

# Analysis of variance when both input and output sets are high-dimensional

Gustavo de los Campos<sup>1,2,3\*</sup>, Torsten Pook<sup>4\*</sup>, Agustin Gonzalez-Raymunde<sup>5</sup>,  
Henner Simianer<sup>4</sup>, George Mias<sup>6,3</sup> & Ana I. Vazquez<sup>1,3</sup>.

1: Epidemiology & Biostatistics, Michigan State University, 909 Wilson Rd Room B601, East Lansing, MI 48824, US.

2: Statistics & Probability, Michigan State University, 619 Red Cedar Rd room c413, East Lansing, MI 48824, US.

3: Institute for Quantitative Health Science and Engineering, 775, Woodlot Dr, East Lansing, MI, 48824, US.

4: Department of Animal Sciences, Center for Integrated Breeding Research, University of Goettingen, Goettingen, 37075, Germany.

5: Genetics and Genome Sciences graduate program, Michigan State University, 775, Woodlot Dr, East Lansing, MI, 48824, US.

6: Biochemistry and Molecular Biology, Michigan State University, 603 Wilson Rd Rm 212, East Lansing, MI 48823, US.

\* To whom correspondence should be addressed.

## Abstract

**Motivation:** Modern genomic data sets often involve multiple data-layers (e.g., DNA-sequence, gene expression), each of which itself can be high-dimensional. The biological processes underlying these data-layers can lead to intricate multivariate association patterns.

**Results:** We propose and evaluate two methods for analysis variance when both input and output sets are high-dimensional. Our approach uses random effects models to estimate the proportion of variance of vectors in the linear span of the output set that can be explained by regression on the input set. We consider a method based on orthogonal basis (Eigen-ANOVA) and one that uses random vectors (Monte Carlo ANOVA, MC-ANOVA) in the linear span of the output set. We used simulations to assess the bias and variance of each of the methods, and to compare it with that of the Partial Least Squares (PLS)—an approach commonly used in multivariate-high-dimensional regressions. The MC-ANOVA method gave nearly unbiased estimates in all the simulation scenarios considered. Estimates produced by Eigen-ANOVA and PLS had noticeable biases. Finally, we demonstrate insight that can be obtained with the of MC-ANOVA and Eigen-ANOVA by applying these two methods to the study of multi-locus linkage disequilibrium in chicken genomes and to the assessment of inter-dependencies between gene expression, methylation and copy-number-variants in data from breast cancer tumors.

39 **Availability:** The Supplementary data includes an R-implementation of each of the proposed  
40 methods as well as the scripts used in simulations and in the real-data analyses.

41 **Contact:** [gustavoc@msu.edu](mailto:gustavoc@msu.edu)

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

43

44 **Keywords:** multi-omic data, high-dimensional data, ANOVA, singular value decomposition,  
45 Monte Carlo Methods, REML, Partial-Least Squares.

46

47

## 48 **1 Introduction**

49

50 Modern genomic data sets often combine information from multiple data-layers, each of which  
51 itself can be high-dimensional. Examples of this include data sets comprising of information  
52 from several omics, or data sets combining genomic information with high-throughput  
53 phenotyping (e.g., crop-imaging imaging, milk infrared spectra data). The biological processes  
54 underlying each of the data-layers can induce complex dependencies between features within  
55 (e.g., linkage disequilibrium among single nucleotide polymorphisms, SNPs) as well as  
56 between layers (e.g., association between DNA and gene expression, GE). The main goal of  
57 this study is to develop and to evaluate methods for quantifying multivariate-associations in  
58 settings in which both the input and output sets are high dimensional.

59 The methods proposed in this study can be used to answer ubiquitous questions such as:  
60 How much of the inter-individual differences in whole-genome sequence genotypes can be  
61 predicted using a low-density SNP array? What proportion of variance in GE can be explained  
62 by differences in DNA methylation (ME)? How much of the variance in image-phenotypes can  
63 be predicted from DNA genotypes?

64 Canonical Correlation Analysis (CCA, Mardia, T., and Bibby, 1979), Multivariate-  
65 Analysis of Variance (MANOVA, Rencher and Christensen, 2012), and Reduced Rank-  
66 Regressions, e.g., Partial Least Squares (PLS, Wold, Sjöström, and Eriksson, 2001) are three  
67 methodologies often used to assess associations in multi-dimensional problems. However,  
68 these approaches have limitations that make some of them inadequate for estimating the  
69 proportion of variance explained when both the output and input layers are high-dimensional.

70 **Canonical Correlation Analysis** extends the concept of the correlation between two  
71 random variables to cases involving two multidimensional data sets; however, correlation is  
72 symmetric by nature. Therefore, CCA cannot address questions regarding proportion of

73 variance explained when the proportion of variance of one set (e.g.,  $\mathbf{X}$ ) that is explained by  
74 another set ( $\mathbf{W}$ ) is not equal to the reciprocal (i.e., proportion of variance of  $\mathbf{W}$  that can be  
75 explained by  $\mathbf{X}$ ). In many multi-layered data sets we do not expect to have a symmetric  
76 variance-decomposition (we illustrate this below using simulations and experimental data).

77 **Multivariate Analyses of Variance** (MANOVA, e.g., Krzanowski and J. 1988) is another  
78 approach that can be considered. However, MANOVA is based on least-squares projections;  
79 therefore, the methodology is not well-suited for cases when data is high dimensional,  
80 including rank-deficient cases. Most of the problems we are focusing on involve high-  
81 dimensional data where the number of features exceeds sample size, thus, making least-squares  
82 methods such as MANOVA inadequate.

83 **Reduced-rank regressions** (e.g., Izenman, 1975) and **penalized multivariate analysis**  
84 **methods** (e.g., Witten, Tibshirani, and Hastie 2009) are often used to analyze high-  
85 dimensional data. However, the results that one can obtain using regularized methods  
86 (including reduced-rank and penalized methods) rely on regularization decisions (e.g., the  
87 number of dimensions used in PLS or CCA, or the sparsity parameters in sparse CCA) which  
88 cannot be made using the likelihood. Thus, these parameters are often tuned to maximize  
89 prediction accuracy in testing sets. This approach is not necessarily optimal for inferences.

90 Therefore, to overcome the limitations of existing methods, we developed two approaches  
91 for estimation of proportion of variance explained when both the input and output sets are high-  
92 dimensional. Our approaches use random effects models to estimate the proportion of variance  
93 of (independent) vectors in the linear span of an output layer that can be explained by regression  
94 on an input layer. We considered two methods for generating a sequence of independent  
95 vectors in the linear span of the output layer: A Monte Carlo method (MC-ANOVA) which  
96 uses random vectors, and one based on eigenvectors (Eigen-ANOVA).

97

## 98 **2 Methods**

99

100 Let  $\mathbf{X}$  and  $\mathbf{W}$  denote two matrices holding data for  $n$  individuals (rows) and  $p$  ( $\mathbf{X}$ ) and  $q$  ( $\mathbf{W}$ )  
101 features in columns. For instance,  $\mathbf{X}$  may be a matrix with genotypes codes at  $p$  SNPs and  $\mathbf{W}$   
102 may be a matrix providing GE levels assessed at  $q$  genes.

