first 5 columns of  $Z$  to  $X$  and the last 5 columns of  $Z$  to  $Y$ .

We applied the aSPC.P, aSPC.Sp, aSPC.dCor, RV, Mantel and dCov tests to each simulated dataset, and compared their empirical Type I error and power estimates. The Mantel and dCov tests were conducted with the Euclidean distance. We set  $B = 1000$  for any permutation-based tests. To save computing time, the empirical Type I error rates and power of aSPC.dCor were based on 1,000 replicates while for all other tests, they were based on 10,000 replicates.

As shown in Table 1, first, the Type I error rates were in general well controlled for each test. Second, among all the tests, GEE-aSPU was most powerful, followed by aSPC.P. Note that, due to the linear association, aSPC.P is expected to be more powerful than aSPC.Sp (and aSPC.dCor). Third, SPC.P(2) gave the results essentially the same as both the asymptotic and permutation-based RV tests, as expected. Fourth, due to the presence many independent columns in the two matrices  $X$  and  $Y$ , a SPC.P test with a larger and finite  $\gamma$  (e.g.  $\gamma = 6$ ) was more powerful than that with a small  $\gamma \leq 4$ ; their power difference increased with the number of independent columns. Fifth, aSPC.dCor gave much higher power than dCov test, due to that SPC.dCor( $\gamma$ ) with larger  $\gamma$  reduced the effects of noise accumulation with independent columns. Moreover, we note the extremely low power of the Mantel test, followed by MANOVA.

To assess the computing time and feasibility for the permutation-based RV, GEE-aSPU and aSPC tests, we changed the number of columns in  $X$  and  $Y$  to 30, 50, 70 and 100 respectively, and with a sample size  $n = 200$ . We then calculated the computing time with a permutation number  $B = 1 \times 10^3$ . Note that, for example, for  $p = q = 300$ , GEE-aSPU needs to construct a large design matrix with dimension  $60,000 \times 90,000$ , requiring about 40GB of memory. The computing time was based on one processor (Intel Haswell E5-2680v3 with 2.5GB of memory on Unix system) from a cluster at the Minnesota Supercomputing Institute (MSI).

As shown in Figure 1, first, our implementation of aSPC.P completely in R was even faster than the RV.perm test, which was surprising given that aSPC.P involved conducting SPC.P( $\gamma$ ) for  $\gamma = 1, \dots, 8, \infty$  and RV.perm is equivalent to SPC.P(2). Second, aSPC.dCor was more computing-intensive than other tests; for data matrices  $X_{n \times p}$  and  $Y_{n \times q}$ , aSPC.dCor required calculating pairwise distance covariances  $pq$  times based on  $p + q$  distance matrices, even if we used more memory space to save the distance matrices in our current implementation in R.

### 3.2 Simulation II: non-linear associations

Now we consider a more challenging case with a non-linear and non-monotonic association. Our simulation set-up was similar to that of Székely et al. (2007).

Data matrix  $X_{n \times 5}$  was simulated from a multivariate standard normal distribution. To calculate the empirical type I error rates, for each replicate a matrix  $Y_{n \times 5}$  was simulated from a multivariate standard normal distribution. For power,  $Y_{n \times p}$  was generated such that each of the first  $p_0$  ( $p_0 = 1, 2, 3, 4$  or  $5$ ) columns  $Y_{ij} = \log(X_{ij}^2)$  for  $j = 1, \dots, p_0$  and  $i = 1 \dots n$ ; when  $p_0 \leq 4$ , each of the other columns of  $Y_{n \times p}$  was independently and identically simulated from a standard normal distribution. We were interested in how the empirical power changed as the number of non-linearly associated column pairs ( $p_0$ ) between  $X$  and  $Y$  varied from 1 to 5. Six tests were applied, including aSPC.dCor, aSPC.Sp, aSPC.P, permutation-based RV test, the Mantel test with the Euclidean distance and Pearson correlation, and dCov. One thousand datasets were simulated to calculate the empirical type I error and power. We used  $B = 1000$  for any permutation-based tests. The simulation results are summarized in the left panel of Figure 2 with sample size  $n = 40$ .

