

covRNA - Discovering covariate associations in large-scale gene expression data

Lara H Urban

European Molecular Biology Laboratory

Christian W Remmele

Department of Bioinformatics, University of Würzburg

Marcus Dittrich

Human Genetik, Julius-Maximilians-Universität Würzburg

Roland F Schwarz

Berlin Institute for Medical Systems Biology, MDC

Tobias Müller (✉ Tobias.Mueller@biozentrum.uni-wuerzburg.de)

<https://orcid.org/0000-0003-3035-7070>

Research note

Keywords: multivariate analysis, fourthcorner analysis, RLQ analysis, transcriptomics, phenotype-44 genotype associations, high-throughput data, visualization, ordination methods, RNA-Seq 45 analysis, microarray analysis

Posted Date: November 22nd, 2019

DOI: <https://doi.org/10.21203/rs.2.17618/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at BMC Research Notes on February 24th, 2020. See the published version at <https://doi.org/10.1186/s13104-020-04946-1>.

1 **covRNA - Discovering covariate associations in large-**
2 **scale gene expression data**

3

4 Lara Urban^{1,2}, Christian W Remmele¹, Marcus Dittrich^{1,3}, Roland F Schwarz⁴, Tobias Müller^{1,*}

5 ¹ Department of Bioinformatics, Biocenter, University of Würzburg, Am Hubland, Würzburg,
6 Germany

7 ² European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome
8 Genome Campus, Hinxton, Cambridge, United Kingdom

9 ³ Institute of Human Genetics, University of Würzburg, Am Hubland, Würzburg, Germany

10 ⁴ Berlin Institute for Medical Systems Biology, Max Delbrück Center, Berlin, Germany

11

12 Email addresses:

13 lara.h.urban@ebi.ac.uk

14 christian.remmele@uni-wuerzburg.de

15 marcus.dittrich@ uni-wuerzburg.de

16 roland.schwarz@mdc-berlin.de

17 tobias.mueller@biozentrum.uni-wuerzburg.de

18

19 * Corresponding author: tobias.mueller@biozentrum.uni-wuerzburg.de

20

21 **Abstract**

22 **Objective**

23 The biological interpretation of gene expression measurements is a challenging task. While
24 ordination methods are routinely used to identify clusters of samples or co-expressed genes,
25 these methods do not take sample or gene annotations into account. We aim to provide a tool
26 that allows users of all backgrounds to assess and visualize the intrinsic correlation structure
27 of complex annotated gene expression data and discover the covariates that jointly affect
28 expression patterns.

29 **Results**

30 The Bioconductor package covRNA provides a convenient and fast interface for testing and
31 visualizing complex relationships between sample and gene covariates mediated by gene
32 expression data in an entirely unsupervised setting. The relationships between sample and
33 gene covariates are tested by statistical permutation tests and visualized by ordination. The
34 methods are inspired by the fourthcorner and RLQ analyses used in ecological research for the
35 analysis of species abundance data, that we modified to make them suitable for the
36 distributional characteristics of both, RNA-Seq read counts and microarray intensities, and to
37 provide a high-performance parallelized implementation for the analysis of large-scale gene
38 expression data on multi-core computational systems. CovRNA provides additional modules
39 for unsupervised gene filtering and plotting functions to ensure a smooth and coherent analysis
40 workflow.

41 **Keywords**

42
43 multivariate analysis, fourthcorner analysis, RLQ analysis, transcriptomics, phenotype-
44 genotype associations, high-throughput data, visualization, ordination methods, RNA-Seq
45 analysis, microarray analysis
46

47 **Introduction**

48
49 The biological interpretation of gene expression measurements and related multivariate
50 datasets is a fundamental yet challenging task in computational biology. Ordination methods
51 like Principal Component Analysis or Correspondence Analysis are routinely used for

52 dimension reduction and visualization to identify clusters of samples or co-expressed genes
53 [2]. These methods do not generally take sample or gene annotations into account. Knowledge-
54 driven approaches such as Gene Ontology Analysis [3] and Gene Set Enrichment Analysis [4]
55 look for differentially regulated sets of genes based on prior information. These methods are
56 powerful but specialized hypothesis-based tools. In functional genomics, it is often desirable to
57 test for associations between extensive categorical and numerical sample and gene covariates.
58 Sample covariates may comprise demographic and clinical data or complex phenotype data
59 derived from imaging. Gene-level covariates often include functional ontology, epigenetic
60 modifications, protein phosphorylation or copy-number state. Methods for the efficient and
61 systematic analysis of the relationship between sample and gene covariates mediated by gene
62 expression are lacking.