103 The columns of  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p\}$  and of  $\mathbf{W} = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_q\}$  can be viewed as axes  
104 spanning two linear spaces ( $L_X$  and  $L_W$ , respectively). The linear spans of  $\mathbf{X}$  and  $\mathbf{W}$  consist of  
105 all the vectors that can be obtained by forming linear combinations of the columns of each of

106 these sets, that is  $L_X = \{\mathbf{x}_s: \mathbf{x}_s = \mathbf{X}\boldsymbol{\alpha}_s = \sum_{j=1}^p \mathbf{x}_j \alpha_{sj}\}$  and  $L_W = \{\mathbf{w}_s: \mathbf{w}_s = \mathbf{W}\boldsymbol{\delta}_s =$   
107  $\sum_{j=1}^q \mathbf{w}_j \delta_{sj}\}$ , for all real-valued vectors  $\boldsymbol{\alpha}_s = \{\alpha_{s1}, \dots, \alpha_{sp}\}$  and  $\boldsymbol{\delta}_s = \{\delta_{s1}, \dots, \delta_{sq}\}$ .

108 For each vector  $\mathbf{x}_s \in L_X$ , one can estimate the proportion of variance that can be explained  
109 by linear regression on  $L_W$  using a model of the form

$$110 \quad \mathbf{x}_s = \mathbf{W}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad [1]$$

111 For cases where  $q$  is large, the proportion of variance of  $\mathbf{x}_s$  that can be explained by  
112 regression on  $L_X$  ( $R_{x_s}^2$ ) can be estimated by regarding both  $\boldsymbol{\beta}$  and  $\boldsymbol{\varepsilon}$  as Gaussian independent  
113 random variables,  $\boldsymbol{\beta} \stackrel{iid}{\sim} N(0, \sigma_{\boldsymbol{\beta}}^2)$  and,  $\boldsymbol{\varepsilon} \stackrel{iid}{\sim} N(0, \sigma_{\boldsymbol{\varepsilon}}^2)$ . Upon appropriate scaling of the columns  
114 of  $\mathbf{X}$ ,  $R_{x_s}^2 = \frac{\sigma_{\boldsymbol{\beta}}^2}{\sigma_{\boldsymbol{\beta}}^2 + \sigma_{\boldsymbol{\varepsilon}}^2}$  can be interpreted as the proportion of variance of the phenotype that could  
115 be explained by regression on the features included in  $\mathbf{X}$ . The variance parameters involved  
116 ( $\sigma_{\boldsymbol{\beta}}^2$  and  $\sigma_{\boldsymbol{\varepsilon}}^2$ ) can be estimated using Bayesian or Likelihood methods (e.g., restricted maximum  
117 likelihood, REML, Patterson and Thompson 1971).

118 In the preceding paragraph we describe how one can estimate the proportion of variance  
119 of a vector in  $L_X$  ( $\mathbf{x}_s$ ) that can be explained by regression on the linear span of  $\mathbf{W}$ . Our goal is  
120 to generalize this to all vectors in  $L_X$ . However,  $L_X$  contains an infinite number of vectors;  
121 therefore, some approximation is needed.

122 Perhaps the most natural approach for estimating the proportion of variance of vectors in  
123  $L_X$  that can be explained by regression in  $L_W$  consists of regressing each of the columns of  $\mathbf{X}$   
124 on  $\mathbf{W}$ . Such an analysis would produce a sequence of  $R^2$  estimates  $\{R_{x_1}^2, R_{x_2}^2, \dots, R_{x_p}^2\}$ , and the  
125 average  $R^2$ ,  $R^2 = p^{-1} \sum_{s=1}^p R_{x_s}^2$ , could be used to estimate the overall proportion of variance of  
126  $\mathbf{X}$  that could be explained by regression on  $\mathbf{W}$ . However, one limitation of this approach is that  
127 the columns of  $\mathbf{X}$  are often not independent. Many features may cluster (e.g., genes may be co-  
128 expressed, or SNPs may be in high linkage-disequilibrium) leading to groups of highly  
129 unbalanced sizes. When some features are highly-correlated, the simple average of individual  
130  $R^2$ -values may lead to inaccurate (even biased) estimates. Furthermore, when  $\mathbf{X}$  is ultra-high  
131 dimensional (e.g., hundreds of thousands or million features) estimating  $R^2$ , on-feature-at-a-  
132 time will be computationally challenging. To address these problems, we discuss two methods  
133 that use independent vectors in the span of the output set.

134

135

136

137 **Method 1: Monte Carlo Analysis of Variance (MC-ANOVA)**

138 Since  $L_X$  is infinite, one we cannot exhaustively estimate  $R_{x_s}^2$  for all vectors in  $L_X$ . However,  
139 we can ‘explore’ the linear span of the output set by generating random vectors in  $L_X$  of the  
140 form  $\mathbf{x}_s = \mathbf{X}\boldsymbol{\alpha}_s$ , where  $\boldsymbol{\alpha}_s$  is sampled from some distribution. This can be repeated for a large  
141 number of vectors in  $L_X$  to produce a sequence of estimates  $\{R_{x_s}^2\}$ , and the resulting sequence  
142 can be used to estimate the average proportion of variance explained as well as other features  
143 of the distribution of the sequence. The method is summarized in Box 1. Importantly, if  $\boldsymbol{\alpha}_s$  and  
144  $\boldsymbol{\alpha}_{s'}$  are independent, so will be  $\mathbf{x}_s$  and  $\mathbf{x}_{s'}$ .

145

**Box 1. Monte Carlo Analysis of Variance (MC-ANOVA)**

- (1) Draw a random vector  $\boldsymbol{\alpha}_s$  from a proper multivariate distribution
- (2) Form the linear combination  $\mathbf{x}_s = \mathbf{X}\boldsymbol{\alpha}_s$
- (3) Estimate the proportion of variance of  $\mathbf{x}_s$  ( $R_{x_s}^2$ ) using a random effects model  
(expression [1]) with variance parameters estimated using either Bayesian or  
likelihood-type methods.
- (4) Repeat 1-3 for a large number of times.
- (5) Obtain a global  $R^2$  estimate by averaging the  $R_{x_s}^2$  's in the sequence.

146

147

148 In Box 1 we did not specify how the  $\boldsymbol{\alpha}_s$  are generated. One possibility is to sample these  
149 coefficients from distributions with support in  $R^p$  (e.g.,  $p$ -variate Gaussian). Alternatively, one  
150 could sample sparse vectors of coefficients from finite mixtures with a point of mass at zero.  
151 The possibility of using different process for generating random vectors in  $L_X$  gives the MC-  
152 ANOVA a great deal of flexibility—we will further explore that flexibility in greater detail in  
153 the case studies presented further below.

154

155 **Method 2: Regression using orthogonal basis (Eigen-ANOVA)**

156 An orthogonal basis for the row-space of  $\mathbf{X}$  can be obtained from the singular-value  
157 decomposition of  $\mathbf{X} = \mathbf{U}_X \mathbf{D}_X \mathbf{V}_X'$ , where  $\mathbf{U}_X$  and  $\mathbf{V}_X$  are the left- and right-singular vectors of  
158  $\mathbf{X}$  respectively, and  $\mathbf{D}_X$  is a diagonal matrix with the singular values of  $\mathbf{X}$  in the diagonal. Both  
159  $\mathbf{U}_X$  and  $\mathbf{V}_X$  are orthonormal, thus  $\mathbf{U}_X' \mathbf{U}_X = \mathbf{I}$  and  $\mathbf{V}_X' \mathbf{V}_X = \mathbf{I}$ . Each vector in  $L_X$  can be  
160 represented as a linear combination of the left-singular vectors of  $\mathbf{X}$ . Therefore, our second

161 method estimates the proportion of variance of each of the left-singular vectors of  $\mathbf{X}$  that can  
162 be explained by regression on  $\mathbf{W}$ , and produces a global  $R^2$  estimate using a weighted sum of  
163 the  $R^2$  estimated for each singular vector (Box 2).