First, the type I error rates were well controlled for all tests. Second, our aSPC.dCor test gave much higher power than the usual dCov test. For example, with only one truly associated pair, the power of aSPC.dCor was 86.5%, much higher than 12.0% of the dCov test. Third, due to the underlying non-monotonic true associations, as expected, none of the RV, aSPC.P, aSPC.Sp and Mantel tests performed well.

To further explore the performance of the tests with increasingly sparse associations, in addition to the above set-up with  $p_0 = 5$ , we added 75, 115, 195, 295 or 395 independent columns to matrix  $Y$ , each of which was simulated from a standard normal distribution. The power curves are shown in the right panel of Figure 2. It is clear that the power of aSPC.dCor remained significantly higher than that of the dCov test, whereas all other tests had no power.

## 4 Real data application

### 4.1 ADNI data

Data used in the preparation of this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database ([adni.loni.usc.edu](http://adni.loni.usc.edu)). The ADNI was launched in 2003 by the National Institute on Aging (NIA), the National Institute of Biomedical Imaging and Bioengi-



neering (NIBIB), the Food and Drug Administration (FDA), private pharmaceutical companies and non-profit organizations, as a \$60 million, 5-year public-private partnership. The primary goal of ADNI has been to test whether serial magnetic resonance imaging (MRI), positron emission tomography (PET), other biological markers, and clinical and neuropsychological assessment can be combined to measure the progression of mild cognitive impairment (MCI) and early Alzheimer's disease (AD). Determination of sensitive and specific markers of very early AD progression is intended to aid researchers and clinicians to develop new treatments and monitor their effectiveness, as well as lessen the time and cost of clinical trials.

The Principal Investigator of this initiative is Michael W. Weiner, MD, VA Medical Center and University of California-San Francisco. ADNI is the result of efforts of many co-investigators from a broad range of academic institutions and private corporations, and subjects have been recruited from over 50 sites across the U.S. and Canada. The initial goal of ADNI was to recruit 800 subjects but ADNI has been followed by ADNI-GO and ADNI-2. To date these three protocols have recruited over 1500 adults, ages 55 to 90, to participate in the research, consisting of cognitively normal older individuals, people with early or late MCI, and people with early AD. The follow up duration of each group is specified in the protocols for ADNI-1, ADNI-2 and ADNI-GO. Subjects originally recruited for ADNI-1 and ADNI-GO had the option to be followed in ADNI-2. For up-to-date information, see [www.adni-info.org](http://www.adni-info.org).

## 4.2 Testing for SNP-gene expression associations

To understand gene regulation, it is important to detect genetic variants like single nucleotide polymorphisms (SNPs) that are associated with gene expression (i.e. transcript) levels, called eQTL (Minas, Curry and Montana, 2013). Due to the relatively small sample size and a severe penalty on multiple testing for a large number of SNP-gene pairs, it is often low-powered to detect many associations at the individual pair level. As an alternative, we may first test the association between a set of SNPs and a set of the genes.

The ADNI genotype data consist of 757 subjects from ADNI-1, two hundred and thirty six of whom also have genome-wide gene expression data based on the whole blood. A pathway for Alzheimer's disease (hsa05010, [http://www.genome.jp/dbget-bin/www\\_bget?hsa05010](http://www.genome.jp/dbget-bin/www_bget?hsa05010)) was downloaded from Kyoto Encyclopedia of Genes and Genomes (KEGG) website (Kanehisa, Sato, Kawashima, Furumichi and Tanabe, 2016). Since the ADNI-1 genotype data are based on the human genome version hg18, we used the hg18 gene coordinate file downloaded from the PLINK

1  
2  
3  
4 391 website (<http://pngu.mgh.harvard.edu/~purcell/plink/>) to identify the starting base pair (bp)  
5 392 and ending bp for each gene. We then extracted two sets of the SNPs for the genes in the AD  
6  
7 393 pathway. In the first, the SNPs within each gene were selected, including possibly both protein  
8  
9 394 coding and regulatory SNPs; in the second, to focus on only regulatory SNPs, only the SNPs within  
10  
11 395 the upstream 20kb of a gene's starting bp or within the downstream 20kb of its ending bp were  
12  
13 396 selected. Since the results were similar, we will discuss only the first dataset.

