

BMJ Open

BMJ Open is committed to open peer review. As part of this commitment we make the peer review history of every article we publish publicly available.

When an article is published we post the peer reviewers' comments and the authors' responses online. We also post the versions of the paper that were used during peer review. These are the versions that the peer review comments apply to.

The versions of the paper that follow are the versions that were submitted during the peer review process. They are not the versions of record or the final published versions. They should not be cited or distributed as the published version of this manuscript.

BMJ Open is an open access journal and the full, final, typeset and author-corrected version of record of the manuscript is available on our site with no access controls, subscription charges or pay-per-view fees (<http://bmjopen.bmj.com>).

If you have any questions on BMJ Open's open peer review process please email info.bmjopen@bmj.com

BMJ Open

A simulation study to demonstrate the biases in the diagnoses of mental illnesses: major depressive episodes, dysthymia, and manic episodes

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2020-037022
Article Type:	Original research
Date Submitted by the Author:	26-Feb-2020
Complete List of Authors:	Chao, Yi-Sheng; Independent researcher, Independent researcher Lin, Kuan-Fu; National Taiwan University Hospital Yun-Lin Branch, Psychiatry Wu, Chao-Jung; UQAM, Département d'informatique Wu, Hsing-Chien; Taipei Hospital, Internal Medicine Hsu, Hui-Ting; Changhua Christian Healthcare System, Pathology Tsao, Lien-Cheng; Changhua Christian Healthcare System, Surgery Cheng, Yen-Po; Changhua Christian Healthcare System, Surgery Lai, Yi-Chun; National Yang Ming University Hospital, Chest Medicine Chen, Wei-Chih; Taipei Veterans General Hospital, Chest Medicine; National Yang-Ming University, Institute of Emergency and Critical Care Medicine
Keywords:	MENTAL HEALTH, Depression & mood disorders < PSYCHIATRY, EPIDEMIOLOGY, STATISTICS & RESEARCH METHODS

SCHOLARONE™
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

A simulation study to demonstrate the biases in the diagnoses of mental illnesses: major depressive episodes, dysthymia, and manic episodes

Yi-Sheng Chao^{1*}, Kuan-Fu Lin,² Chao-Jung Wu³, Hsing-Chien Wu⁴, Hui-Ting Hsu⁵, Lien-Cheng Tsao⁵, Yen-Po Cheng⁵, Yi-Chun Lai⁶, Wei-Chih Chen^{7,8}

¹Independent researcher, Montréal, H2X 0A8 Canada, ²National Taiwan University Hospital Yun-Lin Branch, Yunlin County, 640 Taiwan, ³Département d'informatique Université du Québec à Montréal, Montréal H3B 1B4 Canada, ⁴Taipei Hospital Ministry of Health and Welfare New Taipei city, 242 Taiwan, ⁵Changhua Christian Hospital, Changhua County 526, Taiwan, ⁶National Yang-Ming University Hospital, Yilan 260 Taiwan, ⁷Department of Chest Medicine, Taipei Veterans General Hospital, Taipei 112, Taiwan, ⁸Institute of Emergency and Critical Care Medicine, National Yang-Ming University, Taipei 112, Taiwan
*chaoyisheng@post.harvard.edu

Keywords: Frailty; bias; forward-stepwise regression; the Health and Retirement Study; index mining

20 Abstract

21 Objectives

22 Composite diagnostic criteria are likely to introduce biases to the diagnoses that
23 subsequently have poor relationships with input symptoms. This study aims to understand
24 the magnitudes of biases introduced to the diagnoses of three mental illnesses with large
25 disease burdens (major depressive episodes, dysthymic disorder, and manic episodes) and
26 the relationships between the diagnoses and the input symptoms.

27 Settings

28 Psychiatric care in general

29 Participants

30 Without real-world data available to the public, 100,000 subjects were simulated and the
31 input symptoms were assigned based on the assumed prevalence rates (0.05, 0.1, 0.3, 0.5,
32 and 0.7) and correlations between symptoms (0, 0.1, 0.4, 0.7, and 0.9). The input symptoms
33 were extracted from the diagnostic criteria of three mental illness. The diagnostic criteria
34 were transformed to mathematical equations to convert the input symptoms to diagnoses.