164

### Box 2. Eigen-ANOVA

- (1) Generate an orthogonal basis for  $L_X$ ; for instance, compute the singular-value decomposition of  $\mathbf{X} = \mathbf{U}_X \mathbf{D}_X \mathbf{V}_X'$  where  $\mathbf{U}_X' \mathbf{U}_X = \mathbf{I}$  and  $\mathbf{V}_X' \mathbf{V}_X = \mathbf{I}$  are orthonormal basis for the row- and column-space of  $\mathbf{X}$  respectively, and  $\mathbf{D}_X = \text{Diag}\{d_i\}$  is a diagonal matrix with the singular values of  $\mathbf{X}$  in its diagonal.
- (2) Regress each of the left-singular vectors on  $L_W$  using a linear model such as that in expression [1] with  $\mathbf{u}_i = \mathbf{x}_s$ , and estimate the proportion of variance of each vector that can be explained by regression on  $L_W$ ,  $R_{u_i}^2$ .
- (3) Estimate the global proportion of variance of vectors in  $L_X$  that can be explained by regression on  $L_W$  using  $R^2 = \frac{\sum_{i=1}^r d_i^2 R_{u_i}^2}{\sum_{i=1}^r d_i^2}$ .

165

166

## 167 3 Results

168 In this section we first present results from simulations designed to detect bias on estimates  
169 derived from each of the methods, and to compare the bias of the proposed methods with that  
170 of the PLS—a method commonly used in multivariate-high-dimensional regressions.

### 171 3.1 Statistical properties assessed via simulations

172 **Data** were simulated using genotypes from a wheat data set generated by the International  
173 maize and Wheat Improvement Center (CIMMYT). Briefly, the data set provides genotypes at  
174 1,279 molecular markers assessed in 599 wheat inbred lines. Further details about this data set  
175 are provided by de los Campos et al. (2009). The data set is available with the BGLR R-package  
176 (de los Campos and Perez-Rodriguez, 2015). The scripts used to conduct these simulations are  
177 provided in the Supplementary Data (File S1). The REML estimation algorithm used to  
178 implement the MC- and Eigen-ANOVA is also provided in File S1. (We also tested the  
179 estimator proposed by Schreck et al. 2019; the results were very similar to the REML estimates  
180 presented below.) As benchmark, we used the PLS regression method. Briefly, we regressed  
181 the output matrix ( $\mathbf{X}$ ) on the input data ( $\mathbf{W}$ ) using the `pls` function of the `pls` R-package (Mevik,  
182 Wehrens, and Liland 2019). The number of orthogonal basis used was determined using cross-

183 validation—we choose the number of components that maximized prediction accuracy. The R-  
184 squared was then computed for each feature using the entire data set, and an overall R-squared  
185 was obtained by averaging the R-squared obtained for each of the columns of  $\mathbf{X}$ . The  
186 implementation can be found in File S1.

187

### 188 ***Simulation 1***

189 In our first simulation we use the genotypes of the wheat data set as the input set:  $\mathbf{W} =$   
190  $[\mathbf{w}_1, \dots, \mathbf{w}_q]$ . The ‘output’ set,  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_p]$ , was simulated using:  $\mathbf{x}_i = \mathbf{w}_i + \boldsymbol{\delta}_i$  where  
191  $i=1, \dots, 1,279$  indexes molecular markers and  $\boldsymbol{\delta}_i$  is an  $n$ -dimensional vector of *iid* (independent  
192 and identically distributed) random draws from a normal distribution with zero mean and a  
193 variance parameter, such that the proportion of variance of  $\mathbf{x}_i$  explained by  $\mathbf{w}_i$  was, 0, 0.1, 0.3,  
194 0.5, 0.8, 0.9 and 1. Here, 0 represents the case of complete independence ( $\mathbf{x}_i = \boldsymbol{\delta}_i$ ) and 1 the  
195 case of perfect concordance ( $\mathbf{X}=\mathbf{W}$ ).

196 We conducted a total of 1,000 Monte Carlo (MC) replicates per scenario. The input set ( $\mathbf{W}$ )  
197 did not change across MC samples, however, the output set,  $\mathbf{X}$ , varied across MC replicates  
198 due to the noise term ( $\boldsymbol{\delta}_i$ ). For each simulated data set we then estimated the proportion of  
199 variance of  $\mathbf{X}$  explained by regression of  $\mathbf{W}$  using random effects models, with variance  
200 parameters estimated using REML (Patterson and Thompson 1971).

201 The Monte Carlo method estimated the proportion of variance of  $\mathbf{X}$  explained by  $\mathbf{W}$  without  
202 any noticeable bias (**Table 1**). However, the regression of the eigenvectors of  $\mathbf{X}$  on  $\mathbf{W}$  produced  
203 estimates that, in some scenarios, were downwardly biased. The presence of bias was evident  
204 in scenarios where the true proportion of variance of  $\mathbf{X}$  explained by  $\mathbf{W}$  was greater than 0.5.  
205 Further inspection of the results for individual MC replicates suggested that the bias of the  
206 Eigen-ANOVA method was likely due to a relatively large number of ‘corner’ solutions (zero  
207 estimated proportion of variance) which were common for high-order eigenvectors (i.e., those  
208 with small eigenvalue)—we illustrate this in an analysis of multi-omic cancer data further below.  
209 The  $R^2$  estimates obtained with PLS also had noticeable biases which, in most scenarios, were  
210 larger in absolute value than the bias estimates obtained with Eigen-ANOVA.

211

212

213

214

215

216

217

218

219 **Table 1.** Average (SD) estimate of the proportion of variance explained by simulation scenario  
 220 (first column) and estimation method (Simulation 1).

Simulated Proportion of Variance Explained	Monte Carlo-ANOVA	Eigen-ANOVA	PLS
0.0	0.0082 (0.0028)	0.0081 (0.0006)	0.0017 (0.0001)
0.1	0.1002 (0.0083)	0.0983 (0.0019)	0.0478 (0.0034)
0.3	0.2991 (0.0108)	0.3020 (0.0028)	0.2412 (0.0075)
0.5	0.4992 (0.0102)	0.5054 (0.0028)	0.4865 (0.0076)
0.8	0.8006 (0.0055)	0.7857 (0.0017)	0.8451 (0.0036)
0.9	0.9012 (0.0033)	0.8685 (0.0011)	0.9403 (0.0016)
1.0	1.0000 (<.0001)	0.9377 (<.0001)	0.9988 (<.0001)

221

222 **Simulation 2**

223 We designed a second simulation to consider the case where one of the sets ( $X$ ) was included  
 224 in the other set ( $W$ ). We achieved this as follows:

225 - We set  $X$  to be a matrix containing a subset of the wheat genotypes. Specifically, we  
 226 sampled at random 128 (10%), 256 (20%), 640 (50%), 895 (70%) or 1151 (90%) of the  
 227 available diversity arrays technology (DArT) markers and formed with those genotypes our  $X$   
 228 matrix.

229 - Subsequently,  $W=[X,Z]$ , was formed by combining in a single matrix the columns of  $X$   
 230 and a matrix ( $Z$ ) whose entries were filled with *iid* draws from standard normal distributions.  
 231 The columns of  $X$  and  $W$  were all centered and scaled to unit variance; therefore, the proportion  
 232 of variance of  $W$  that can be explained by  $X$  equals the ratio of the number of columns of  $W$   
 233 that are shared with  $X$ . On the other hand, since  $L_X \in L_W$ , the true proportion of variance of  $X$   
 234 explained by  $W$  was always 1.