14 397 To account for possible effects of age and gender on gene expression, we used a linear regression  
15  
16 398 model to regress each gene's expression level on the two covariates, then used the residuals as the  
17  
18 399 gene's adjusted expression levels in the subsequent analysis. In the end, there were 441 probes  
19  
20 400 corresponding to 151 genes, and 2,483 SNPs (after excluding those with a minor allele frequency  
21  
22 401 less than 0.05) in the first dataset.

23 402 To demonstrate the effects of association patterns, especially the signal sparsity levels, on the  
24  
25 403 testing results, we screened the SNP-gene pairs using each pair's P-value for their marginal associa-  
26  
27 404 tion, which was based on a simple linear regression of each gene's adjusted expression level on each  
28  
29 405 SNP in the set. The expression level of each gene was calculated as the average of its corresponding  
30  
31 406 probes for those genes with more than one probe. We used various threshold values to select subsets  
32  
33 407 of the SNP-gene pairs, with a marginal P-value smaller than a given threshold. Then we pooled  
34  
35 408 the SNPs and the probes in the genes surviving such a screening into a SNP set and a probe set  
36  
37 409 respectively, then tested their associations using various methods. For any permutation-based test,  
38  
39 410 we used a permutation number  $B = 1 \times 10^4$  (unless specified otherwise). As the dimensions of the  
40  
41 411 probes and the SNPs were high (i.e. in hundreds to thousands), it would be infeasible to run the  
42  
43 412 GEE-aSPU test as it required a too large memory space. The results are summarized in Table 2.

44 413 We have the following observations. First, when we included all the SNPs and the probes (with  
45  
46 414 a P-value threshold 1), the aSPC tests (i.e. aSPC.P, aSPC.Sp, and aSPC.dCor) all gave significant  
47  
48 415 P-values; in contrast, none of the other tests, including the RV test, the Mantel test and dCov test,  
49  
50 416 gave any significant P-value less than the nominal level 0.05. Second, most strikingly, regardless  
51  
52 417 of the dimensions  $(p, q)$  with various threshold values, the aSPC tests consistently gave small and  
53  
54 418 significant P-values (e.g.  $< 0.001$ ), showing their robustness to the varying association patterns  
55  
56 419 (e.g. signal sparsity levels); in contrast, as fewer and fewer, but more significant, SNPs and probes  
57  
58 420 were included, other global tests gradually gave more and more significant P-values, suggesting  
59  
60 421 their loss of power in the presence of sparse signals due to their none-adaptiveness. Third, among  
61  
62 422 the SPC tests, those SPC.P( $\gamma$ ) tests with larger  $\gamma$  (e.g.  $\gamma \geq 4$ ) gave more significant P-values

1  
2  
3 423 than those with smaller  $\gamma$  (e.g.  $\gamma < 4$ ), indicating sparse signals as expected (i.e. most SNP-probe  
4 424 pairs were not associated).

