1  **MINOTAUR: A platform for the analysis and visualization of multivariate**
2  **results from genome scans with R Shiny**

3  Robert Verity*[1§], Caitlin Collins*[1], Daren C. Card[2], Sara M. Schaal[3], Liuyang Wang[4], Katie E.
4  Lotterhos[3]

5  *These authors contributed equally to this work.*

6  [1] Imperial College London, MRC Centre for Outbreak Analysis and Modelling, Department of
7  Infectious Disease Epidemiology, St. Mary's Campus, Norfolk Place, Imperial College
8  London, London W21PG, UK

9  [2] Department of Biology, 501 S. Nedderman Drive, The University of Texas at Arlington,
10  Arlington, TX 76019, USA

11  [3] Department of Marine and Environmental Sciences, 430 Nahant Road, Northeastern
12  University, Nahant, MA 01908, USA

13  [4] Department of Molecular Genetics and Microbiology, School of Medicine, 213 Research
14  Drive, Duke University, Durham, NC 27710, USA

15

16  [§]**Correspondence**: r.verity@imperial.ac.uk

17  **Running Head:** Multivariate outliers in genomics

18  **Abstract**

19  Genome scans are widely used to identify "outliers" in genomic data: loci with different

20  patterns compared with the rest of the genome due to the action of selection or other

21  non-adaptive forces of evolution. These genomic datasets are often high-dimensional,

22  with complex correlation structures among variables, making it a challenge to identify

23  outliers in a robust way. The Mahalanobis distance has been widely used for this

24  purpose, but has the major limitation of assuming that data follow a simple parametric

25  distribution. Here we develop three new metrics that can be used to identify outliers in

26  multivariate space, while making no strong assumptions about the distribution of the

27  data. These metrics are implemented in the R package MINOTAUR, which also includes

28    an interactive web-based application for visualizing outliers in high-dimensional

29    datasets. We illustrate how these metrics can be used to identify outliers from

30    simulated genetic data, and discuss some of the limitations they may face in application.

31    *Keywords*: genomic scans, Mahalanobis, kernel density

32    **Introduction**

33    Knowledge of the genetic architecture of biological traits —the number of loci that

34    affect a phenotype, the magnitude of their effect, and their distribution across the

35    genome—not only illuminates the evolutionary processes that shape genomes, but also

36    has important implications for complex diseases (McCarthy and Hirschhorn 2008),

37    conservation (Kohn et al. 2006; Allendorf et al. 2010; Funk et al. 2012), and breeding

38    programs (Goddard et al. 2009; Varshney et al. 2009). With the advent of next-

39    generation sequencing we now have the ability to examine genomes at a fine scale; and,

40    as a result, we have identified a large number of genomic variants that are implicated in

41    complex diseases (Carlson et al. 2004; Hindorff et al. 2009) and adaptation to the local

42    environment (Savolainen et al. 2013). This wealth of data is likely to yield new insights,

43    but it also brings with it the challenge of extracting the relevant signal from noisy,

44    complex, multi-dimensional data sets. This is perhaps one reason why most of the

45    variants detected so far have only managed to explain a very small proportion of the

46    observable phenotypic variation (Yang et al. 2010; Brachi et al. 2011).

47    The preferred method for detecting genomic variants is via genome scans. There are

48    many different approaches toward scanning genomes, but all are based on the same

49    premise: that the loci of interest to the investigator are likely to be statistical outliers

50    when compared with the rest of the genome. The particular choice of statistic will

51    depend on the question being asked and the experimental design, and may include one

52    or more statistics from the following categories: tests for genetic differentiation

53    (Lotterhos & Whitlock 2014; Hoban et al. *in revision*), scans for strong positive selection

54    and/or selective sweeps (Hohenlohe 2010; Vatsiou et al. 2016), genome-wide

55    association studies for phenotype-associated loci (GWAS, reviewed in Carlson et al.

56    2004 and McCarthy et al. 2008), linkage mapping for quantitative trait loci (QTL,

57    Savolainen et al. 2013), genetic-environment associations (reviewed in Rellstab et al.

58    2015), and scans for differentially expressed genes (Wang et al. 2009). A number of

59    different genome-scan test statistics may be calculated for a single genomic dataset and

60    these are usually examined one-at-a-time (i.e., in univariate analyses). Some test

61    statistics may be highly correlated, while the power of other test statistics may vary for

62    different regions of the genome depending on the details of selection, recombination,

63    mutation, and migration rates (Tiffin and Ross-Ibarra 2014). Additionally, the power of

64    different approaches may vary among species because of demographic history, and

65    within a species because of sampling design (De Mita et al. 2013; de Villemereuil et al.

66    2014; Lotterhos and Whitlock 2015). Finally, loci with intermediate probabilities of

67    detection will often exhibit the highest variance in results from genome scans

68    (Lotterhos et al. *in review*).

69    Given the complex evolutionary histories of most species, it is doubtful whether any

70    single statistic can fully capture the genomic signal of interest in the majority of cases

71    (Verity and Nichols 2014). Furthermore, the uncertainty in demographic history,

72    coupled with the variation in statistical outcomes in different scenarios, makes it

73    difficult to know which statistics have the greatest power to detect selection and which

74    have the highest false positive rates. These issues point to a need for composite,

75    multivariate outlier methods that integrate information across multiple test statistics.

76    Multivariate methods have been utilized extensively in many biological applications,

77    although in application to genome scans the power of the multivariate approach for

78    detecting outliers has not yet been fully evaluated. Because some dimension reduction

79    methods such as Principal Component Analysis rely on assumptions about the data that

80    may be unjustifiable in the context of genome scans (O'reilly et al. 2012), these methods

81    are not ideally designed for the identification of multivariate outliers (Pattterson et al.

82    2006). Some GWAS analyses have successfully employed multivariate approaches to

83    identify genetic associations with multiple phenotypes (O'reilly et al. 2012; Galesloot et

84    al. 2014).  Additionally, multivariate approaches have also been used in GWAS meta-

85    analysis to simultaneously consider multiple genetic or phenotypic variables (reviewed

86    in Evangelou and Ioannidis 2013). It is evident, however, that more opportunities exist

87    for the use of multivariate approaches in outlier detection than are currently being

88    capitalized on.

89    While there are dedicated software tools for calculating a variety of test statistics, there

90    does not currently exist a unified platform for the filtering, visualization, and integration

91    of test statistics in multivariate space. Here we describe a new R package called

92    MINOTAUR (Multivariate vIsualisatioN and OuTlier Analysis Using R) built specifically

93    for this purpose. This software package - initiated during a hackathon for population

94    genetics in R (https://github.com/NESCent/r-popgen-hackathon) - provides functions

95    for detecting outliers in multivariate space alongside procedures to manipulate,

96    summarize, and visualize these data. The R software environment (R Core Team 2015)

97    is free, open-source, and hosts a large collection of tools for statistical analysis, making

98    it the ideal host for the development and uptake of such a platform. Furthermore,

99    because data visualization is an important part of verifying and identifying outliers, the

100    R Shiny and Shiny Dashboard environments (Chang 2015; Chang et al. 2016) have been

101    employed to provide MINOTAUR users with an interactive interface that streamlines

102    the process of data input, statistical analysis, and graphical exploration. Together, these

103    tools have the potential to increase the efficiency with which the results of genome

104    scans are interrogated.

105    **Approaches to identifying multivariate outliers**

106    In the MINOTAUR package we implement four composite measures that can be used to

107    integrate information over multiple univariate statistics: the Mahalanobis distance,

108    harmonic mean distance, nearest neighbor distance, and kernel density deviance. We

109    developed the latter three measures, which are related to Mahalanobis distance but

110    make no strong assumptions about the parametric form of the data, meaning they can

111    be applied to multivariate statistics that have complex correlated or even multimodal

112    distributions. Some of these measures are heavily influenced by the distance of points

113    from the multivariate centroid (Mahalanobis and harmonic mean distance) while others

114    are mainly influenced by the sparseness of points in the local vicinity (nearest neighbor

115    distance and kernel density deviance), and so we would expect the measures to behave

116    differently from one another, and to vary in their behavior depending on the data at

117    hand.