235

236



237 In our second simulation study the MC-ANOVA method rendered nearly unbiased estimates  
 238 of the proportion of variance of one set explained by the other (**Table 2**). However, the Eigen-  
 239 ANOVA method and the PLS produced noticeable biases. In most cases, both methods were  
 240 downwardly biased; however, the PLS had an upward bias in a few scenarios.

241

242 **Table 2.** Average (SD) REML estimates of the proportion of variance explained by simulation  
 243 scenario (first column) and estimation method (second simulation).

Scenario  <i># Columns of X</i> <i># Columns of W</i>	X regressed on W			W regressed on X		
	MC-ANOVA	Eigen-ANOVA	PLS	MC-ANOVA	Eigen-ANOVA	PLS
0.05	0.9960 (0.0039)	0.9085 (0.0051)	0.8885 (0.0069)	0.0505 (0.0050)	0.0548 (0.0012)	0.0244 (0.0029)
0.10	0.9972 (0.0030)	0.8891 (0.0041)	0.9193 (0.0036)	0.1000 (0.0072)	0.1061 (0.0018)	0.0652 (0.0038)
0.30	0.9964 (0.0025)	0.8835 (0.0024)	0.9781 ( $<.0001$ )	0.2999 (0.0106)	0.3060 (0.0028)	0.2656 (0.0068)
0.50	0.9943 (0.0028)	0.8989 (0.0019)	0.9954 ( $<.0001$ )	0.4996 (0.0102)	0.4977 (0.0030)	0.4902 (0.0072)
0.80	0.9965 (0.0013)	0.9223 (0.0010)	0.997 ( $<.0001$ )	0.8000 (0.0061)	0.7714 (0.0025)	0.8259 (0.0047)
0.90	0.9992 (0.0005)	0.9302 (0.0008)	0.9979 ( $<.0001$ )	0.9008 (0.0039)	0.8593 (0.0019)	0.9277 (0.0035)
0.95	0.9998 (0.0002)	0.9345 (0.0008)	0.9984 ( $<.0001$ )	0.9511 (0.0026)	0.9016 (0.0013)	0.9746 (0.0025)

244

### 245 **3.2 Applications to experimental data**

246 We used the MC-ANOVA and Eigen-ANOVA to quantify the proportion of variance explained  
 247 in two experimental data sets. The first data set contains a set of ultra-high-density (UHD,  
 248 1million+SNPs) SNPs derived from a combination of whole-genome sequencing (WGS) and  
 249 imputation, and a set of (in-silico created) low-density SNP panels. We use this data set to  
 250 assess the proportion of variance of UHD genotypes that can be captured and predicted using  
 251 low-density SNP sets. The second data set involved three omic-layers (GE, ME, and copy-  
 252 number-variants, CNV) of female breast cancer tumors. We used this data set to assess the  
 253 proportion of variance at one omic that can be explained by another omic.

254

255 **Case-study 1: Multi-locus linkage disequilibrium between ultra-high-density SNP**  
256 **genotypes and low-density SNP sets in chicken genomes**

257 The continued reduction of genotyping and sequencing costs has led to a sustained increase in  
258 the number of loci that can be genotyped. In plant and animal breeding four typical genotyping  
259 options include: (i) customized low density arrays with anywhere between hundreds or a few  
260 thousand SNPs, (ii) commercial arrays of common SNPs with circa 50K (K=1000) SNPs, (iii)  
261 high density SNP arrays with a number of markers of ~0.5M (M=million) SNPs and (iv) whole-  
262 genome sequence-derived SNP genotypes. The number of SNPs that can be derived from WGS  
263 varies between populations, but is usually of the order of tens of millions (UHD SNP  
264 genotypes). In recent years, several projects have produced large volumes of fully sequenced  
265 genomes for various agricultural species and model organisms. However, generating, storing,  
266 and fitting models with UHD-genotypes can be logistically, economically, and  
267 computationally challenging. Moreover, empirical evidence seems to suggest that using UHD  
268 SNP-genotypes does not lead to substantial gains in prediction accuracy relative to models  
269 trained using tens of thousands of SNPs (Erbe et al. 2013; Ober et al. 2012). This often leads  
270 researchers and the industry wondering: *How many SNPs are needed to capture (almost all)*  
271 *the information contained in UHD SNP genotypes?* We used the MC- and Eigen-ANOVA  
272 methods to address precisely this question.

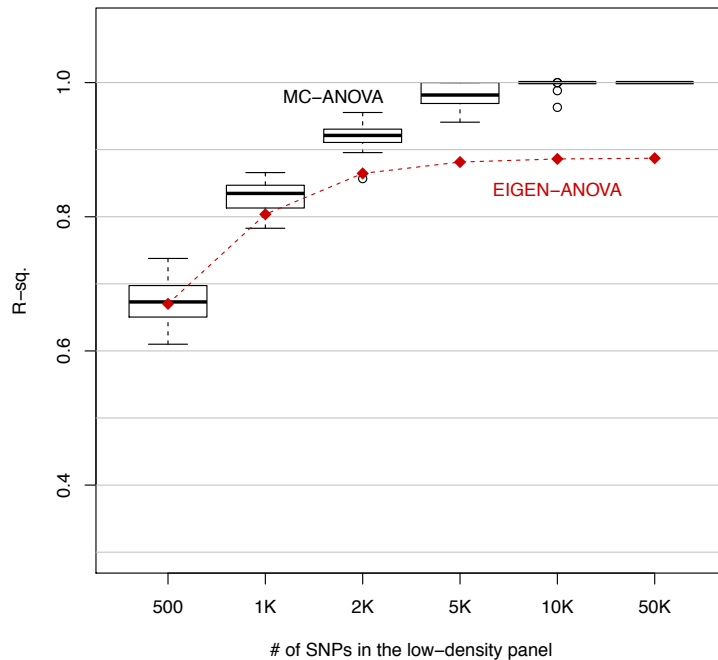
273 **Data** consisted of UHD SNP genotypes of 892 female and male chickens from six  
274 generations of a purebred commercial brown layer line of Lohmann Tierzucht GmbH. The  
275 genomes of 25 layers were sequenced at 8x read-depth, from the sequence data 4.92M  
276 (M=million) SNPs were derived. The remaining layers (n=867) were genotyped using the  
277 Affymetrix Axiom Chicken Genotyping Array (Kranis et al. 2013) which contains ~600K  
278 (580,961) SNPs, and 336,224 passing quality control and filtering. The genotypes of these  
279 layers were then combined with those derived from WGS (4.92M SNPs) by first imputing the  
280 HD-data from the chip using BEAGLE 3.3.2 (Browning and Browning 2009) and afterwards  
281 imputing to UHD (4.92M SNPs) via MiniMac3 (Howie et al. 2012) with pre-phasing  
282 performed in BEAGLE 4.0. For details on the imputing procedure we refer to Ni et al. (2015).  
283 From the 4.92M SNPs we remove those with minor-allele-frequency smaller than 0.005 (0.5%)  
284 and pruned SNPs to guarantee a maximum  $R^2$  between adjacent SNPs of 0.99; these editing  
285 criteria left a total of 1.79M SNPs which were used for the subsequent analysis.