## 8 425 **5 Discussion**

10 426 We have proposed an adaptive and powerful association test called aSPC for two moderate- to high-  
11 427 dimensional random vectors. It has been shown to be more powerful in a variety of simulations  
12 428 than several commonly used tests. In an application to a real genotype-gene expression dataset,  
13 429 under various moderately high dimensions for the SNPs and genes, the proposed test robustly and  
14 430 consistently gave more significant P-values than other existing tests, which appeared to lose power  
15 431 dramatically for larger sets of the SNPs and genes. The proposed aSPC test can be regarded as a  
16 432 generalization of the standard RV test from low-dimensional data to moderate- to high-dimensional  
17 433 data with the incorporation of data-adaptive weighting on each variable pair. The main idea is  
18 434 that, for moderate- to high-dimensional data, often there will be many variable pairs that are not  
19 435 associated; treating these null pairs equally as other truly associated pairs will simply accumulate  
20 436 noises, leading to substantial power loss as in most other existing tests like the RV test. Hence,  
21 437 this main idea is related to the GEE-aSPU test in genetics. Indeed the aSPC test (more precisely,  
22 438 the version denoted aSPC.P with Pearson's correlation) is a special case of the GEE-aSPU test.  
23 439 However, due to its simplicity, the aSPC.P test has some computational advantage over the GEE-  
24 440 aSPU test, which in its currently implementation is not applicable to high-dimensional data. More  
25 441 importantly, the aSPC.P test can be easily extended by replacing the Pearson correlation coefficient  
26 442 with other coefficient, which may be more suitable for other non-linear associations. For example,  
27 443 if the distance correlation is used as in aSPC.dCor, it can detect non-monotonic associations.  
28 444 Compared to the usual dCov (or dCor) test, again due to its adaptiveness, the aSPC.dCor test  
29 445 is much more powerful for less dense or sparse signals for high-dimensional data, as shown in our  
30 446 simulations.

31 447 In the current implementation of the new tests, we have resorted to permutations to calculate  
32 448 their P-values, which seems feasible and satisfactory in many applications. However, it would be  
33 449 interesting to establish their asymptotics as both  $p$  and  $q$  diverge with  $n$  (Xu et al., 2016), which  
34 450 may be challenging due to the dependencies among the individual correlation coefficients in each  
35 451 SPC test statistic. Nevertheless, an asymptotic theory will be useful in facilitating speedy P-value  
36 452 calculations, especially for a high significance level.

1  
2  
3 The various versions of the aSPC test are implemented in R package aSPC, freely available on  
4  
5 453 CRAN or at <https://github.com/jasonzyx/aSPC>.  
6  
7

## 8 455 **Acknowledgment**

9  
10  
11 456 The authors are grateful to the reviewers for constructive comments, and thank Dr. Baolin Wu for  
12  
13 457 his question that motivated the development in section 2.3. This research was supported by NIH  
14  
15 458 grants R01GM113250, R01HL105397 and R01HL116720, and by the Minnesota Supercomputing  
16  
17 459 Institute at the University of Minnesota. ZX was supported by a University of Minnesota MnDRIVE  
18  
19 460 Fellowship.

20 461 Data collection and sharing for this project was funded by the Alzheimer's Disease Neuroimaging  
21  
22 462 Initiative (ADNI) (National Institutes of Health Grant U01 AG024904) and DOD ADNI (Depart-  
23  
24 463 ment of Defense award number W81XWH-12-2-0012). ADNI is funded by the National Institute  
25  
26 464 on Aging, the National Institute of Biomedical Imaging and Bioengineering, and through generous  
27  
28 465 contributions from the following: Alzheimer's Association; Alzheimers Drug Discovery Foundation;  
29  
30 466 Araclon Biotech; BioClinica, Inc.; Biogen Idec Inc.; Bristol-Myers Squibb Company; Eisai Inc.;  
31  
32 467 Elan Pharmaceuticals, Inc.; Eli Lilly and Company; EuroImmun; F. Hoffmann-La Roche Ltd and  
33  
34 468 its affiliated company Genentech, Inc.; Fujirebio; GE Healthcare; IXICO Ltd.; Janssen Alzheimer  
35  
36 469 Immunotherapy Research & Development, LLC.; Johnson & Johnson Pharmaceutical Research &  
37  
38 470 Development LLC.; Medpace, Inc.; Merck & Co., Inc.; Meso Scale Diagnostics, LLC.; NeuroRx  
39  
40 471 Research; Neurotrack Technologies; Novartis Pharmaceuticals Corporation; Pfizer Inc.; Piramal  
41  
42 472 Imaging; Servier; Synarc Inc.; and Takeda Pharmaceutical Company. The Canadian Institutes of  
43  
44 473 Rev December 5, 2013 Health Research is providing funds to support ADNI clinical sites in Canada.  
45  
46 474 Private sector contributions are facilitated by the Foundation for the National Institutes of Health  
47  
48 475 (www.fnih.org). The grantee organization is the Northern California Institute for Research and  
49  
50 476 Education, and the study is coordinated by the Alzheimer's Disease Cooperative Study at the Uni-  
51  
52 477 versity of California, San Diego. ADNI data are disseminated by the Laboratory for Neuro Imaging  
53  
54 478 at the University of Southern California.