63 **Main Text** 64

65 Here we present covRNA, a Bioconductor package [5] providing a convenient and fast interface
66 for testing and visualizing the relationship between sample and gene covariates mediated by
67 gene expression in an entirely unsupervised setting. The methods are inspired by the
68 fourthcorner and RLQ analyses used in ecological research for the analysis of species
69 abundance data [6, 7]. While the scope of these analyses is comparable to knowledge-based
70 approaches like GSEA, their inherently unsupervised and hypothesis-free nature provides a
71 huge advantage if no prior knowledge is available. In addition, while approaches like GSEA are
72 based on parametric distributions like the hypergeometric distribution, the here presented
73 analyses are based on simulated distributions to capture and account for respective dataset-
74 specific data structures and modalities.

75 The RLQ analysis of the ade4 package [7] has previously been applied for the analysis of
76 microarray data describing the time-course effect of steroids on the growth of human lung
77 fibroblasts [8]. Within the covRNA package, we have modified the fourthcorner and RLQ
78 algorithms to make the methods inherently suitable for the distributional characteristics of both
79 RNA-Seq read counts and microarray intensities. We provide a parallelized high-performance
80 implementation to make the method suitable for the analysis of large-scale multivariate gene
81 expression data on multi-core computational systems, with additional modules for unsupervised

82 gene filtering and plotting functions to ensure a smooth and coherent analysis workflow. Here,
83 we demonstrate the analysis of a microarray dataset of the immune response of human
84 dendritic cells to fungal infection [9]. In addition, in order to show the applicability of our
85 approach to a more complex RNA-Seq data, a detailed vignette integrated in our Bioconductor
86 package [1] demonstrates the analysis of a well-established RNA-Seq dataset of *Bacillus*
87 *anthracis* [10].

88 **Methods**

89
90 covRNA takes as input three data frames: (i) a n times m gene expression data frame L of n
91 genes in m samples, (ii) a m times p sample annotation data frame Q of p sample covariates
92 for m samples and (iii) a n times s gene annotation data frame R of s gene covariates for n
93 genes. covRNA then performs a test for association between each sample and gene covariate
94 pair following the fourthcorner procedure. Data frames R , L and Q are multiplied to yield the s
95 times p test data frame $T = R'LQ$, where T_{ij} reduces to a pairwise Pearson correlation
96 coefficients weighted by the gene expression values of L . If both variables of a covariate pair
97 (i,j) are categorical, the entry T_{ij} is normalized by the sum over L to yield a Chi²-statistic. These
98 statistics potentially follow a non-normal and non-symmetric distribution due to the distributional
99 characteristics of RNA-Seq data. To account for violations of normality and symmetry, covRNA
100 uses a permutation test to calculate two-sided p-values and makes use of Fisher's assumption
101 of doubling the one-sided p-value in non-symmetric distributions [11]. We then use permutation
102 of the data frames to test for significant association between the covariates of R and Q .
103 Specifically, we adopt the permutation scheme according to Ter Braak *et al.* (2012) [12] to
104 ensure that all associations between gene and samples covariates are perturbed: First, the
105 rows of L are permuted and p-values p_1 between all covariates of R and Q are calculated. Then,
106 the columns of L are permuted and p-values p_2 between all covariates of R and Q are
107 calculated. After false discovery rate correction according to Benjamini and Hochberg [13] of
108 p_1 and p_2 , respectively, the actual p-values are obtained by $p = \max(p_1, p_2)$ [12]. Taking the
109 most conservative p-values hereby assures to model dependencies between samples and
110 genes correctly.

111 The high-performance implementation of this statistical analysis in covRNA allows for
 112 straightforward parallelization on multiple available cores and significant speed-up of the
 113 analysis of large-scale datasets (Table 1).

Permutations	10 ³	10 ⁴	10 ⁵	10 ⁶	10 ⁷
1 Core	9.1	52.9	5.3 × 10 ²	6.8 × 10 ³	6.9 × 10 ⁴
10 Cores	8.5	15.7	84.7	7.8 × 10 ²	7.7 × 10 ³
Speed-up	1.1	3.4	6.3	8.2	9.0

114

115 **Table 1.** Speed-up of the fourthcorner analysis implemented in covRNA due to parallelization
 116 across multiple cores. The fourthcorner analysis is performed on the Bacillus anthracis example
 117 dataset on 1 and 10 cores for different numbers of permutations as indicated in the first row.
 118 The next rows indicate the required user time in seconds while the last row indicates the speed-
 119 up of the multi-threading approach.