118    The calculation of these composite measures has been optimized for genome-scale data

119    by using precompiled routines, written in C++ and integrated into R using the package

120    Rcpp (Eddelbuettel and Francois 2011; Eddelbuettel 2013). Several packages devoted

121    to multivariate statistics that may be appropriate for genome-scale data already exist in

122    R (see Supplementary Table 1), and thus users are free to utilize both existing statistical

123    methods and the more targeted functions included within the MINOTAUR package.

124    *Mahalanobis distance.* The Mahalanobis distance is a multidimensional measure of the

125    number of standard deviations that a point lies from the mean of a distribution. The

126    Mahalanobis distance of a $d$-dimensional observation $x_i = (x_{i1}, x_{i2}, \ldots, x_{id})^T$ from a

127    distribution of $N$ variables with mean $\bar{x} = (\overline{x_1}, \overline{x_2}, \ldots, \overline{x_d})^T$ and covariance matrix $S$ is

128    defined as follows (Mahalanobis 1936):

129    $$D_M(x_i) = \sqrt{(x_i - \bar{x})^T S^{-1} (x_i - \bar{x})} \ . \tag{1}$$

130    This distance differs from the ordinary Euclidean distance due to the correction for

131    covariance among observations, making it a better distance measure for genome scan

132    summary statistics because it does not assume that statistics are independent (i.e.,

133    Euclidean distance equals Mahalanobis distance when $S$ is a diagonal matrix). However,

134    this distance does make the assumption that points disperse smoothly from a single

135    multivariate centroid, and so it will tend to perform poorly when observations have a

136    complex or multimodal distribution.

137    *Harmonic mean distance*. In this context the "harmonic mean distance" of an

138    observation $x_i$ refers to the harmonic mean of the distances between this point and all

139    other points. The distance measure used here is the Euclidian distance normalized by

140    multiplying by the inverse covariance matrix. This ensures that results are not

141    dominated by a few statistics with a large spread, and also accounts for any potential

142    correlation between statistics, analogously to the Mahalonobis distance. Mathematically

143    we can define the harmonic mean distance as follows:

144 $$D_H(x_i) = N \left[ \sum_{j \neq i} \left[ (x_i - x_j) S^{-1} (x_i - x_j) \right]^{-1/2} \right]^{-1} . \tag{2}$$

145 The harmonic mean is heavily influenced by small values, which in this context means

146 local effects are amplified. However, more distant points also have some effect on the

147 final value (unlike the nearest neighbor distance described below), and so the harmonic

148 mean strikes a balance between local and global effects. This has some advantages in

149 outlier detection, as observations that are both distant from the main mass of the data

150 and have few neighbors in the local vicinity will tend to be outliers.

151 *Nearest neighbor distance.* The nearest neighbor distance of the observation $x_i$ gives the

152 minimum distance between this point and any other point. As with the harmonic mean

153 distance, we use the Euclidian distance normalized by the inverse covariance matrix.

154 Mathematically we can write

155 $$D_N(x_i) = \min_{j \neq i} \left( \sqrt{(x_i - x_j) S^{-1} (x_i - x_j)} \right) . \tag{3}$$

156 This statistic exclusively measures local effects, being largest when an observation is a

157 long way from any other point. Because this distance is only based on two points (the

158 focal point and its nearest neighbor), it is not influenced by the global distribution of the

159 data, unlike the harmonic mean distance.

160 *Kernel density deviance.* Kernel density-based methods attempt to capture

161 mathematically the distribution of the data as the sum of a number of simple parametric

162 distributions. Here we apply these methods to identifying multivariate outliers, defined

163 as those points with a low density of data around them in multivariate space. We

164 assume a multivariate normal kernel $G(x_i \mid x_j, \lambda^2 S)$ centered at the point $x_j$, where $\lambda$ is

165    the bandwidth of the kernel, which is scaled in each dimension by the covariance matrix

166    of the data. We then calculate the leave-one-out log-likelihood (Leiva-Murillo and Artés-

167    Rodríguez, 2012) of the point $x_i$ as follows:

168    $L(x_i \mid \lambda) = \log\left(\frac{1}{N-1}\Sigma_{j \neq i}\, G(x_i \mid x_j, \lambda^2 S)\right)$ .    (4)

169    In other words, this is equal to the log-probability density of the point $x_i$ under the

170    kernel density distribution constructed from all points *apart from $x_i$*. Our final density-

171    based measure is defined as follows:

172    $D_K(x_i) = -2L(x_i \mid \lambda)$ ,    (5)

173    which is sometimes referred to as the Bayesian deviance. This will be large whenever

174    the density of the point $x_i$ is low, and so the kernel density deviance can be thought of as

175    a measure of the sparseness of points around the focal point.

176    One challenge when using kernel density methods is choosing an appropriate value for

177    the bandwidth. Here we simply use the bandwidth for which the total deviance of all

178    points is minimized, i.e.

179    $\lambda^* = argmin_\lambda\left(\sum_{i=1}^{N} -2L(x_i \mid \lambda)\right)$ .    (6)

180    It can be shown that this is equivalent to the maximum-likelihood value of $\lambda$ under the

181    leave-one-out criterion. The value $\lambda^*$ can be found using the MINOTAUR function

182    `kernelDeviance()`, which takes a vector of bandwidths as input and returns the total

183    deviance of each. This function can be used to search for the minimum value of $\lambda$

184    manually, or via an optimization routine such as `optim()`. Users are also free to use any

185    other bandwidth, entered manually, or in the absence of a user-defined bandwidth a

186     simple method based on Silverman's rule is implemented as a default (this assumes that

187     data is normally distributed, and is a simple function of the standard deviation of the

188     samples (Silverman 1986)).

189     **The MINOTAUR R package - an R Shiny graphical user interface for multivariate**

190     **outlier analysis and visualization**

191     The MINOTAUR package performs two main functions: (1) it calculates the compound

192     multivariate outlier statistics described above and (2) it enables users to harness the

193     interactive graphical power of the R Shiny environment to manipulate and visualize

194     their data within the MINOTAUR graphical user interface (GUI). The GUI allows users to

195     perform the former task with the click of a button; however, outlier identification can

196     also be performed on the R command line using stand-alone functions available in

197     MINOTAUR, if preferred. Directions for downloading and installing the package can be

198     found at the end of this manuscript.

199     The MINOTAUR GUI is designed to streamline the process of genomic data analysis and

200     outlier identification, taking users from data input to graphical output within a single

201     platform. Distinct panels are used for each stage of the analysis, including data input

202     and filtering, outlier detection via the methods described above, and plotting results

203     (e.g., histograms, scatterplots, and Manhattan plots). An overview of the MINOTAUR GUI

204     workflow is show in Figure 1.

205     In the *Data* panel, the MINOTAUR GUI allows users to either upload their own datasets

206     or select among a set of four in-built example datasets. Data can be uploaded in a

207     number of file formats, including comma- or tab-separated text files, and Rdata.

208     Regardless of the file format, MINOTAUR expects all incoming datasets to be arranged in

209   data frames, with each row representing a different genetic locus and each column

210   representing a different univariate genome scan statistic (e.g., $F_{ST}$, Tajima's $D$, etc.) or

211   other piece of locus-specific metadata (e.g., SNP identifiers, chromosomes/scaffolds and

212   positions, etc.). Raw data objects can be filtered within the GUI, meaning, for example,

213   that columns not related to outlier analysis can be dropped at an early stage.

214   Four example datasets are made available to users within the MINOTAUR package and

215   GUI. The "HumanGWAS" dataset contains example output from an unpublished human

216   Genome-Wide Association Study. The simulated "NonParametricInverse" and

217   "NonParametricMultimodal" datasets each contain an example of nonparametric data,

218   one with an inverse relationship (Figure 3) and one that is highly multimodal

219   (Supplemental Figure S1). The "TwoRefSim" dataset contains population genetic data

220   simulated under a model of expansion from two refugia (Lotterhos and Whitlock 2015).