286  
287  
288

289 ***Proportion of variance at ultra-high-density genotypes explained by low-density SNP-sets***

290 In a first analysis, the output space was the linear space ( $L_X$ ) spanned by the UHD SNP  
291 genotypes. The input set, ( $L_W$ ), consisted of low-density genotypes obtained by selecting  $p$   
292 ( $p=500, 1K, 2K, 3K, 5K$  and  $10K$ ) evenly-spaced (in variant counts) SNPs. We estimated the  
293 proportion of variance captured by low-density panels using the Eigen- and MC-ANOVA. For  
294 the MC methods we sampled weights from *iid* standard normal distribution,  $\alpha_{js} \sim N(0,1)$  and  
295 then form a random vector in  $L_X$  using  $\mathbf{x}_s = \mathbf{X}\alpha_s$ , where  $\mathbf{X}$  is the matrix of UHD SNP-  
296 genotypes. These random vectors were then regressed on the lower-density SNP sets, and the  
297 proportion of variance explained was estimated using REML. This was repeated 1,000 times  
298 to estimate the distribution of proportion of variance of vectors in  $L_X$  explained by each of the  
299 low-density SNP-sets. For the Eigen-ANOVA we regressed each of the eigenvectors of the  
300 UHD SNP genotypes on the low-density panels.

301 According to the MC-ANOVA method, the panel containing 500 evenly-spaced SNPs  
302 captured about two thirds of the variance spanned by the UHD SNP genotypes (**Figure 1**). The  
303 proportion of variance of the UHD SNPs explained by low-density panels increased with the  
304 number of SNPs in the low-density panels reaching 100% with  $p=10K$  SNPs. The variance in  
305 the proportion of variance captured by low-density panels also decreased with the number of  
306 SNPs in the array (Figure 1). The Eigen-ANOVA yielded a very similar estimate of proportion  
307 of variance explained as the MC-ANOVA for  $p=500$ . However, for SNP-panels with more than  
308 500 SNPs, the estimated proportion of variance obtained with the Eigen-ANOVA was  
309 systematically lower than the one obtained with MC-ANOVA. This agrees with what we found  
310 in the simulations where for high R-sq. the Eigen-ANOVA method gave downwardly biased  
311 estimates. (Note that while the MC-ANOVA yields both a point estimate and measures of  
312 dispersion (across random vectors) of  $R^2$ , the Eigen-ANOVA only yields the point-estimates  
313 which are shown in Figure 1.)



**Figure 1.** Proportion of the variance of whole-genome-sequence-derived SNPs (1.79 million) explained by SNP-panels consisting of 500, to 50K (K=1000) evenly-spaced SNPs.

314

315

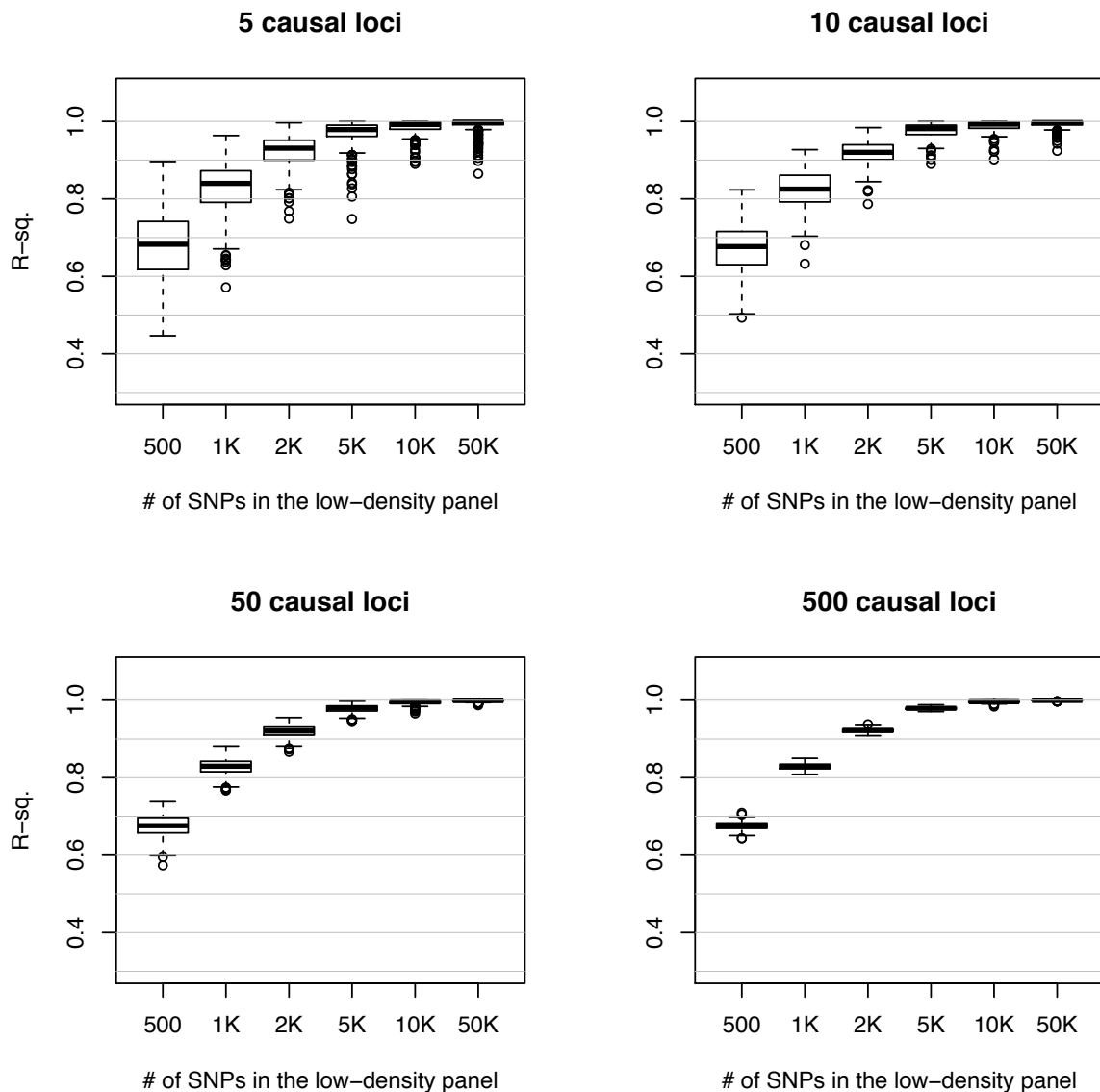
316

### 317 ***Quantifying the effect of trait-complexity***

318 In the previous application in the MC methods we drew random effect vectors that had weights  
319 (drawn from a normal distribution) on all the SNPs of the UHD set. However, for any trait, the  
320 vast majority of variants in the genome are expected to have no effect. The number of variants  
321 affecting any trait could vary from very few (simple traits) to hundreds or thousands (complex  
322 traits). Therefore, to explore the effect of the trait architecture on the distribution of the  
323 proportion of genetic variance of those traits that could be captured by low-density SNP panels,  
324 we repeated the previous analyses using random vectors that had  $q$  (with  $q=5, 10, 50, 500$ ) non-  
325 zero weights – the set of SNPs with non-zero weight were randomly sampled from the UHD-  
326 genotypes, and the weights of those SNPs were *iid* normal.

327 The estimated proportion of variance explained by regression on lower-density SNP panels  
328 were, on average, the same across “trait-architectures” (Figure 2). However, the dispersion  
329 about the estimated means was, as expected, much larger for simple traits. For “complex traits”  
330 with 500 “causal variants” the proportion of variance explained by regression on 10K or more  
331 SNPs was greater than 95% for all MC replicates. However, for simpler traits (e.g., 5 ‘causal  
332 variants’) we had some random vectors with proportion of variance explained smaller than 0.8.