120

121 To visualize the relationship within and between sample and gene covariates we perform
 122 singular value decomposition on T, following the standard RLQ approach. This creates two-
 123 dimensional ordinations for both, sample and gene covariates, which are then combined into a
 124 joint ordination plot. In this plot, the covariates that are significantly associated with each other
 125 according to the statistical tests are connected by lines whose color reflect the type of the
 126 association (positive or negative).

127 Results

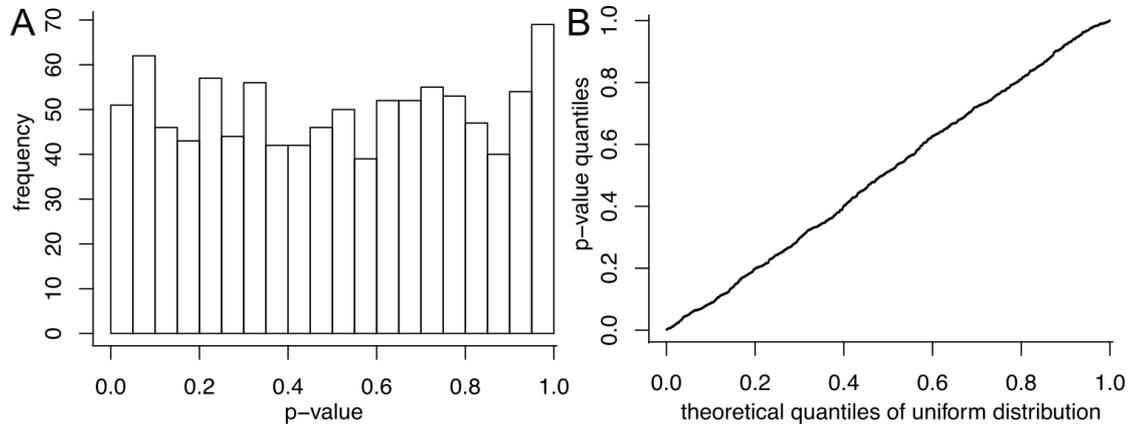
128

129 We applied our method to a microarray dataset of the immune response of human dendritic
 130 cells to *Aspergillus fumigatus* infection (GSE69723, GSE77969) [8]. The ExpressionSet Expr
 131 contains gene expression data under different stimuli ('control', 'LPS', 'A. fumigatus') and at
 132 different time points ('6h', '12h'). The genes are annotated by immune-related hallmark gene
 133 sets (n=7 gene sets) of the MSigDB collection [4].

134 We firstly tested if our statistical analyses were calibrated. We therefore chose an association
 135 between sample and gene annotations, and randomly permuted the gene annotation labels

136 n=1,000 times. The resulting p-values were uniformly distributed, affirming calibration of the
137 statistical tests (**Figure 1** for one sample annotation-gene annotation association).

138



139

140

141 **Figure 1.** covRNA's statistical test is shown to control the type I error rate correctly. A p-value
142 distribution under the null hypothesis of covRNA's statistical test between sample and gene
143 annotations for n=1,000 permutations is generated. The results of the permutation of one
144 random sample annotation-gene annotation association are shown here. **A.** Histogram of the
145 resulting p-values. **B.** Q-Q plot of the p-values.