## 51 479 **References**

52  
53  
54 480 Colantuoni, C., Lipska B. K., Ye, T., Hyde, T. M., Tao, R., Leek, J. T., ... Kleinman J. E.  
55  
56 481 (2011). Temporal dynamics and genetic control of transcription in the human prefrontal  
57  
58

1  
2  
3  
4 482 cortex. *Nature* 478(7370), 519-523.

5  
6 483 Escoufier, Y. (1970). Echantillonnage dans une population de variables aléatoires réelles. Depart-  
7 484 ment de math.; Univ. des sciences et techniques du Languedoc.

8  
9  
10 485 Escoufier, Y. (1973). Le traitement des variables vectorielles. *Biometrics* 29(4), 751-760.

11  
12 486 Fan, R., Chiu, C., Jung J., Weeks, D. E., Wilson, A. F., Bailey-Wilson, J. E., ... Xiong, M. (2016).  
13 487 A comparison study of fixed and mixed effect models for gene level association studies of  
14 488 complex traits. *Genetic Epidemiology* 40(8), 702-721.

15  
16  
17  
18 489 He Z., Zhang M., Lee S., Smith J. A., Guo X., Palmas W., ... Mukherjee B. (2015). Set-based  
19 490 tests for genetic association in longitudinal studies. *Biometrics* 71(3), 606-615.

20  
21  
22 491 Hua, W. Y., & Ghosh, D. (2015). Equivalence of kernel machine regression and kernel distance  
23 492 covariance for multidimensional phenotype association studies. *Biometrics* 71(3), 812-820.

24  
25  
26 493 Josse, J., Holmes, S. (2014). Measures of dependence between random vectors and tests of  
27 494 independence. Literature review. Ithaca, NY: Cornell University Library. Available at  
28 495 <http://arxiv.org/pdf/1307.7383v3.pdf> (accessed November 22th, 2014).

29  
30  
31  
32 496 Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M., & Tanabe, M. (2016). KEGG as a  
33 497 reference resource for gene and protein annotation. *Nucleic Acids Research* 44(D1), D457-  
34 498 D462 .

35  
36  
37  
38 499 Kim, J., Zhang, Y., & Pan, W. (2016). Powerful and adaptive testing for multi-trait and multi-  
39 500 SNP associations with GWAS and sequencing data. *Genetics* 203(2), 715-731.

40  
41  
42 501 Liang, K. Y., & Zeger, S. L. (1986). Longitudinal data analysis using generalized linear models.  
43 502 *Biometrika* 73(1), 13-22.

44  
45  
46 503 Maity, A., Sullivan, P. F., & Tzeng, J. Y. (2012). Multivariate phenotype association analysis by  
47 504 marker-set kernel machine regression. *Genetic Epidemiology* 36(7), 686-695.

48  
49  
50 505 Mantel, N. (1967). The detection of disease clustering and a generalized regression approach.  
51 506 *Cancer Research* 27(2 Part 1), 209-220.