221   Note that the example datasets can also be accessed outside the GUI by running the

222   `data()` command with the appropriate dataset name. For example, to load the

223   "HumanGWAS" dataset, type `data(HumanGWAS)` and hit ENTER. To learn more about a

224   dataset while in the R terminal, add a question mark before the dataset name to load the

225   relevant Help page; for example, type `?HumanGWAS` and hit ENTER.

226   In the *Outlier Detection* panel, multiple univariate statistics can be integrated to produce

227   the compound distance measures described above. These measures can be appended to

228   the data frame and visualized interactively in the *Produce Plots* panel, which includes

229   several submenus with useful plots for visualizing high-dimensional datasets, including

230   Manhattan plots, 1D histograms and density-based 2D Scatterplots. The plotting

231   methods are designed with large genomic datasets in mind; for example the `plot2d()`

232   function included with the package calculates the density of points for a given bin size

233 and shades bins according to the density of points within them, and then optionally

234 adds user-supplied points (ideally a small subset of points, for example the outliers

235 only) to the plot. Additional options allow users to log-scale statistics and control

236 various other visual settings commonly used when plotting data in R (Figure 2).

**Example applications of multivariate outliers**

237

238 *Evaluation of computational speed*. First, we evaluated the speed of calculating the four

239 compound distance measures for datasets with increasing numbers of loci (rows) and

240 univariate statistics (columns). For this example, variables were randomly generated

241 from a multivariate normal distribution. Table 1 gives the "order" of complexity of these

242 algorithms, together with measured run-times for a dataset composed of 50,000 loci

243 and 10 variables (see Supplementary Table S2 for extended run-time analyses). Overall,

244 the Mahalanobis distance is calculated in a matter of seconds, even with particularly

245 large datasets. The harmonic mean distance, nearest neighbor distance, and kernel

246 density deviance each scale approximately equally with increasing dataset sizes, though

247 the maximum likelihood estimate of the ideal bandwidth for the latter measure can add

248 significant computation time.

249 *Example on simulated nonparametric distributions*. Some kinds of genomic data - for

250 example gene expression data - may generate complex nonparametric distributions.

251 Genes that have high expression in one environment may have low expression in

252 another environment, while investigators may be interested in identifying genes that

253 have moderate expression in both environments. To test the performance of the

254 multivariate outlier statistics in nonparametric situations, we simulated two examples

255 of nonparametric distributions.

256    In the first example, we simulated a distribution of two variables that follow an inverse

257    relationship, with some additional noise. We used contour plots to visualize the

258    different ways in which each of the compound distance measures changes over the two-

259    dimensional plane (Figure 3). In these plots, the darker red lines indicate less-

260    significant values of the test statistic and lighter yellow lines indicate more-significant

261    values of the test statistic. We also looked at two manually chosen points on the plane -

262    indicated by a blue square and triangle - chosen to represent different sorts of outliers.

263    The blue triangle would not be considered an outlier from the perspective of either one-

264    dimensional distribution despite being a clear outlier from the two-dimensional

265    distribution, while the blue square would be considered an outlier in the first dimension

266    but not the second. In this example, the nonparametric distribution affects the relative

267    ability of the four statistics to identify each of these outliers (Figure 4). The blue triangle

268    would not have the largest value (i.e., not be the most outlying point) by the

269    Mahalanobis or the harmonic mean distance, while it would have the largest value by

270    nearest neighbor distance or kernel density deviance. In contrast, the blue square has

271    the largest value of the test statistic by all four methods.

272    In the second example, we simulated a highly multimodal distribution from a normal

273    mixture model. In this example, it can be seen how the parametric assumption of the

274    Mahalanobis distance fails to capture the complexity of the data (Supplementary Figure

275    S1). In contrast to the previous example, the harmonic mean distance behaves similarly

276    to the kernel density deviance, and nearest neighbor distance has the most complex

277    contour landscape.

278    *Example on simulated genomic data*. To test the power of multivariate statistics for

279    genome scans, we applied them to a published simulated dataset that was used to test

280    different genome scan methods (Lotterhos and Whitlock 2014, 2015). Briefly, a

281    landscape simulator was used to simulate haploid neutral and selected loci that adapted

282    to an environmental cline (Lotterhos and Whitlock 2015). The landscape consisted of

283    360 x 360 demes and the allele frequency of each deme changed each generation

284    according to recurrence equations for mutation, migration, selection (if applicable), and

285    drift (Lotterhos and Whitlock 2015). For the dataset used in this example, a total of

286    9900 neutral and 100 selected loci (simulated under varying strengths of selection: 12

287    loci with s = 0.1, 38 loci with s = 0.01, and 50 loci with s = 0.005) were simulated under

288    a two-refuge demographic expansion. Individuals were then sampled from the

289    landscape according to the allele frequency in each deme at 30 randomly chosen

290    locations on the landscape at 20 individuals per location. For additional details see

291    Lotterhos and Whitlock (2014, 2015).

292    The simulated data were used to create a single nucleotide polymorphism (SNP) table

293    and this data was used to perform genome scans in the programs Bayenv2 (Günther

294    and Coop 2013) and LFMM (Frichot et al. 2013, now implemented in the R package LEA:

295    Frichot and François 2015). A total of four univariate statistics from these two

296    programs were used in the search for multivariate outliers: (i) log-Bayes Factor (log-BF,

297    a measure of the association between allele frequency and the environment in

298    Bayenv2), (ii) Spearman's rho (a measure of the association between allele frequency

299    and the environment in Bayenv2), (iii) $X^T X$ (a measure of genetic differentiation among

300    populations in Bayenv2), and (iv) $Z$-score (a measure of the association between

301    genotype and the environment in LFMM). These four univariate statistics, plotted in

302    Figure 4, were previously shown to have different strengths and weaknesses depending

303    on sampling design and demographic history (Lotterhos and Whitlock 2015).

304    To illustrate the flexibility of the outlier functions implemented in MINOTAUR, we

305    calculated multivariate outliers in two ways, corresponding to two different ways of

306    calculating the covariance matrix $S$ in equations (1) to (4). First, we used the traditional

307    method of calculating the covariance matrix based on all the data. For high-dimensional

308    data, estimation of the multivariate mean and covariance (location and scatter) are

309    expected to be robust to outliers as long as the proportion of outliers in the data is less

310    than $1/(k+1)$, where $k$ is the number variables in the dataframe (Ro et al. 2015).

311    However, we found that even in this small dataframe of only 4 variables and 10,000 loci,

312    the 1% of selected loci (a fraction of which were true outliers) affected the estimation of

313    the covariance matrix. For this reason, our MINOTAUR functions are designed to allow

314    the user to input their own covariance matrix. To illustrate this use of the function, we

315    also calculated a robust multivariate location and scatter estimate with a high

316    breakdown point, using the 'Fast MCD' (Minimum Covariance Determinant) estimator

317    with the function `CovNAMcd` in the R package `rrcovNA` (Rousseeuw et al 1999; Todorov

318    et al. 2011).

319    To compare the ability of the univariate statistics and the multivariate statistics to

320    separate neutral from selected loci, we calculated the empirical power. The empirical

321    power is based on using all known neutral loci to generate a null distribution, and then

322    for each locus an empirical $p$-value is calculated based on its cumulative frequency in

323    this null distribution. To control for false discovery rate, empirical $p$-values were

324    converted to $q$-values (in the R package `qvalue`: Dabney and Storey 2014) and loci with

325    a $q$-value less than 0.05 were retained as positive hits (a $q$-value of 0.05 has a desired

326    rate of 5 false positives out of 100 positive hits).