35 Primary and secondary outcomes

36 Biases due to data censoring or categorization introduced to the intermediate variables and
37 the three diagnoses were measured. The relationships between the input symptoms and
38 diagnoses were interpreted using forward stepwise linear regressions.

39 Results

40 The prevalence rates of the diagnoses were lower than those of the input symptoms and
41 proportional to the assumed prevalence rates and the correlations between the input
42 symptoms. Certain input or bias variables consistently explained the diagnoses better than
43 the others. Except for zero assumed correlations and 0.7 prevalence rates of the input
44 symptoms for the diagnosis of dysthymic disorder, the input variables could not fully explain
45 the diagnoses.

46 Conclusions

47 There are biases introduced to the diagnoses of three mental illnesses, major depressive
48 episodes, dysthymic disorder, and manic episodes. The design of the diagnostic criteria
49 determines the prevalence of the diagnoses, the relationships between the input symptoms
47 and the diagnoses, and the biases introduced. The importance of the input variables has
48 been largely distorted by the diagnostic criteria.

52 Trial registration

53 Not applicable

54 Strength and limitation

- 55 1. The prevalence of three mental illnesses were determined by the prevalence of the
56 input symptoms and modified by the diagnostic criteria and correlations between the
57 input variables in simulated populations.
- 58 2. Biases due to data censoring or categorization were introduced to the intermediate
59 variables and the three diagnoses of mental illnesses in simulated populations.

3. The diagnostic criteria modified the importance of the input variables and certain input or bias variables were given more weights than expected in simulated populations.
4. The design of diagnostic criteria influenced the prevalence. Even with the same input variable prevalence, dysthymic disorder was the most prevalent and major depressive episodes were the least prevalent in simulated populations.
5. This study is based on simulated data and needs to be verified with real-world data.

For peer review only

68 Background

69 The diagnoses of several mental illnesses in patients are made often based on a
70 variety of criteria. These criteria often involve symptoms complained by the patients.[1, 2]
71 For example, the diagnosis of major depressive disorder defined in the Diagnostic and
72 Statistical Manual of Mental Disorders, 4th Edition, Text Revision (DSM-IV-TR) requires at
73 least one major depressive episodes.[1, 2] For each major depressive episode, the major
74 criteria are “depressive mood and/or loss of interest or pleasure in life activities for at least 2
75 weeks”. [1, 2] In addition to qualify the major criteria, the patients need to report at least five
76 of the nine symptoms that “cause clinically significant impairment in social, work, or other
77 important areas of functioning almost every day”, including insomnia or hypersomnia and
78 fatigue or loss of interest.[1, 2] In other words, patients need to match both the major and
79 minor criteria before being diagnosed with a major depressive episode.

80 Historically this symptom-based diagnostic approach developed by Feighner et al.
81 has been widely accepted.[3, 4] Since then, mental illnesses can be diagnosed through
82 different sets of criteria. This approach is important because clinicians become capable of
83 screening important symptoms before diagnosing and treating patients accordingly. In fact,
84 these criteria can also be seen as composite measures that use multiple measures to
85 capture disorders that may not be quantified with single variables.[5, 6] Recent studies on
86 composite measures have found that composite measures are problematic because biases
87 can be introduced while aggregating information from input variables.[6] The biases emerge
88 while the sums of input variables are censored or while input variables are transformed
89 inadequately.[6, 7] These biases have been proven vital to the diagnosis of frailty syndrome,
90 a condition that often occurs in the elderly and is significant for several health outcomes.[6]
91 For the diagnosis of frailty syndrome using the Biological Syndrome Model,[8] biases alone
92 can explain more than 71% of the variances of the frailty diagnosis.[6]

93 Designed as composite measures, it is possible that the diagnostic criteria of mental
94 illnesses also introduce biases to diagnoses so that the diagnoses could not be fully
95 explained by the input symptoms listed in the criteria. This study aims to first understand the
96 relationships between mental symptoms and diagnoses and then quantify the potential role
97 of the biases regarding the diagnoses by simulating populations with different prevalence
98 rates and between-variable correlations of mental symptoms.