52  
53  
54 507 Minas, C., Curry, E., & Montana, G. (2013). A distance-based test of association between paired  
55 508 heterogeneous genomic data. *Bioinformatics* 29(20), 2555-2563.

- 1  
2  
3 509 Sejdinovic, D., Sriperumbudur, B., Gretton, A., & Fukumizu, K. (2013). Equivalence of distance-  
4 based and RKHS-based statistics in hypothesis testing. *Annals of Statistics* 41(5), 2263-2291.  
5 510  
6  
7 511 Spearman, C. E. (1904a). The proof and measurement of association between two things. *Amer-*  
8 *ican Journal of Psychology* 15(1), 72-101.  
9 512  
10  
11 513 Székely, G. J., Rizzo, M. L., & Bakirov, N. K. (2007). Measuring and testing dependence by  
12 correlation of distances. *Annals of Statistics* 35(6), 2769-2794.  
13 514  
14  
15 515 Wang, X., Lee, S., Zhu, X., Redline, S., & Lin, X. (2013). GEE-based SNP set association test  
16 for continuous and discrete traits in family-based association studies. *Genetic Epidemiology*  
17 516  
18 37(8), 778-786.  
19 517  
20  
21 518 Wang, Y., Liu, A., Mills, J. L., Boehnke, M., Wilson, A. F., Bailey-Wilson, J. E., ... Fan, R.  
22 (2015). Pleiotropy analysis of quantitative traits at gene level by multivariate functional  
23 519  
24 linear models. *Genetic Epidemiology* 39(4), 259-275.  
25 520  
26  
27 521 Wang Z., Xu K., Zhang X., Wu X., & Wang Z. (2017). Longitudinal SNP-set association analysis  
28 of quantitative phenotypes. *Genetic Epidemiology* 41(1), 81-93.  
29 522  
30  
31 523 Xu, G., Lin, L., Wei, P., & Pan, W. (2016). An adaptive two-sample test for high-dimensional  
32 means. *Biometrika* 103(3), 609-624.  
33 524  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

525 **Tables**

For Peer Review

Author Manuscript



Table 1: Simulation I: empirical Type I error and power rates when the number of independent columns (denoted as “No. ind”) is 25, 45, and 65 respectively. “RV.asy” and “RV.perm” stand for the asymptotic and permutation-based RV tests, respectively

No.Ind		SPC,P( $\gamma$ )								Inf	aSPC.P	aSPC.Sp	aSPC.dCor	RV.asy	RV.perm	Mantel	dCov	GEE-aSPU	MANOVA
		$\gamma = 1$	2	3	4	5	6	7	8										
25	Type I	0.047	0.053	0.049	0.050	0.050	0.051	0.052	0.053	0.054	0.046	0.049	0.049	0.053	0.055	0.050	0.052	0.055	0.046
	Power	0.417	0.844	0.886	0.932	0.917	0.908	0.879	0.852	0.589	0.933	0.893	0.828	0.840	0.838	0.098	0.819	0.955	0.378
45	Type I	0.055	0.052	0.052	0.052	0.052	0.053	0.050	0.050	0.048	0.052	0.049	0.050	0.052	0.051	0.050	0.053	0.061	0.045
	Power	0.196	0.538	0.587	0.753	0.732	0.759	0.710	0.700	0.425	0.749	0.674	0.587	0.539	0.538	0.074	0.522	0.832	0.174
65	Type I	0.056	0.055	0.050	0.052	0.054	0.050	0.050	0.051	0.049	0.050	0.050	0.047	0.056	0.055	0.052	0.054	0.057	0.041
	Power	0.118	0.352	0.371	0.581	0.558	0.627	0.576	0.578	0.328	0.594	0.506	0.450	0.355	0.354	0.072	0.345	0.702	0.110

Table 2: The analysis results for the ADNI data.  $p$  and  $q$  denote the numbers of SNPs and of probes surviving the P-value cut-off based on the corresponding univariate SNP-gene expression associations