327    For the univariate statistics, the empirical power was highest for log-BF (0.54) and

328    lowest for $Z$-score (0.15), with Spearman's rho (0.46) and $X^T X$ (0.42) also showing

329    moderate power. For the multivariate statistics with the default covariance estimation,

330    the empirical power was high for harmonic mean distance and Mahalanobis distance

331    (0.41 for both), with kernel density and nearest neighbor distance performing poorly in

332    this case (0.09 for both) (Supplementary Figure S2). For the user-input covariance

333    matrix estimated with a high breakdown point (i.e., less influenced by outliers), the

334    empirical power was highest for harmonic mean distance and Mahalanobis distance

335    (0.58 for both), with kernel density and nearest neighbor distance still performing

336    poorly (Figure 5). This final example illustrates the potential of Mahalanobis and

337    harmonic mean distance to improve the signal-to-noise ratio in genome scans, because

338    the empirical power in this case was higher than any univariate statistic alone.


339    **Discussion**


340    Although the number of packages for population genetic data analysis in the R software

341    is rapidly increasing (http://popgen.nescent.org/PACKAGES.html), basic tools for

342    manipulating and visualizing genome-scale datasets have so far been lacking.

343    MINOTAUR fills this gap using the R Shiny Dashboard package to implement a GUI that

344    makes it easy to upload, manipulate, analyze, and visualize genomic data.


345    The multivariate metrics calculated in MINOTAUR contribute to a growing number of

346    multivariate tools implemented in the R environment (see Supplementary Table S1).

347    Methods that are influenced heavily by the distance of a point from the centroid in

348    multivariate space (such as Mahalanobis and the harmonic mean distance) will perform

349    differently compared with methods that are influenced mainly by the sparseness of

350    points in multivariate space (such as nearest neighbor distance and kernel density), as

351    illustrated in the examples here. However, depending on how the data are distributed,

352    the harmonic mean distance may be influenced by both these factors. For a single

353    simulated dataset, we found that robust use of the Mahalanobis or harmonic mean

354    distance (i.e., when the covariance matrix used was estimated with a high breakdown

355    point) could have higher power than any single univariate statistic alone. Although

356    nearest neighbor distance and kernel density deviance performed poorly on the

357    simulated genomic data, they may be useful in application to other kinds of

358    nonparametric data, as illustrated in our examples (Figures 3 and S1). Further

359    evaluation, however, will be needed on both simulated and empirical data to determine

360    whether multivariate outlier approaches will improve the signal-to-noise ratio in

361    genome scans.

362    The MINOTAUR package is designed to complement existing tools for the analysis and

363    integration of genome-scan data. Thus, in addition to providing its own tools for

364    genome-scale analyses, MINOTAUR can serve as a platform for the further analysis and

365    visualization of results generated by other R packages. Examples include results from

366    differential gene expression (LIMMA: Ritchie et al. 2015; DESeq: Anders and Huber

367    2010; SeqGSEA: Wang and Cairns 2014), outliers for genetic differentiation (OutFLANK:

368    Whitlock and Lotterhos 2015; PCAdapt: Luu and Blum 2015), genetic-environment

369    associations (LEA: Frichot and François 2015), or genome-wide association studies (e.g.

370    GenABEL: Aulchenko et al. 2007; BlueSNP: Huang et al. 2013).

371    Recent developments such as the R Shiny and Shiny Dashboard environments (Chang

372    2015; Chang et al. 2016) dramatically aid in the development of R-based user-friendly

373    web interfaces. Taking advantage of these tools, MINOTAUR is able to offer a new

374    platform for visualizing and integrating genomic data that may appeal to molecular

375    ecologists, modellers, statisticians, and public health agencies.


376    **Resources**


377    **Availability:** Upon acceptance for publication, MINOTAUR will be distributed on CRAN
378    (http://cran.r-project.org/) and be available for R on Windows, Mac OSX, and Linux
379    platforms. Currently, MINOTAUR can be accessed via the following steps:

380    • `install.packages("devtools", dependencies=TRUE)`
381    • `library(devtools)`
382    • `install_github("NESCent/MINOTAUR", build_vignettes=TRUE)`
383    • `library(MINOTAUR)`
384    • `MINOTAUR()`

385    Note to reviewers: If you are facing issues with installation, try updating to the newest
386    version of R and reinstalling devtools from source.  MINOTAUR has been tested on R
387    version 3.3.0.

388    **Licence:** GNU General Public Licence (GPL) >= 2.

389    **Documentation:** Besides the usual package documentation, MINOTAUR is released
390    with a tutorial which can be opened by typing: `vignette("MINOTAUR")`.

391    **Development:** The development of MINOTAUR is hosted on GitHub:
392    (https://github.com/NESCent/MINOTAUR).

393
394    **Acknowledgments**

395    The resource reported in this paper began at the Population Genetics in R Hackathon,

396    which was held in March 2015 at the National Evolutionary Synthesis Center (NESCent)

397    in Durham, NC, with the goal of addressing interoperability, scalability, and workflow

398    building challenges for the population genetics package ecosystem in R. The authors

399    were participants in the hackathon, and are indebted to NESCent (NSF #EF-0905606)

400    for hosting and supporting the event. RV, CC, DCC, SMS benefited from travel support to

401    attend this Hackathon.


402    **Author Contributions**

403 RV and KEL conceptualized the study. RV derived the compound outlier measures and

404 implemented them in R. RV and CC conceptualized and implemented the Shiny

405 Dashboard GUI. CC managed R package development, package structure and

406 documentation. All authors contributed R code for plotting in the Shiny Dashboard. DCC

407 wrote the vignette for the package and estimated function computational time. KEL

408 performed analysis of the performance of compound measures on simulated data. All

409 authors contributed to writing this manuscript.

## References

410

411 Allendorf FW, Hohenlohe PA, Luikart G (2010) Genomics and the future of conservation
412     genetics. *Nature Reviews Genetics*, **11**, 697–709.
413 Anders S, Huber W (2010) Differential expression analysis for sequence count data.
414     *Genome Biology*, **11**, R106
415 Aulchenko YS, Ripke S, Isaacs A, Duijn CMV (2007) GenABEL: an R library for genome-
416     wide association analysis. *Bioinformatics*, **23**, 1294–1296.
417 Brachi B, Morris GP, Borevitz JO (2011) Genome-wide association studies in plants: the
418     missing heritability is in the field. *Genome Biology*, **12**, 232.
419 Carlson CS, Eberle MA, Kruglyak L, Nickerson DA (2004) Mapping complex disease loci
420     in whole-genome association studies. *Nature*, **429**, 446–452.
421 Chang W (2015). shinydashboard: Create Dashboards with 'Shiny'. R package version
422     0.5.1. https://CRAN.R-project.org/package=shinydashboard
423 Chang W, Cheng J, Allaire JJ, Xie Y, McPherson J (2016). shiny: Web Application
424     Framework for R. R package version 0.13.1. https://CRAN.R-
425     project.org/package=shiny
426 Dabney A, Storey JD (2014) qvalue: Q-value estimation for false discovery rate control.
427     R package, version 1.38.0. https://github .com/jdstorey/qvalue.
428 De Mita S, Thuillet AC, Gay L *et al.* (2013) Detecting selection along environmental
429     gradients: analysis of eight methods and their effectiveness for outbreeding and
430     selfing populations. *Molecular Ecology*, **22**, 1383–1399.
431 De Villemereuil PD, Frichot É, Bazin É, François O, Gaggiotti OE (2014) Genome scan
432     methods against more complex models: when and how much should we trust
433     them? *Molecular Ecology*, **23**, 2006–2019.
434 Eddelbuettel D (2013) Seamless R and C++ Integration with Rcpp. Springer, New York.
435     ISBN 978-1-4614-6867-7.
436 Eddelbuettel D, François R (2011) Rcpp : Seamless R and C Integration. *Journal of*
437     *Statistical Software*, **40**, 1-18
438 Evangelou E, Ioannidis JP (2013) Meta-analysis methods for genome-wide association
439     studies and beyond. *Nature Reviews Genetics*, **14**, 379–389.
440 Frichot E, François O (2015) LEA : An R package for landscape and ecological
441     association studies. *Methods in Ecology and Evolution*, **6**, 925–929.