99 Methods

100 Assumptions and simulation parameters

101 Simulated populations with mental symptoms of different prevalence rates and
102 between-variable correlations were created to interpret the diagnoses and understand the
103 potential magnitudes of biases that could be introduced via data processing (reproducible
104 using data sets in the S 1 and S 2). Three diagnoses of mental illnesses were chosen for the
105 leading associated disease burdens:[2] major depressive episodes for the diagnosis of major
106 depressive disorder, dysthymic disorder, and manic episodes for the diagnosis of bipolar
107 disorder.[1]

108 There were assumptions made to simulate the populations (Table 1). First, for each
109 simulation the prevalence rates of the input symptoms were assumed to be similar for the

1
2
3 110 three diagnoses evaluated in this study. Second, the input symptoms for the diagnoses of
4 111 major depressive episodes and dysthymic disorder correlated with the same correlation
5 112 coefficients and those for the diagnosis of manic episodes correlated to one another.[9]
6 113 Third, the input symptoms for the diagnosis of manic episodes were created independently
7 114 of those for the diagnosis of the other two mental illnesses. The assumptions of the
8 115 prevalence rates and between-variable correlations were made because there was no
9 116 acceptable-quality data on the symptoms of mental illnesses published. There were studies
10 117 on the prevalence of mental illnesses,[10, 11] but the information on the prevalence of
11 118 mental symptoms was very limited. There were variables about depression or anxiety
12 119 collected in national surveys, such as the items collected through the Center for
13 120 Epidemiologic Studies Depression Scale.[6, 12-18] However, these variables were not the
14 121 symptoms used in the DSM-IV-TR. Lastly, we assumed that the diagnoses were made
15 122 accurately based on the presence of the input symptoms and the diagnostic criteria in the
16 123 DSM-IV-TR. However, this did not hold in the real world.[19] For simplicity and practical
17 124 reasons, we assumed perfect diagnostic quality by physicians and accurate reporting of the
18 125 input symptoms by patients in the simulated populations.

126 **Diagnostic criteria as mathematical functions**

127 The input symptoms were extracted from the major and minor criteria of the diagnoses and listed in
128 Table 2 to Table 4. The input symptoms, major and minor criteria, and the diagnoses were labelled
129 with new variable names. All input symptoms, items or domains in the major or minor criteria, and
130 the diagnoses were binomial variables, presenting zero and one for the absence and presence of the
131 symptoms, criteria, and the diagnoses respectively. For example, “insomnia” and “hypersomnia”
132 were extracted from one of the minor criteria for the diagnosis of major depressive episodes. “More
133 talkative than usual” and “pressure to keep talking” were extracted from one of the minor criteria
134 for the diagnosis of manic episodes.

135 Mathematical functions were generated based on the criteria to convert input symptoms into
136 diagnoses. For example, one of the minor criteria of dysthymic disorder was “poor appetite or
137 overeating”. This required two input symptoms and one bias variable to generate the criterion.[6]
138 “Poor appetite or overeating” equaling the sum of two input variables, “poor appetite” and
139 “overeating”, and a bias variable to achieve censoring of the sum of both variables.[6] The sum of
140 two binomial variables could exceed one and the bias variable had values of -1 for certain subjects to
141 obtain values less than or equal to one in all subjects.[6] In addition to adding variables together to
142 derive an intermediate variable or a diagnosis, multiplication, categorization, and other more
143 complicated methods were used in the diagnostic criteria to generate diagnosis variables and
144 domain variables in the major or minor criteria. For example, the diagnosis of dysthymic disorder
145 required the confirmation of both the major criteria, “depressed mood most of the day for more
146 days than not, for at least 2 years” and the minor criteria, “the presence of two or more of the
147 following symptoms”, at the same time. This was the same as multiplying two binomial variables to
148 obtain the diagnosis of dysthymic disorder. Other equations to generate the intermediate variables
149 and the diagnoses were listed and explained in Table 2 to Table 4.

150 **Generation of bias variables**

151 Bias variables were generated while binomial input symptoms were summed or multiplied
152 to obtain binomial intermediate or diagnosis variables.[6] Therefore the number of bias variables

1
2
3 153 depended on the complexity of how the diagnoses were made. For example, six of the nine items or
4 154 domains in the minor criteria for the diagnosis of major depressive episodes were the censored
5 155 sums of the input symptoms and six bias variables were derived along with the intermediate
6 156 variables that represented the items in the minor criteria. The other bias variables were described in
7 157 Table 2 to Table 4.