333



334

335 **Figure 2.** Proportion of variance of random vectors derived from ultra-high-density SNPs  
336 explained by regression on low-density SNP-panels, by number of loci used to form “genetic  
337 traits”.

338

339 **Case Study 2: Proportion of variance explained in multi-omic data from breast cancer**  
340 **tumors**

341

342 Cancerous processes involve the deregulation of signaling pathways controlling cell fate and  
343 progression, arising from the accumulation of genomic and epigenomics alterations across  
344 multiple genes (Vogelstein et al. 2013; Witte, Plass, and Gerhauser 2014). Genetic and

345 epigenetic modifications can lead to changes in GE, which in turn can lead to changes in  
346 downstream (e.g., protein expression) and upstream (e.g., DNA, ME) processes, thus resulting  
347 in complex multivariate association patterns between multiple omic-layers.

348 **Data:** We used GE, ME and CNV data from breast cancer tumors to study multivariate  
349 associations between those three omics. Data (n=593) was from The Cancer Genome Atlas  
350 (TCGA), and consisted of samples from frozen primary breast cancer tumors from female  
351 patients.

352 Gene expression data (RNA-Sequencing counts) were determined using the Illumina HiSeq  
353 RNA V2 platform and DNA methylation profiles were determined using the Illumina HM450  
354 platform. RNA-sequencing data were transformed using the natural logarithm and individual  
355 CpG site  $\beta$ -values were summarized at the CpG island level, using the maximum connectivity  
356 approach implemented in the WGCNA R package (Langfelder and Horvath 2008). The CpG  
357 island summaries were transformed into M-values ( $M=\beta/(1-\beta)$ ; Du et al. 2010) CNV profiles  
358 corresponded to gene-level copy number intensity derived from Affymetrix SNP Array 6.0  
359 platform, using hg19 as reference.

360 **Data editions:** From each of the three omics we removed features with coefficient of  
361 variation smaller than 1% and those with proportion of missing values greater than 20%. The  
362 missing values that remained were imputed using a k-nearest neighbors clustering algorithm,  
363 with  $k = 3$  (Hastie et al. 2016). After imputation, each feature was adjusted by batch effects  
364 using ComBat (Lazar et al. 2013). After applying the steps described above, the data set used  
365 in the analyses consisted of the (log-transformed) expression of 20,319 genes, 11,552 CVN-  
366 sites and ME intensity at 28,241 ME CpG islands.

367 **Results:** We used the MC- and the Eigen-ANOVA methods to estimate the proportion of  
368 variance of one omic that can be explained by regression on another omic; we did this for all  
369 pairwise omics combinations (GE~ME, GE~CNV, ME~GE, ME~CVN, CNV~GE and  
370 CVN~ME). Our results with the MC-ANOVA method indicate that the CNV data were  
371 completely explained by both GE and ME (Table 3). About 70% of the variance spanned by  
372 ME was explained by GE and vice versa. Finally, CNV explained a relatively small fraction of  
373 the variance spanned by either GE or ME. These results suggest that most CNVs have effects  
374 in both ME and GE and therefore, variation in CNV can be predicted by ME and GE. However,  
375 although there is association between CNV and both ME and GE, many other factors (e.g.,  
376 environmental effects) seem intervene, thus, making the proportion of GE or ME explained by  
377 CNV relatively small (~20%). Overall the MC- and Eigen-ANOVA methods yielded similar

378 results. However, in cases involving high  $R^2$  (CNV~ME, CNV~GE, GE~ME and ME~GE) the  
 379 Eigen-ANOVA method gave  $R^2$  estimates that were lower than those of the MC method. This  
 380 pattern is consistent with what we observed in the simulation and in the analyses of chicken  
 381 genomes.

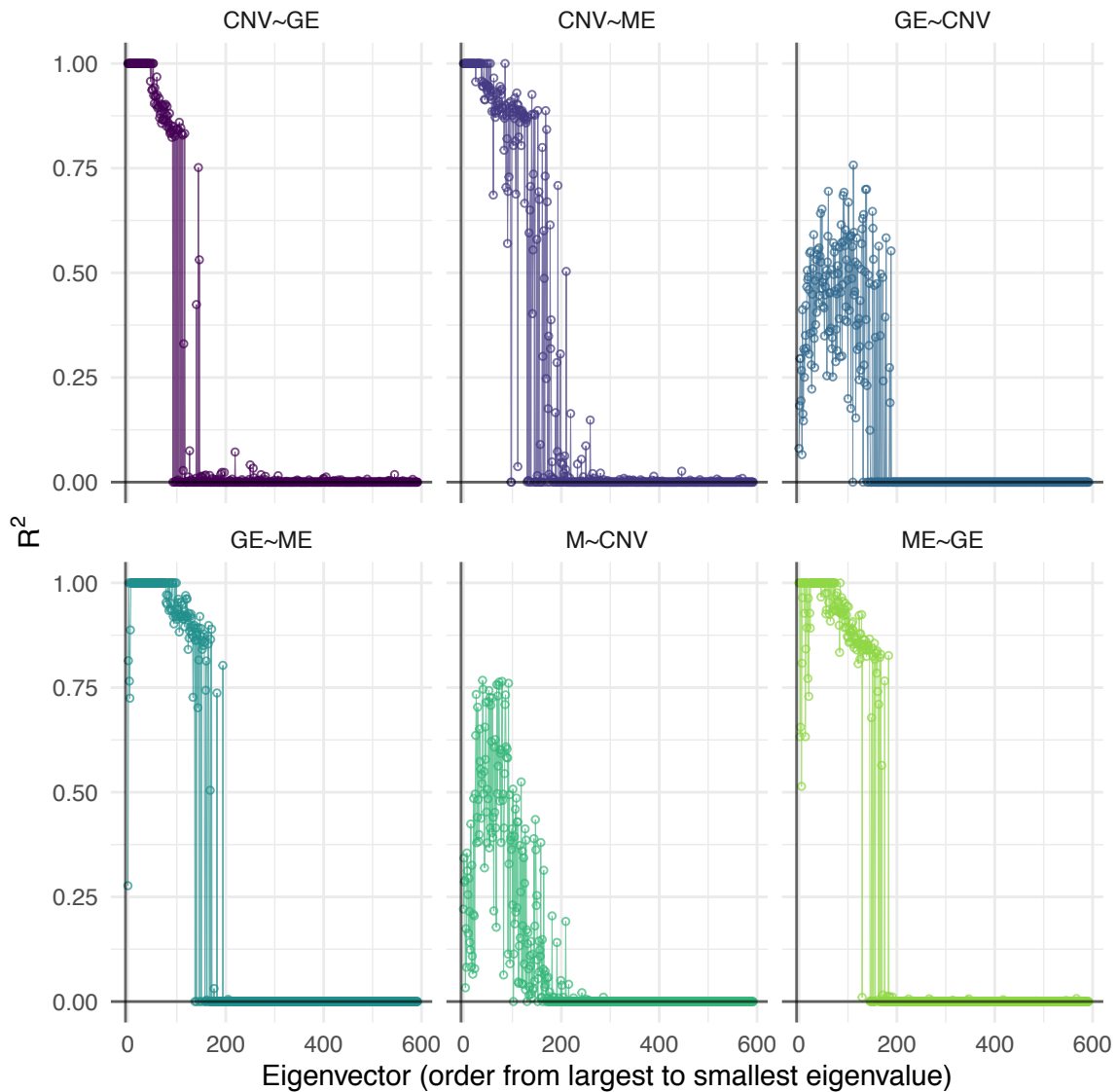
382

383 **Table 3.** Proportion of variance of one omic explained by regression of the omic in each row  
 384 on the omic in each column obtained with MC-ANOVA (Eigen-ANOVA).