442  Frichot E, Schoville SD, Bouchard G, Francois O (2013) Testing for Associations between
443      Loci and Environmental Gradients Using Latent Factor Mixed Models. *Molecular*
444      *Biology and Evolution*, **30**, 1687–1699.
445  Funk WC, Mckay JK, Hohenlohe PA, Allendorf FW (2012) Harnessing genomics for
446      delineating conservation units. *Trends in Ecology & Evolution*, **27**, 489–496.
447  Galesloot TE, van Steen K, Kiemeney LALM, Janss LL, Vermeulen SH (2014). A
448      comparison of multivariate genome-wide association methods. PLoS ONE, **9**,
449      e95923.
450  Goddard ME, Hayes BJ (2009) Mapping genes for complex traits in domestic animals
451      and their use in breeding programmes. *Nature Reviews Genetics*, **10**, 381–391.
452  Hindorff LA, Sethupathy P, Junkins HA *et al.* (2009) Potential etiologic and functional
453      implications of genome-wide association loci for human diseases and traits.
454      *Proceedings of the National Academy of Sciences*, **106**, 9362–9367.
455  Hoban S, Kelley JK, Lotterhos KE, Bradburd G, Lowry DB, Poss ML, Reed LK, Storfer A,
456      Whitlock MC (2016) Finding the genetic basis of local adaptation: problems,
457      pitfalls, and future direction. *American Naturalist* (In revision)
458  Hohenlohe PA, Phillips PC, Cresko WA (2010) Using population genomics to detect
459      selection in natural populations: key concepts and methodological considerations.
460      *International Journal of Plant Sciences*, **171**, 1059–1071.
461  Huang H, Tata S, Prill RJ (2012) BlueSNP: R package for highly scalable genome-wide
462      association studies using Hadoop clusters. *Bioinformatics*, **29**, 135–136.
463  Kohn MH, Murphy WJ, Ostrander EA, Wayne RK (2006) Genomics and conservation
464      genetics. *Trends in Ecology & Evolution,* 21, 629-637.
465  Leiva-Murillo JM, Artés-Rodríguez A (2012) Algorithms for maximum-likelihood
466      bandwidth selection in kernel density estimators. *Pattern Recognition Letters*, **33**,
467      1717–1724.
468  Lotterhos KE, Francois O, Blum M (2016) Not just methods: User expertise explains the
469      variability of outcomes of genome-wide studies. *Molecular Ecology* (In review).
470  Lotterhos KE, Whitlock MC (2014) Evaluation of demographic history and neutral
471      parameterization on the performance of F ST outlier tests. *Molecular Ecology*, **23**,
472      2178–2192.
473  Lotterhos KE, Whitlock MC (2015) The relative power of genome scans to detect local
474      adaptation depends on sampling design and statistical method. *Molecular Ecology*,
475      **24**, 1031–1046.
476  Luu K, Blum MGB (2015). pcadapt: Principal Component Analysis for Outlier Detection.
477      R package version 2.1. https://CRAN.R-project.org/package=pcadapt
478  Mahalanobis PC (1936) On the generalised distance in statistics. *Proceedings of the*
479      *National Institute of Sciences of India,* 2, 49–55.
480  Mccarthy MI, Abecasis GR, Cardon LR *et al.* (2008) Genome-wide association studies for
481      complex traits: consensus, uncertainty and challenges. *Nature Reviews Genetics*, **9**,
482      356–369.
483  Mccarthy MI, Hirschhorn JN (2008) Genome-wide association studies: past, present and
484      future. *Human Molecular Genetics*, **17**, R100-101.
485  O'Reilly, P. F. *et al.* MultiPhen: joint model of multiple phenotypes can increase
486      discovery in GWAS. *PLoS One* **7,** e34861 (2012).
487  Patterson, N., Price, A. L. & Reich, D. Population structure and eigenanalysis. *PLoS Genet.*
488      **2,** e190 (2006).
489  R Core Team (2015) R: A language and environment for statistical computing. R
490      Foundation for Statistical Computing, Vienna, Austria. https://www.R-project.org/.

491 Rellstab C, Gugerli F, Eckert AJ, Hancock AM, Holderegger R (2015) A practical guide to
492      environmental association analysis in landscape genomics. *Molecular Ecology*, **24**,
493      4348–4370.
494 Ritchie ME, Phipson B, Wu D *et al.* (2015) limma powers differential expression
495      analyses for RNA-sequencing and microarray studies. *Nucleic Acids Research*, **43**, 1-
496      13.
497 Rousseeuw PJ, Driessen KV (1999) A Fast Algorithm for the Minimum Covariance
498      Determinant Estimator. *Technometrics*, **41**, 212–223.
499 Savolainen O, Lascoux M, Merilä J (2013) Ecological genomics of local adaptation.
500      *Nature Reviews Genetics*, **14**, 807–820.
501 Silverman BW (1986) *Density estimation for statistics and data analysis*. Chapman and
502      Hall, London.
503 Tiffin P, Ross-Ibarra J (2014) Advances and limits of using population genetics to
504      understand local adaptation. *Trends in Ecology & Evolution*, **29**, 673–680.
505 Todorov V, Templ M, Filzmoser P (2011) Detection of multivariate outliers in business
506      survey data with incomplete information. *Advances in Data Analysis and*
507      *Classification*, **5**, 37–56.
508 Varshney RK, Nayak SN, May GD, Jackson SA (2009) Next-generation sequencing
509      technologies and their implications for crop genetics and breeding. *Trends in*
510      *Biotechnology*, **27**, 522–530.
511 Vatsiou AI, Bazin E, Gaggiotti OE (2016) Detection of selective sweeps in structured
512      populations: a comparison of recent methods. *Molecular Ecology*, **25**, 89–103.
513 Verity R, Nichols RA (2014) What is genetic differentiation, and how should we measure
514      it- $G$ST, D, neither or both? *Molecular Ecology*, **23**, 4216–4225.
515 Wang X, Cairns MJ (2014) SeqGSEA: a Bioconductor package for gene set enrichment
516      analysis of RNA-Seq data integrating differential expression and splicing.
517      *Bioinformatics*, **30**, 1777–1779.
518 Wang Z, Gerstein M, Snyder M (2009) RNA-Seq: a revolutionary tool for transcriptomics.
519      *Nature Reviews Genetics*, **10**, 57–63.
520 Whitlock MC, Lotterhos KE (2015) Reliable Detection of Loci Responsible for Local
521      Adaptation: Inference of a Null Model through Trimming the Distribution of Fst.
522      *The American Naturalist*, **186**, S24–36.
523 Yang J, Benyamin B, Mcevoy BP *et al.* (2010) Common SNPs explain a large proportion
524      of the heritability for human height. *Nature Genetics*, **42**, 565–569.

525 **Tables**

526 **Table 1.** Multivariate outlier detection methods implemented in MINOTAUR and
527 associated computational run times. Computational complexity is given in "big O"
528 notation, with $N$ referring to the number of observations and $k$ the number of statistics
529 (dimensions). Run times were determined using an Apple iMac with a 3.1 GHz Intel
530 Core i5 processor and 32 GB of RAM running Apple OSX 10.9.5 and R version 3.2.3. Note
531 that for computation time the kernel density deviance includes both the maximum
532 likelihood estimation of the optimal bandwidth and the density calculations based on
533 the optimal bandwidth.