158 **Simulation parameters and simulated populations**

159 We simulated populations of 100,000 subjects. There were five prevalence rates to simulate
160 the input symptoms for the diagnosis of major depressive episodes, dysthymic disorder, and manic
161 episodes: 0.05, 0.1, 0.3, 0.5, and 0.7. The correlations between the input symptoms were
162 hypothesized to be 0, 0.1, 0.4, 0.7, and 0.9. There were 25 combinations of the assumed prevalence
163 rates and between-variable correlations. The presence of the input symptoms were randomly
164 assigned to the subjects after specifying the prevalence rates and between-variable correlations
165 between the input symptoms.[20, 21] The intermediate and diagnosis variables were derived
166 according to the equations in Table 2 to Table 4. For each combination of prevalence rates and
167 between-variable correlations, the populations were simulated for 100 times to obtain the mean
168 values and 95% confidence intervals (CIs) of derived prevalence rates, as well as the adjusted R
169 squared and p values to approximate the diagnosis variables.

170 **Diagnosis approximation**

171 Due to the existence of the biases, the input symptoms were not likely to fully explain the
172 diagnoses.[6] Therefore, the diagnoses were approximated by the input, bias, and intermediate
173 variables individually or collectively.[6, 12, 14, 16] The approximation was conducted using forward-
174 stepwise linear regressions.[6, 12, 14, 16, 22] The interpretability of the diagnoses by the input
175 symptoms and bias variables was assessed via adjusted R square: zero suggesting that the input
176 symptoms unrelated to the diagnosis; one suggesting that the input symptoms perfectly explained
177 the diagnosis.[14, 15, 23-26]

178 All statistical analyses were conducted within R environment (v3.4.1)[27] and RStudio
179 (v1.0.153).[28] P values less than 0.05 were considered statistical significant, two-tailed.

180 **Results**

181 The derived prevalence rates of the input symptoms of the three mental illnesses
182 matched the assumed rates in Figure 1. The derived correlations between the input
183 symptoms were close to assumed levels in S 3. The simulations were successful and
184 accurate based on the assumed prevalence rates and correlations.

185 **Prevalence of intermediate variables**

186 The items in the major and minor criteria were the intermediate variables necessary to
187 create the diagnoses. The methods to generate the intermediate variables were as important to the
188 prevalence rates of the intermediate variables as the prevalence rates and correlations of the input
189 symptoms in Figure 2. The intermediate variable, significant unintentional weight loss or gain, was
190 created by summing and censoring two binomial variables with values of zero and one (significant
191 unintentional weight loss; significant unintentional weight gain). The prevalence rates of the

1
2
3 192 intermediate variables were larger than those of the two input symptoms regardless of the assumed
4 193 prevalence rates or between-variable correlations of the input symptoms.

6 194 In contrast, the diagnosis of dysthymic disorder was a multiplication product of two binomial
7 195 variables, the major and minor criteria, and the prevalence rates of dysthymic disorder were lower
8 196 than those of the major and minor criteria under all combinations of assumed correlations and
9 197 prevalence rates in Figure 3.

12 198 Prevalence of mental illnesses

14 199 The prevalence rates of three diagnoses were plotted against the assumed prevalence rates
15 200 and correlations of the input symptoms in Figure 3 to Figure 5 and listed in Table 5. None of the
16 201 three diagnoses had prevalence rates exceeding those of the input symptoms. In general, higher
17 202 prevalence rates or between-variable correlations of the input symptoms were associated with
18 203 higher prevalence rates in the three diagnoses, except for manic episodes that had higher
19 204 prevalence rates (0.692) assuming zero correlations and 0.7 prevalence rates than the prevalence
20 205 rate (0.679) assuming 0.1 correlations and 0.7 prevalence rates of the input symptoms. When
21 206 compared across Figure 3 to Figure 5, given the same assumed prevalence rates and between-
22 207 variable correlations of the input symptoms, the diagnostic criteria of dysthymic disorder
23 208 consistently generated diagnoses of the highest prevalence rates and the criteria of major
24 209 depressive episodes created diagnoses of the least prevalence rates (see Table 5 for details).