146

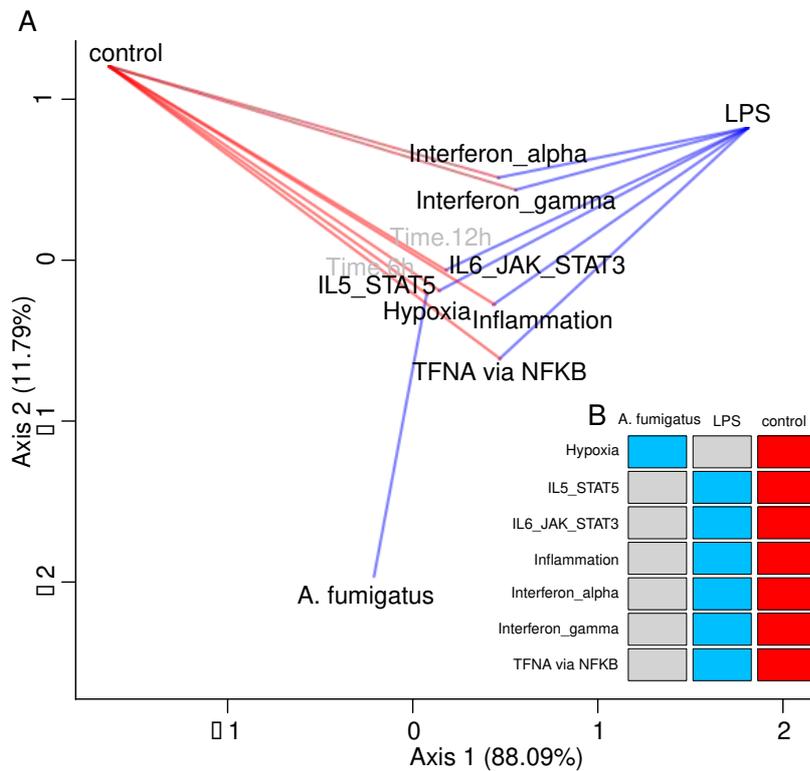
147 Having established the calibration of covRNA's statistical tests, we applied the covRNA
148 methods to the microarray dataset of *Aspergillus fumigatus* infections. The following R code
149 applied to the ExpressionSet Expr produces the results shown in **Figure 2.**

150

```
151 statobj <- stat(Expr)    # statistical tests  
152 ordobj <- ord(Expr)     # ordination parameters  
153 vis(statobj, ordobj)    # visualization (Figure 1A)  
154 plot(statobj)           # visualization of tests (Figure 1B)
```

155

156



157 **Figure 2.** Visualization of covRNA analysis of microarray data of human dendritic cells infected
158 with *A. fumigatus*. **A.** Ordination of sample and gene covariates. The lines between the
159 covariates denote significant negative (red) and positive (blue) associations (at a significance
160 level $\alpha=0.05$). Gray covariates are not involved in any significant association. **B.** Results of the
161 association test. Consistently, red, blue and gray colors denote significant negative, positive or
162 no significant associations (at a significance level $\alpha=0.05$).

163

164 **Figure 2** illustrates the concordance of both analysis approaches. Non-associated covariates,
165 here the two time points (6h, 12h) cluster around the origin of the ordination while
166 positively/negatively associated covariates are situated at different angles from the origin (at a
167 significance level $\alpha=0.05$; **Figure 2A**). The significant associations are also summarized in a
168 table (here $n=14$ significant associations; **Figure 2B**). This combined statistical and
169 visualization analysis allows researchers to obtain a quick overview of regulatory patterns in
170 their gene expression experiment: Here, the overview plot shows that the LPS infection of
171 dendritic cells elicits typical bacterial infection responses like interferon activation, while a
172 fungal infection by *A. fumigatus* leads to hypoxia in the cells. This overview confirms the

173 successful infection of the dendritic cells in the experiment, and allows for building first
174 hypotheses about the different molecular responses between bacterial and fungal infections.

175 **Discussion**

176

177 The Bioconductor package covRNA provides a coherent workflow to systematically test for and
178 visualize associations between sample and gene covariates mediated by gene expression.

179 With only a few lines of R code, users can assess and visualize the intrinsic correlation structure
180 of complex annotation data and discover the covariates that jointly affect the gene expression
181 patterns. Further, experimental biologists are provided with a quick tool to validate their
182 experiments, e.g. to assess if their stimulation assays have been successful.

183 The adaptation of the fourthcorner and RLQ methods, which are frequently applied in ecological
184 landscape analyses, to the distributional characteristics of gene expression data makes the
185 analyses accessible to a wider community. The efficient implementation and parallelization on
186 multiple cores further allows for the analysis and visualization of large-scale multivariate gene
187 expression datasets.

188 **Limitations**

189

190 While one of the benefits of the covRNA package is the efficient implementation that allows
191 scaling analyses up to thousands of genes, the analysis of too many gene and sample
192 annotations will lead to an unclear ordination visualization with too many annotations
193 overlapping each other. In such a case, we recommend to firstly consider the data frame
194 visualization, to then select interesting annotations for visualization.