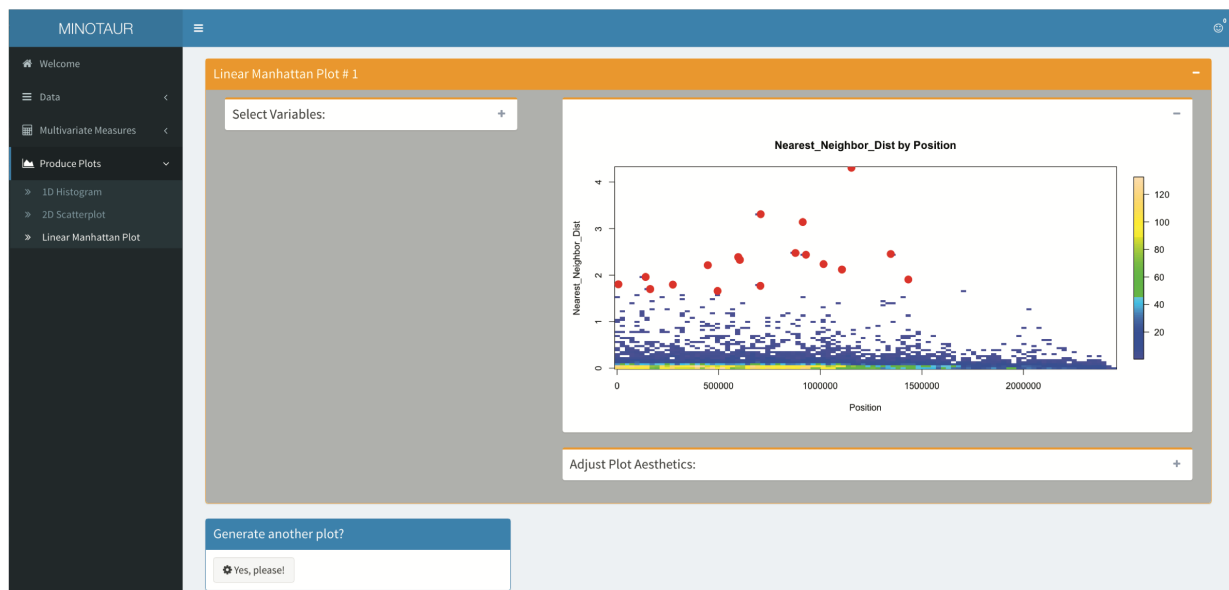| Compound measure | Description | R Function | Computational complexity (big O notation) | Computation Elapsed Time for 50,000 loci & 10 variables (hh:mm:ss.ms) |
|---|---|---|---|---|
| Mahalanobis distance | Distance from multivariate centroid | Mahalanobis() | $O(Nk^2)$ | 00:00:00.095 |
| Harmonic mean distance | Inverse-weighted distance from all other points | harmonicDist() | $O(Nk^2)$ | 00:04:13.620 |
| Kernel density deviance | Local density of points | kernelDist() | $O(Nk^2)$ | 01:40:03.600 |
| Nearest neighbor distance | Distance to nearest neighbor | neighborDist() | $O(Nk^2)$ | 00:04:07.020 |

534

535

536 **Figures**

537 **Figure 1.** Graphical overview of the MINOTAUR GUI workflow.
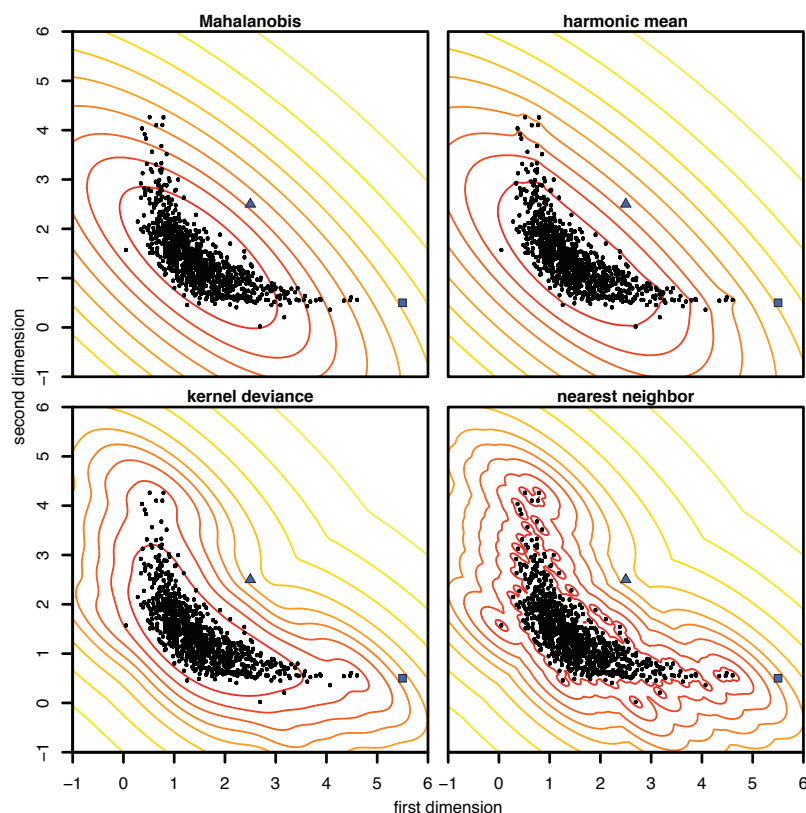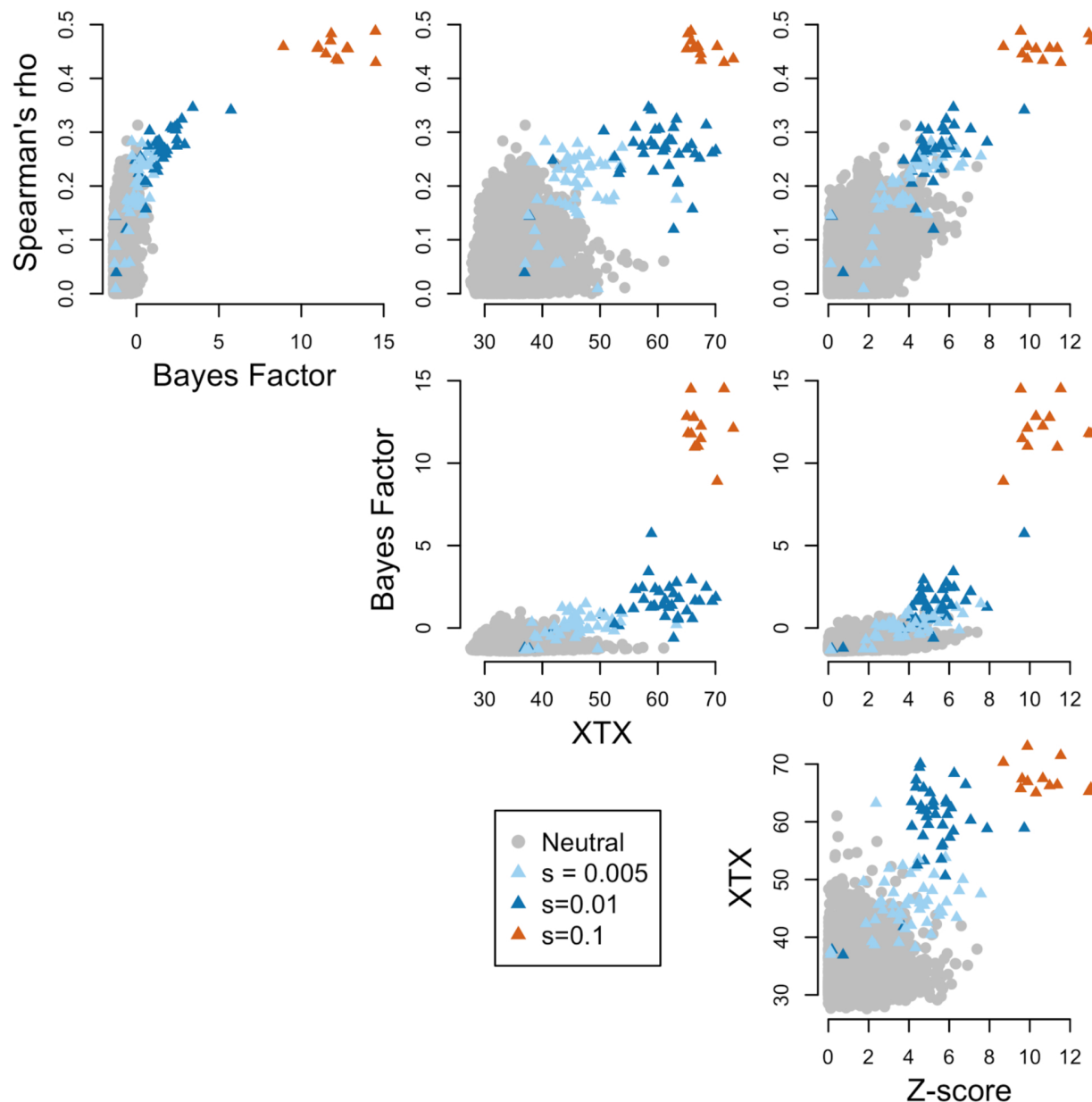