29 210 Associations between the diagnoses and individual input symptoms and bias 30 211 variables

32 212 The diagnoses were interpreted by the input symptoms (including intermediate variables) and the
33 213 bias variables individually first. Take dysthymic disorder for example, the diagnosis was interpreted
34 214 with the input symptoms, the bias variables, and both in Figure 6. For each simulation, the diagnosis
35 215 of dysthymic disorder was approximated with an increasing number of the input symptoms, the bias
36 216 variables, or both. After selecting the variables that best approximated the diagnosis based on
37 217 adjusted R-squared, the input symptoms could explain a proportion of 0.955 of the diagnosis
38 218 variance and the bias variables could explain at most a proportion of 0.405 of the diagnosis
39 219 variance in Figure 6. With all variables used in the regression, the diagnosis could be perfectly explained by
40 220 the input symptoms and bias variables (adjusted R-squared = 1). By repeating the same procedures
41 221 to the diagnoses, the individual input symptoms and the bias variables that individually best
42 222 explained the diagnoses were listed in **Error! Reference source not found.** and Table 7 respectively.

47 223 For the diagnosis of major depressive episodes, the first and second items in the major
48 224 criteria (variable names: mde_ma1 for or mde_ma2 in Table 2) individually explained the
49 225 diagnosis the best depending on the assumed prevalence rates and correlations in **Error! Reference**
50 226 **source not found.** For the diagnosis of dysthymic disorder, the major criteria (dys_ma in Table 3)
51 227 consistently and individually explained the diagnosis the best. For the diagnosis of manic episodes,
52 228 the third item of the major criteria (man_ma3 in Table 4) individually explained the diagnosis the
53 229 best in all combinations of assumed prevalence rates and correlations. However, the proportions of
54 230 diagnosis variances best explained by individual input symptoms varied in a large range between
55 231 0.001 to 0.974 depending on the assumed prevalence rates and between-variable correlations.
56 232 Based on the adjusted R-squared for individual input symptoms, certain input variables were more
57 233 important than other symptoms due to high correlation with the diagnoses, such as the major

234 criteria for the diagnosis of dysthymic disorder. The prevalence rates and between-variable
235 correlations were important to determine the relationships between input symptoms and diagnoses.

236 Similarly, there were bias variables that consistently explained the diagnoses the best in Table 7. For
237 the diagnosis of major depressive episodes, the bias due to categorization of the numbers of
238 confirmed input symptoms up to three or four (mde_bias1 or mde_bias2 respectively in Table 2)
239 were the leading bias variable. The diagnosis of major depressive episodes not explained by the
240 input symptoms or information censoring (mde_bias in Table 2) were the leading bias variable in
241 two combinations of the assumed prevalence rates and correlations. For the diagnosis of dysthymic
242 disorder, the residual of the diagnosis not explained by the major and minor criteria (dys_bias in
243 Table 3) and the bias due to categorization of the confirmed input symptoms in the minor criteria
244 (dys_mi_bias) were the leading bias variables. For the diagnosis of manic episodes, the bias due to
245 categorization of the number of confirmed input symptoms in the minor criteria up to three
246 (man_bias1 in Table 4) was the leading bias variables, except for two combinations of the assumed
247 prevalence rates and correlations, in which the bias due to categorization of the confirmed input
248 symptoms in the minor criteria up to four (man_bias2 in Table 4) best explained the diagnosis.
249 However, the proportions of diagnosis variances explained by individual bias variables varied in a
250 wide range from zero to 0.87. Depending on the assumed prevalence rates and between-variable
251 correlations of the input symptoms, certain bias variables were more important than other bias
252 variables and even some input variables. The assumed prevalence rates and between-variable
253 correlations were important factors for the relationships between the bias variables and the
254 diagnoses.