Cut-off	(p, q)	SPC.P( $\gamma$ )					aSPC.P	RV.asy	RV.perm	Mantel	dCov	aSPC.Sp	aSPC.dCor
		$\gamma = 1$	2	3	4	5-8, $\infty$							
1	(2483, 382)	5.85e-02	3.68e-02	2.87e-01	3.00e-04	1.00e-04	8.00e-04	6.89e-02	7.17e-02	8.69e-02	5.61e-02	9.00e-04	5.00e-04
0.9	(2274, 380)	3.63e-02	4.50e-02	2.16e-01	5.00e-04	1.00e-04	8.00e-04	6.47e-02	6.37e-02	9.76e-02	4.99e-02	8.00e-04	5.00e-04
0.8	(2069, 371)	2.29e-02	2.09e-02	2.22e-01	1.00e-04	1.00e-04	8.00e-04	3.87e-02	3.86e-02	7.40e-02	2.62e-02	8.00e-04	5.00e-04
0.7	(1871, 357)	3.26e-02	8.60e-03	2.39e-01	1.00e-04	1.00e-04	7.00e-04	2.01e-02	2.05e-02	5.81e-02	1.27e-02	8.00e-04	5.00e-04
0.6	(1647, 353)	1.39e-02	4.40e-03	1.74e-01	1.00e-04	1.00e-04	9.00e-04	9.22e-03	8.90e-03	4.71e-02	6.50e-03	7.00e-04	6.00e-04
0.5	(1435, 351)	1.62e-02	2.90e-03	2.80e-01	1.00e-04	1.00e-04	6.00e-04	6.69e-03	7.70e-03	5.96e-02	3.10e-03	6.00e-04	6.00e-04
0.4	(1228, 340)	1.99e-02	8.00e-04	2.54e-01	1.00e-04	1.00e-04	9.00e-04	1.91e-03	2.30e-03	1.49e-02	1.60e-03	9.00e-04	4.00e-04
0.3	(999, 306)	5.95e-02	1.20e-03	4.62e-01	1.00e-04	1.00e-04	8.00e-04	1.48e-03	1.50e-03	7.30e-03	9.00e-04	8.00e-04	1.00e-04
0.2	(756, 286)	7.54e-02	6.00e-04	5.93e-01	1.00e-04	1.00e-04	7.00e-04	6.07e-04	2.00e-04	2.00e-03	4.00e-04	9.00e-04	4.00e-04
0.1	(485, 245)	2.93e-01	1.00e-04	3.34e-01	1.00e-04	1.00e-04	7.00e-04	8.29e-05	3.00e-04	4.00e-04	2.00e-04	8.00e-04	4.00e-04