538

539

540 **Figure 2.** Screenshot of MINOTAUR GUI highlighting the overall interface and the ability
541 to visualize multivariate distributions. The plot is a Manhattan plot of the nearest
542 neighbor distance across loci for all traits in the "HumanGWAS" example dataset
543 provided as part of MINOTAUR. The base scatter plot demonstrates the binned
544 visualization, where the density of data in an area is apparent from the color. 99.5%
545 percentile outliers are indicated with solid orange circles. Visualization menus have
546 been collapsed to simplify the image. Additional plots can also be stacked below to
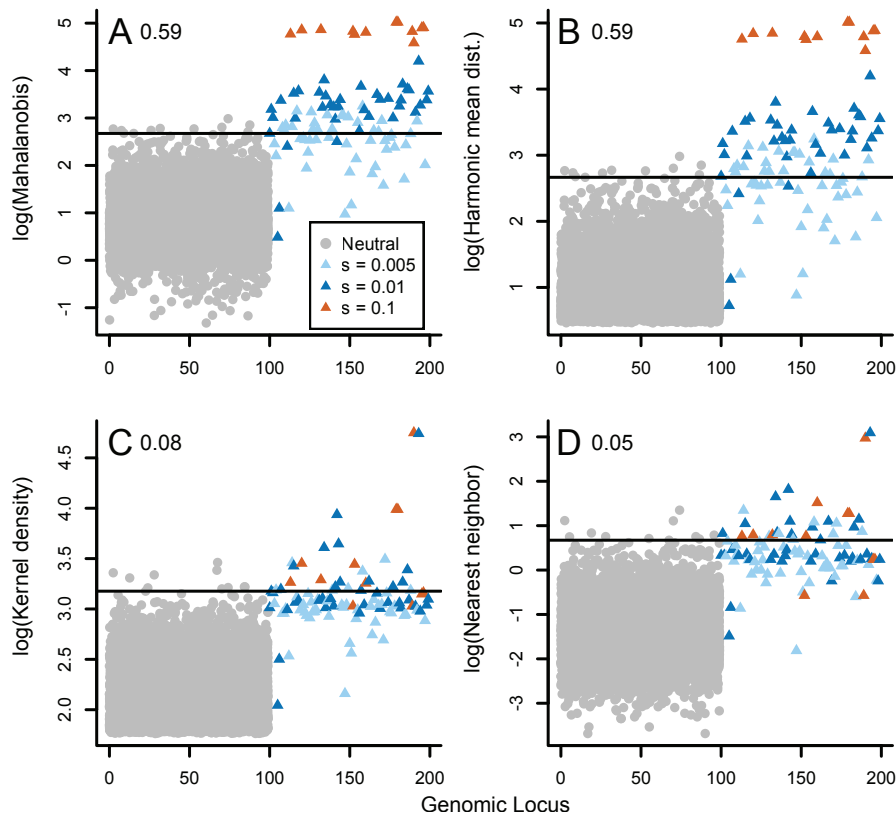547 enable comparisons across multiple plots (not shown).

**Figure 3.** Comparison of multivariate distance measures for nonparametric example data. Black dots show the simulated data, in which the two statistics (dimensions) are assumed to follow an inverse relationship with some additional noise. Solid lines show the distance measure computed at each point in the plane, arranged in 10% quantiles (e.g. the inner ring shows the 10% of locations with the smallest distance). The blue square and triangle show particular outlier points referred to in the main text.

557

**Figure 4.** Distributions of four univariate statistics from the two refuge dataset from Lotterhos and Whitlock (2015).



560

561

**Figure 5.** Distributions of the four multivariate compound statistics applied to the four univariate statistics shown in Figure 2. The MCD calculation of the covariance matrix was used. All 9900 neutral loci are plotted on indexes 0-100, and the selected loci are plotted on indexes 100-200. Note log transformation of each variable on the y-axis for: A) Mahalanobis distance, B) Harmonic mean distance, C) Kernel density, and D) Nearest Neighbor distance. The empirical power of the statistic to discriminate neutral from selected loci (see main text for details) is shown in the upper left hand corner.

569

## Supplementary Material for the Paper *MINOTAUR: A platform for the analysis and visualization of multivariate results from genome scans with R Shiny*

**Table S1.** Table of multivariate outlier statistics in other R packages that could be used in the context of genomic scans.

**Supplementary Table S1.** Table of multivariate outlier statistics available in other R packages

| Multivariate Outlier Statistic | R Package | R Function | Brief Description | Reference |
|---|---|---|---|---|
| Hierarchical Clustering Ranks | DMwR | outliers.ranking() | Uses an agglomerative hierarchical clustering algorithm to rank outlierness. | Torgo 2011 |
| Projection Congruent Subset | FastPCS | FastPCS() | Computes fast and robust multivariate outlyingness index. | Vakili & Schmitt 2014 |
| Kernel Density Estimator | ks | kde() | Computes the kernel density estimate for up to 6 dimensional datasets. | Duong 2007 |
| Mahalanobis Distance | mvoutlier | locoutNeighbor() | Computes global and pairwise Mahalanobis distances for outlier visualization with number of neighbors varying and fraction of neighbors fixed. | Filzmoser & Gschwandtner 2015 |
| Mahalanobis Distance | mvoutlier | locoutSort() | Computes global and pairwise Mahalanobis distances for interactive outlier visualization. | ' ' |
| Mahalanobis Distance | mvoutlier | locoutPercent() | Computes global and pairwise Mahalanobis distances for outlier visualization with number of neighbors fixed and varying fraction of neighbors. | ' ' |
| Principal Components Distance | mvoutlier | pcout() | Principal components distances are used to identify weighted location and scatter of outliers. | ' ' |
| Mahalanobis Distance | mvoutlier | sign1() | Principal components are used to calculate Mahalanobis distance covariance matrix and a critical value cutoff is used to determine outliers from chi-squared distribution. | ' ' |
| Principal Components Distance | mvoutlier | sign2() | Principal components distances are computed and transformed to approach a chi-squared distribution and a critical value cutoff is used to detect outlier. | ' ' |
| Adjusted Mahalanobis Distance | mvoutlier | arw() | Adjusts outlier rejection thresholds by using an adaptive reweighting estimator and determines outliers by the supremum of the difference between Mahalanobis distance and the theoretical distribution function. | ' ' |