195 While covRNA tests the statistical association of annotations, it does not include a test of
196 causality of associations. Instead, it provides a first insight into the internal structure of gene
197 expression data.

198 **Declarations**

199 **Ethics approval and consent to participate**

200

201 Not applicable

202 **Consent for publication**

203

204 Not applicable

205 **Availability of data and material**

206

207 The dataset analysed in the current manuscript is available from [8]. The dataset analysed in
208 the vignette of the Bioconductor package [1] is available from [9] and accessible via the
209 covRNA package.

210 Bioconductor package availability:

211 Project home page: <https://bioconductor.org/packages/release/bioc/html/covRNA.html>

212 Operating system(s): Platform independent; multi-core systems

213 Programming language: R

214 License: GPL version 2 or later

215 **Competing interests**

216

217 Not applicable

218 **Funding**

219

220 The Collaborative Research Center / Transregio 124 - FungiNet Transregio provided financial
221 support (B2). LU was supported by the Friedrich-Ebert-Foundation (ref. nr. 1451239). Open
222 access publishing was supported by DFG and University of Würzburg.

223 **Author's contributions**

224

225 LU developed the Bioconductor package, analyzed the datasets, and wrote the manuscript.

226 TM created and supervised the project, and contributed to writing the manuscript. CWR pre-

227 processed several datasets and contributed to writing the manuscript. RFS contributed to

228 developing the Bioconductor package and to writing the manuscript. MD contributed to writing

229 the manuscript. All authors read and approved the final manuscript.

230 **Acknowledgements**

231

232 We would like to thank the Bioconductor core team and community for providing feedback

233 and support.

References

- 234
235 [1] <https://bioconductor.org/packages/release/bioc/html/covRNA.html>. Accessed 21
236 December 2018.
- 237 [2] Alter O, Brown PO, Botstein D. Singular value decomposition for genome-wide expression
238 data processing and modeling. *Proc Natl Acad Sci* 2000;97(18).
- 239 [3] Beissbarth T, Speed TP. GOstat: find statistically overrepresented gene ontologies within
240 a group of genes. *Bioinformatics* 2004;20(9).
- 241 [4] Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al. Gene
242 set enrichment analysis: a knowledge-based approach for interpreting genome-wide
243 expression profiles. *Proc Natl Acad Sci* 2005;102(43).
- 244 [5] Gentleman RC, Carey VJ, Bates DM, Bolstad B, Dettling M, Dudoit S, et al. Bioconductor:
245 open software development for computational biology and bioinformatics. *Genome biology*
246 2004;5(10).
- 247 [6] Dray S, Choler P, Dolédec S, Peres-Neto PR, Thuiller W, Pavoine S, et al. Combining the
248 fourth-corner and the rlq methods for assessing trait responses to environmental variation.
249 *Ecology* 2014;95(1)
- 250 [7] Dray S, Dufour AB. The ade4 package: implementing the duality diagram for ecologists.
251 *Journal of statistical software* 2007;22(4)
- 252 [8] Baty F, Ruediger J, Miglino N, Kern L, Borger P, Brutsche M. Exploring the transcription
253 factor activity in high-throughput gene expression data using RLQ analysis. *BMC*
254 *Bioinformatics* 2013;14:178.
- 255 [9] Czakai K, Leonhardt I, Dix A, Bonin M, Linde J, Einsele H, et al. Krüppel-like factor 4
256 modulates interleukin-6 release in human dendritic cells after in vitro stimulation with
257 *Aspergillus fumigatus* and *Candida albicans*. *Sci. Rep.* 2016; 6:27990.
- 258 [10] Passalacqua KD, Varadarajan A, Weist C, Ondov BD, Byrd B, Read TD, et al. Strand-
259 Specific RNA-Seq Reveals Ordered Patterns of Sense and Antisense Transcription in
260 *Bacillus anthracis*. *PLoS ONE* 2012;7(8): e43350.
- 261 [11] Yates F. Tests of significance for 2 × 2 contingency tables (with discussion). *Journal of*
262 *the Royal Statistical Society* 1984;147.
- 263 [12] Ter Braak CJF, Cormont A, Dray S (2012) Improved testing of species traits–
264 environment relationships in the fourth-corner problem. *Ecology* 2012;93(7).

265 [13] Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful
266 approach to multiple testing. *Journal of the Royal Statistical Society* 1995.