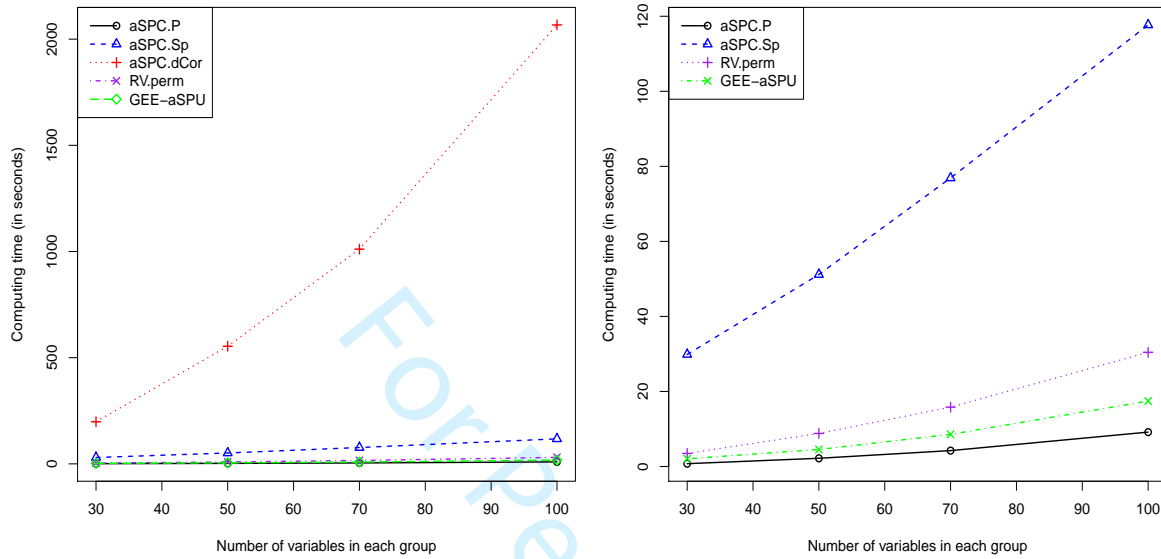
526 **Figures**

Figure 1: The computing time of the permutation-based RV, GEE-aSPU, aSPC.P, aSPC.Sp and aSPC.dCor tests. The left panel shows the computing time of aSPC.dCor test as compared to that of all the other tests, while the right panel is a zoom-in for all the tests except aSPC.dCor

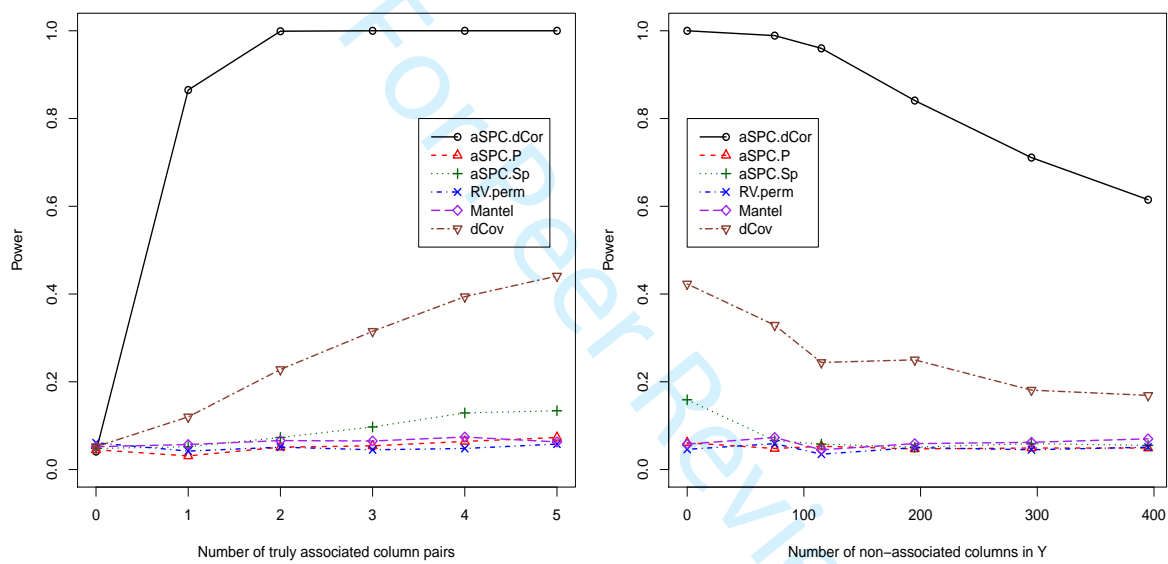


Figure 2: Simulation II results. The left panel: when the number of columns in  $X$  and  $Y$  are 5, the empirical type I error and power curves of the tests as the number of truly non-linearly associated column pairs between  $X$  and  $Y$  ranges from 0 (type I error) to 5. Right panel: when the number of non-linearly associated column pairs in  $X$  and  $Y$  is fixed at 5, the power curves of the tests as more and more non-associated columns are added to  $Y$ . The nominal significance level is 0.05