**Supplementary Table S1.** Table of multivariate outlier statistics available in other R packages

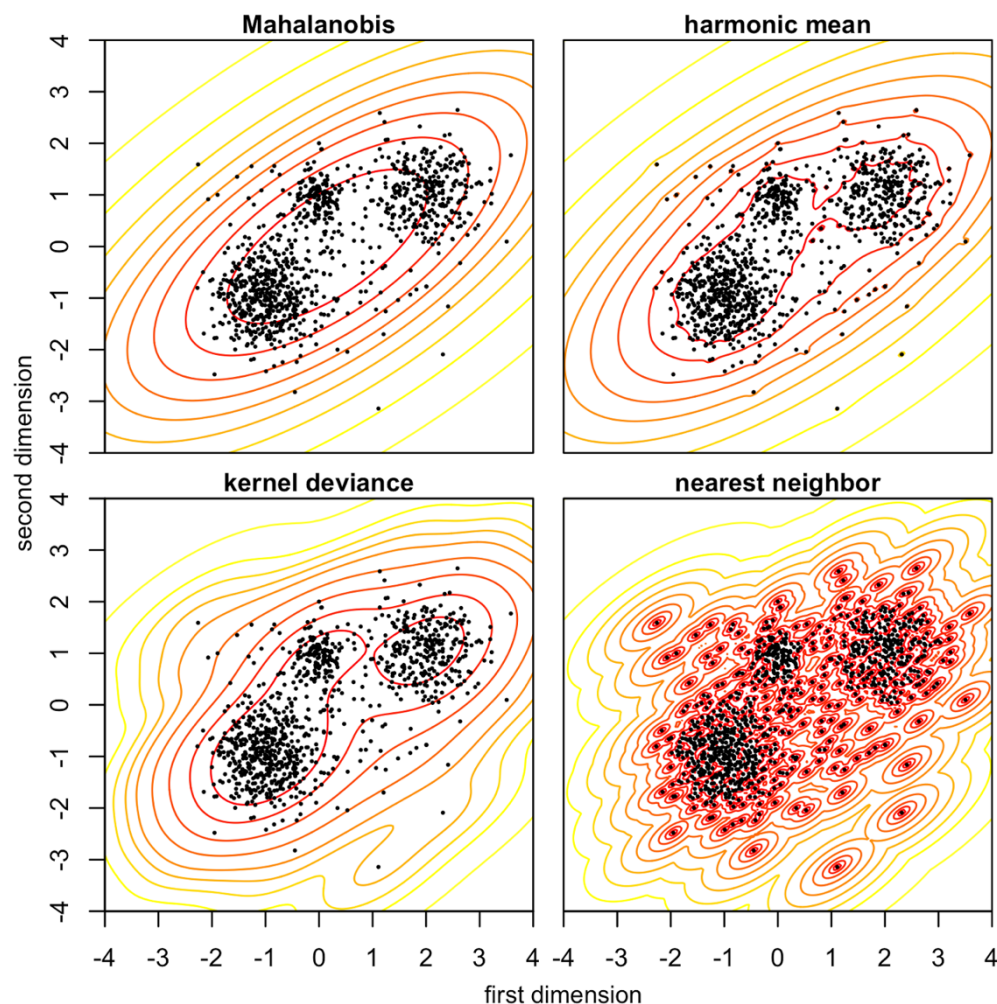| Multivariate Outlier Statistic | R Package | R Function | Brief Description | Reference |
|---|---|---|---|---|
| Mahalanobis Distance | rrcovHD | OutlierMahdist() | Calculates Mahalanobis distance and determines outliers based on a critical value of the chi-squared distribution. | Todorov 2016 |
| Mahalanobis Distance | CerioliOutlierDetection | cerioli2010.fsrmcd.test() | Calculates Mahalanobis distance based on the finite-sample reweighted Minimum Covariance Determinant (MCD) dispersion estimate. | Cerioli 2010 |
| Mahalanobis Distance | CerioliOutlierDetection | cerioli2010.irmcd.test() | Calculates Mahalanobis distances based on an iterated reweighted MCD dispersion estimate. | '' |
| Mahalanobis Distance & Adjusted Mahalanobis Distance | MVN | mvOutlier() | Calculates Mahalanobis distance or adjusted Mahalanobis distance and determines outliers based on the 97.5 percent quantile critical value of the chi-square distribution. | Korkmaz, Goksuluk & Zararsiz 2015 |

**References:**

Cerioli A (2010) Multivariate Outlier Detection With High-Breakdown Estimators. *Journal of the American Statistical Association*, **105**, 147–156.

Duong T (2007) ks: Kernel Density Estimation and Kernel Discriminant Analysis for Multivariate Data in R. *Journal of Statistical Software*, **21**, 1–16.

Filzmoser P, Gschwandtner M (2015) Mvoutlier: Multivariate Outlier Detection Based on Robust Methods. R package version 2.0.6. https://cran.r-project.org/web/packages/mvoutlier/mvoutlier.pdf

Korkmaz S, Goksuluk D, Zararsiz G (2014) MVN: an R package for assessing multivariate normality. *R Journal*, **6**, 151-162.

Todorov V (2016) Robust Multivariate Methods for high Dimensional Data. R package version 0.2-4. https://cran.r-project.org/web/packages/rrcovHD/rrcovHD.pdf

Torgo L (2011) *Data mining with R: learning with case studies*. Chapman & Hall/CRC, Boca Raton.

Vakili K, Schmitt E (2014) Finding multivariate outliers with FastPCS. *Computational Statistics & Data Analysis*, **69**, 54–66.

**Table S2.** Computation times for the four multivariate outlier detection methods in MINOTAUR for datasets up to 100,000 loci (rows) and 20 variables (columns) in hh:mm:ss.ms format. Run times were determined using an Apple iMac with a 3.1 GHz Intel Core i5 processor and 32 GB of RAM running Apple OSX 10.9.5 and R version 3.2.3. Note that the kernel density deviance includes both the maximum likelihood estimation of the optimal bandwidth and the density calculations based on the optimal bandwidth.

| No. Loci | No. Variables | Mahalanobis distance | Harmonic mean distance | Kernel density deviance | Nearest neighbor distance |
|---|---|---|---|---|---|
| 1000 | 5 | 00:00:00.001 | 00:00:00.040 | 00:00:01.233 | 00:00:00.034 |
| 1000 | 10 | 00:00:00.002 | 00:00:00.098 | 00:00:02.366 | 00:00:00.094 |
| 1000 | 15 | 00:00:00.003 | 00:00:00.188 | 00:00:04.303 | 00:00:00.185 |
| 1000 | 20 | 00:00:00.005 | 00:00:00.318 | 00:00:07.548 | 00:00:00.317 |
| 5000 | 5 | 00:00:00.003 | 00:00:00.986 | 00:00:30.215 | 00:00:00.829 |
| 5000 | 10 | 00:00:00.008 | 00:00:02.431 | 00:00:57.794 | 00:00:02.382 |
| 5000 | 15 | 00:00:00.014 | 00:00:04.809 | 00:01:49.900 | 00:00:04.622 |
| 5000 | 20 | 00:00:00.024 | 00:00:07.968 | 00:02:46.393 | 00:00:07.821 |
| 10000 | 5 | 00:00:00.006 | 00:00:03.922 | 00:02:00.694 | 00:00:03.328 |
| 10000 | 10 | 00:00:00.017 | 00:00:09.728 | 00:03:52.081 | 00:00:09.310 |
| 10000 | 15 | 00:00:00.169 | 00:00:18.773 | 00:06:53.011 | 00:00:18.487 |
| 10000 | 20 | 00:00:00.041 | 00:00:32.415 | 00:11:10.482 | 00:00:31.735 |
| 50000 | 5 | 00:00:00.045 | 00:01:37.808 | 00:50:19.078 | 00:01:24.120 |
| 50000 | 10 | 00:00:00.095 | 00:04:13.621 | 01:40:03.603 | 00:04:07.027 |
| 50000 | 15 | 00:00:00.161 | 00:09:12.221 | 00:19:36.847 | 00:08:51.240 |
| 50000 | 20 | 00:00:00.240 | 00:16:38.589 | 05:45:21.387 | 00:16:12.965 |
| 100000 | 5 | 00:00:00.085 | 00:06:35.265 | 03:21:34.758 | 00:05:35:996 |
| 100000 | 10 | 00:00:00.184 | 00:17:01.708 | 09:28:30.378 | 00:16:24:535 |
| 100000 | 15 | 00:00:00.324 | 00:36:19.473 | 13:28:52.158 | 00:37:20.698 |
| 100000 | 20 | 00:00:00.487 | 01:14:47.783 | 24:45:25.372 | 01:14:24.686 |

**Figure S1.** Comparison of multivariate distance measures for multi-modal example data. Black dots show the simulated data, drawn from a bivariate normal mixture model. Solid lines show the distance measure computed at each point in the plane, arranged in 10% quantiles, equivalently to Figure 3 in the main paper.

**Figure S2.** Analogue to Figure 5 in the main paper, but with a default estimate of covariance using all the data.