Dependent	Explanatory		
	CNV	Methylation	Gene Expression
CNV	---	1.00 (0.929)	1.00 (0.904)
Methylation	0.164 (0.228)	---	0.715 (0.685)
Gene Expression	0.204 (0.238)	0.738 (0.660)	---

385

386 Eigen-vector-specific  $R^2$  values obtained with the Eigen-ANOVA method (Figure 3)  
 387 showed that the  $R^2$  values were, in most cases (except GE~CNV and M~CNV) very high (and  
 388 in many cases very close to one) for the top-eigenvectors (i.e., those with high eigenvalue), and  
 389 very small for eigenvectors associated with low eigenvalues. The transition in the  $R^2$  profile of  
 390 individual eigenvectors showed a relatively sharp phase transition from  $R^2$  values near one to  
 391 near-zero values. Overall, our results suggest a relatively good agreement in the patterns  
 392 captured by the top-eigenvectors across omics.



393

394

395

396

397

398

399

400

401

402

403

404

405

**Figure 3.** Proportion of variance of omic-derived eigenvectors of an omic-set explained by regression on a different omic-set. Points give the proportion of variance for individual eigenvectors. GE=Gene Expression, ME=Methylation, CNV=Copy-number variants (global  $R^2$  estimates, derived from random vectors and from the Eigen-ANOVA method are shown in Table 3).



## 406 **4 Discussion**

407

408 Modern genomic data sets often combine information from multiple data-layers (e.g., pedigree,  
409 DNA-sequence, epigenomic information, gene expression, proteomics, metabolomics). The  
410 biological processes underlying these data can lead to complex dependencies between data-  
411 layers. MANOVA can be used to quantify the proportion of variance explained in multivariate  
412 settings. However, MANOVA is based on least-squares projections which are not-well suited  
413 for analysis of high-dimensional data. Reduced-rank regression methods (e.g. PLS, CCA) and  
414 penalized regressions (e.g., Lasso-type methods) can be used to confront the problems  
415 emerging when the number of features exceed sample size ( $p \gg n$ ). However, these methods  
416 are not adequate for estimation of proportion of variance explained, because they rely on  
417 regularization decisions (e.g., choosing the number of dimensions, or selecting a penalization  
418 parameter that controls sparsity) which are often based on cross-validation procedures that are  
419 not well-suited for inferences.

420 To overcome the limitations of existing methods, we developed two procedures to estimate  
421 the proportion of variance explained in settings where both the input and output sets are high-  
422 dimensional. The proposed approach uses random effects Gaussian models to estimate the  
423 proportion of variance of (independent) vectors in the linear span of an output set ( $\mathbf{X}$ ) that can  
424 be explained by regression on an input set ( $\mathbf{W}$ ). The resulting  $R^2$  estimate is a weighted average  
425 of the  $R^2$  values obtained for independent vectors. We considered two approaches to generate  
426 independent vectors in the linear span of the output set: The first one (MC-ANOVA) is a Monte  
427 Carlo method that uses randomly generated (independent) vectors in the linear span of the  
428 output set. The second one (Eigen-ANOVA) uses an orthogonal basis for the linear span of  $\mathbf{X}$ .

429 The two proposed methods have four important features: (i) Both methods can be used to  
430 perform analysis of variance when both explanatory and dependent data are high-dimensional;  
431 (ii) Estimates are entirely based on the likelihood function and there is no need to make  
432 regularization decisions (number of dimensions, penalty parameters). (iii) For any pair of  
433 information sets, the analysis of variance is not necessarily symmetric; therefore, the approach  
434 accommodates cases where the proportion of variance of  $\mathbf{W}$  explained by  $\mathbf{X}$  is not equal than  
435 the reciprocal. (iv) Finally, in addition to producing an  $R^2$  estimate, the proposed methods can  
436 shed light on important aspects of the underlying association patterns (e.g., decomposition of  
437 the global  $R^2$  on eigen-vector specific  $R^2$ 's, distribution of  $R^2$  over possible vectors in the linear  
438 span of the output set).

439

440 Our simulations suggest that MC-ANOVA renders nearly unbiased estimates of the  
441 proportion of the variance of one set that can be explained by another. However, the Eigen-  
442 ANOVA exhibited small but systematic biases in scenarios in which the true proportion of  
443 variance was either too low or very high. The biases of Eigen-ANOVA were comparable, in  
444 magnitude, with those of the PLS regression. Therefore, for estimation of proportion of  
445 variance explained we recommend using MC-ANOVA.

446 The MC-ANOVA has clear computational advantages relative to Eigen-ANOVA and PLS  
447 because this method can render relatively accurate estimates of  $R^2$  based on a few hundreds of  
448 random vectors. Therefore, the number of regressions that one may need to perform can be  
449 much smaller than with PLS and the Eigen-ANOVA.

450 Consistent with our simulation results, the analyses of experimental data showed that in  
451 problems involving a high  $R^2$  the Eigen-ANOVA method yielded lower estimates of the  
452 proportion of variance explained than those obtained with the MC-ANOVA (e.g., see Figure 1  
453 and Table 3). Inspection of the results of the Eigen-ANOVA for individual eigenvectors  
454 suggests that the downward bias of the method may originate from corner solutions (zero-  
455 estimates of  $R^2$ ) for eigenvectors associated with small eigenvalues. Therefore, if the only goal  
456 is to estimate the proportion of variance of one set explained by another set, we recommend  
457 using the MC-ANOVA method.

458 The Eigen-ANOVA method yields  $R^2$ -values for each of the eigenvectors of the output set.  
459 This information can help elucidate whether global patterns (e.g., those associated with the top-  
460 eigenvectors) in one information set can be predicted from information contained in another  
461 information set. For instance, our analysis of the multi-omic breast cancer revealed that the  
462 patterns described in the top-eigenvectors derived from GE and ME are very similar; therefore,  
463 one should not expect big differences in tumor classifications that are based on the top-  
464 eigenvectors derived from either set. Interestingly, we found that in the analyses of omic data  
465 the  $R^2$  of individual eigenvectors showed a very sharp phase transition, suggesting that  
466 eigenvectors associated with intermediate and small eigenvalues may describe omic-specific  
467 patterns, or perhaps measurement error associated to each of the techniques.

468 The MC-ANOVA method can be used to characterize the distribution the  $R^2$  estimates  
469 across vectors in the linear span of the output set. We used this feature to study the effect of  
470 the trait-architecture; on the distribution of the  $R^2$  estimates. Our results indicate that while the  
471 average  $R^2$  does not depend on the distribution of the coefficients used to form random vectors  
472 (i.e., the  $\alpha$ 's), the dispersion of the distribution is highly dependent on the process used to  
473 generate the weights. More specifically, random vectors that have non-zero weights on a small

474 number of vectors have a larger dispersion in the distribution of the  $R^2$ , compared to the  
475 dispersion observed when the random vectors have non-zero weights for all the basis in the  
476 linear span.

477 An important feature of the methods proposed in this study is that the  $R^2$  measure is not  
478 symmetric, in contrast to CCA. Our simulation study shows that if the underlying patterns are  
479 non-symmetric (e.g., when one of the linear spaces is a subspace of the other) the proposed  
480 estimation methods (in particular the MC-ANOVA) can detect the lack of symmetry very well  
481 (see Table 2). Interestingly, our analysis of multi-omic data from breast cancer patients  
482 exhibited cases where  $R^2$  was rather symmetric (e.g., the regression ME~GE and the regression  
483 ME~GE) and others that were highly asymmetric (e.g., CNV~GE and GE~CNV). The  
484 asymmetric cases suggest that almost all the variability in CNV can be predicted from GE (and  
485 ME as well); however, only a fraction of the GE variance can be explained by differences in  
486 CNV patters. This result is consistent with the hypothesis that most CNV have an impact on  
487 GE, but GE is also affected by factors other than CNV (e.g., methylation, environmental  
488 effects).

489 **In conclusion:** We developed two methods for estimating the proportion of variance  
490 explained in problems in which both the input and output sets are high-dimensional. The MC-  
491 ANOVA method provided nearly unbiased estimates across a range of simulation scenarios. In  
492 addition to providing estimates of proportion of variance explained, the two methods can yield  
493 useful insight into the association patterns underlying multi-layered high-dimensional data.

494

495 **Funding:** Part of the study was conducted while GdlC was visiting the Research Training  
496 Group “Scaling Problems in Statistics” (RTG 1644) at the University of Goettingen, funded by  
497 the German Research Association (DFG).

498

499

## REFERENCES

500

- 501 Du, Pan, Xiao Zhang, Chiang-Ching Huang, Nadereh Jafari, Warren A Kibbe, Lifang Hou,  
502 Simon M Lin, et al. 2010. “Comparison of Beta-Value and M-Value Methods for  
503 Quantifying Methylation Levels by Microarray Analysis.” *BMC Bioinformatics* 11 (1).  
504 BioMed Central:587. <https://doi.org/10.1186/1471-2105-11-587>.
- 505 Erbe, Malena, Birgit Gredler, Franz Reinhold Seefried, Beat Bapst, and Henner Simianer.  
506 2013. “A Function Accounting for Training Set Size and Marker Density to Model the  
507 Average Accuracy of Genomic Prediction.” Edited by Zhanjiang Liu. *PLoS ONE* 8 (12).

- 508 Public Library of Science:e81046. <https://doi.org/10.1371/journal.pone.0081046>.
- 509 Hastie, Trevor, R Tibshirani, B Narasimhan, and Gilbert Chu. 2016. “Impute: Imputation for  
510 Microarray Data” 17:520–25.
- 511 Izenman, Alan Julian. 1975. “Reduced-Rank Regression for the Multivariate Linear Model.”  
512 *Journal of Multivariate Analysis* 5 (2). Academic Press:248–64.  
513 [https://doi.org/10.1016/0047-259X\(75\)90042-1](https://doi.org/10.1016/0047-259X(75)90042-1).
- 514 Kranis, Andreas, Almas A Gheyas, Clarissa Boschiero, Frances Turner, Le Yu, Sarah Smith,  
515 Richard Talbot, et al. 2013. “Development of a High Density 600K SNP Genotyping  
516 Array for Chicken.” *BMC Genomics* 14 (1). BioMed Central:59.  
517 <https://doi.org/10.1186/1471-2164-14-59>.
- 518 Krzanowski, W. J., and W. J. 1988. *Principles of Multivariate Analysis : A User’s*  
519 *Perspective*. Clarendon Press. <https://dl.acm.org/citation.cfm?id=59560>.
- 520 Langfelder, Peter, and Steve Horvath. 2008. “WGCNA: An R Package for Weighted  
521 Correlation Network Analysis.” *BMC Bioinformatics* 9 (1). BioMed Central:559.  
522 <https://doi.org/10.1186/1471-2105-9-559>.
- 523 Lazar, C., S. Meganck, J. Taminau, D. Steenhoff, A. Coletta, C. Molter, D. Y. Weiss-Solis,  
524 R. Duque, H. Bersini, and A. Nowe. 2013. “Batch Effect Removal Methods for  
525 Microarray Gene Expression Data Integration: A Survey.” *Briefings in Bioinformatics*  
526 14 (4). Oxford University Press:469–90. <https://doi.org/10.1093/bib/bbs037>.
- 527 los Campos, G de, H Naya, D Gianola, J Crossa, A Legarra, E Manfredi, K Weigel, and J M  
528 Cotes. 2009. “Predicting Quantitative Traits with Regression Models for Dense  
529 Molecular Markers and Pedigree.” *Genetics* 182:375–85.
- 530 los Campos, Gustavo de, and Paulino Perez-Rodriguez. 2015. “BGLR: Bayesian Generalized  
531 Linear Regression.”
- 532 Mardia, K. V., J. T. (John T.) Kent, and J. M. (John M.) Bibby. 1979. *Multivariate Analysis*.  
533 Academic Press.
- 534 Mevik, Bjørn-Helge, Ron Wehrens, and Kristian Hovde Liland. 2019. “Pls: Partial Least  
535 Squares and Principal Component Regression.”
- 536 Ni, Guiyan, Tim M. Strom, Hubert Pausch, Christian Reimer, Rudolf Preisinger, Henner  
537 Simianer, and Malena Erbe. 2015. “Comparison among Three Variant Callers and  
538 Assessment of the Accuracy of Imputation from SNP Array Data to Whole-Genome  
539 Sequence Level in Chicken.” *BMC Genomics* 16 (1):824.  
540 <https://doi.org/10.1186/s12864-015-2059-2>.
- 541 Ober, Ulrike, Julien F. Ayroles, Eric A. Stone, Stephen Richards, Dianhui Zhu, Richard A.

- 542 Gibbs, Christian Stricker, et al. 2012. “Using Whole-Genome Sequence Data to Predict  
543 Quantitative Trait Phenotypes in *Drosophila Melanogaster*.” Edited by Naomi R. Wray.  
544 *PLoS Genetics* 8 (5). Public Library of Science:e1002685.  
545 <https://doi.org/10.1371/journal.pgen.1002685>.
- 546 Patterson, H. D., and R. Thompson. 1971. “Recovery of Inter-Block Information When Block  
547 Sizes Are Unequal.” *Biometrika* 58 (3). Narnia:545–54.  
548 <https://doi.org/10.1093/biomet/58.3.545>.
- 549 Rencher, Alvin C., and William F. Christensen. 2012. *Methods of Multivariate Analysis*.  
550 Wiley Series in Probability and Statistics. Hoboken, NJ, USA: John Wiley & Sons, Inc.  
551 <https://doi.org/10.1002/9781118391686>.
- 552 Vogelstein, Bert, Nickolas Papadopoulos, Victor E Velculescu, Shibin Zhou, Luis A Diaz,  
553 and Kenneth W Kinzler. 2013. “Cancer Genome Landscapes.” *Science (New York, N.Y.)*  
554 339 (6127):1546–58. <https://doi.org/10.1126/science.1235122>.
- 555 Witte, Tania, Christoph Plass, and Clarissa Gerhauser. 2014. “Pan-Cancer Patterns of DNA  
556 Methylation,” 1–18.
- 557 Witten, D. M., R. Tibshirani, and T. Hastie. 2009. “A Penalized Matrix Decomposition, with  
558 Applications to Sparse Principal Components and Canonical Correlation Analysis.”  
559 *Biostatistics* 10 (3):515–34. <https://doi.org/10.1093/biostatistics/kxp008>.
- 560 Wold, Svante, Michael Sjöström, and Lennart Eriksson. 2001. “PLS-Regression: A Basic  
561 Tool of Chemometrics.” *Chemometrics and Intelligent Laboratory Systems* 58 (2).  
562 Elsevier:109–30. [https://doi.org/10.1016/S0169-7439\(01\)00155-1](https://doi.org/10.1016/S0169-7439(01)00155-1).
- 563