255 In general, the proportions of the diagnosis variances could be explained by either individual input
256 symptoms or single bias variables were low when the prevalence rates and between-variable
257 correlations of the input symptoms were assumed low. With higher assumed prevalence rates or
258 correlations, the proportions of the diagnoses explained by the single input symptoms or bias
259 variables were higher. Across three diagnoses, the diagnosis of dysthymic disorder could be better
260 explained by its own single input variables (higher adjusted R-squared) and the diagnosis of major
261 depressive episodes was associated with the least adjusted R-squared. The bias variables of the
262 diagnosis of manic episodes could explain the diagnosis individually better than the bias variables of
263 the other two diagnoses.

264 **Approximation of the diagnoses with input symptoms**

265 When the diagnoses were approximated with all input symptoms of their own in Table
266 8, there were always some diagnosis variances that could not be explained by the input
267 symptoms. In other words, the input symptoms could not fully explain the diagnoses, except
268 for the diagnosis of dysthymic disorder that could be fully explained by the input symptoms
269 (adjusted R-squared = 1) assuming zero between-variable correlations and 0.7 prevalence
270 rates for the input symptoms. In Table 8, the proportions of diagnosis variances explained by
271 input symptoms increased with higher assumed prevalence rates or between-variable
272 correlations of the input symptoms in general. The input symptoms of dysthymic disorder
273 explained the diagnosis better than those of the other two diagnoses under all combinations
274 of assumed prevalence rates and between-variable correlations. The diagnosis of major
275 depressive episodes was the worst approximated with its own input symptoms in terms of
276 adjusted R-squared. However, the proportions of diagnosis variances explained by own

1
2
3 277 input symptoms varied in a wide range between 0.003 to 1.0. The assumed prevalence rates
4 278 and between-variable correlations of the input symptoms and the design of the diagnostic
5 279 criteria were all important for the relationships between input symptoms and diagnoses.

8 280 **Approximating the diagnoses with bias variables**

9 281 The diagnoses were approximated with the bias variables of their own. The bias
10 282 variables always explained some of the diagnosis variances, except for the diagnosis of
11 283 dysthymic disorder assuming zero between-variable correlations and 0.7 prevalence rates
12 284 for the input symptoms (adjusted R-squared = 0). With increasing assumed between-
13 285 variable correlations for the input symptoms, the adjusted R-squared increased. However,
14 286 given the same assumed between-variable correlations, the proportions of diagnosis
15 287 variances explained by the bias variables might increase or decrease with the assumed
16 288 prevalence rates. Compared to the adjusted R-squared in Table 8, the proportions of the
17 289 diagnosis variances explained by the bias variables were always smaller than those
18 290 explained by the input symptoms in Table 9. However, the proportions of diagnosis variances
19 291 explained by bias variables also varied in a wide range from zero to 0.89. The assumed
20 292 prevalence rates and between-variable correlations of input symptoms and the design of the
21 293 diagnostic criteria were important for the relationship between the bias variables and the
22 294 diagnoses. Only when the input symptoms for the diagnosis of dysthymic disorder were
23 295 randomly and independently prevalent to 70% of the simulated populations, the bias
24 296 variables became irrelevant to the diagnosis.

30 31 **Discussion**

32 298 This study is the first attempt to understand the relationships between the input
33 299 symptoms and the diagnoses of three mental illnesses: major depressive episodes (at least
34 300 one episode required for the diagnosis of major depressive disorder), dysthymic disorder,
35 301 and manic episodes. The diagnostic criteria of three mental illnesses have been reviewed
36 302 and rewrote as mathematical functions. Simulated populations, 100,000 for each of 100
37 303 simulations, with input symptoms of the three diagnoses were created. For simplicity and
38 304 practicality reasons, the presence of the input symptoms was randomly assigned and the
39 305 input symptoms were assumed to have uniform prevalence rates and between-variable
40 306 correlations. There were 25 combinations of assumed prevalence rates and between-
41 307 variable correlations simulated.

42 308 Mathematically, the diagnostic criteria are functions and composite measures to
43 309 transform the information from the input variables to diagnoses. There are bias variables
44 310 created in the process of information transformation.[6] There are three major mechanisms
45 311 of introducing biases, censoring, data categorization[7] and multiplication of input symptoms
46 312 with values of zero and one presenting the absence and presence of the symptoms.[6]
47 313 These mechanisms introduce information or biases that cannot be fully explained by the
48 314 input symptoms.[6] The biases introduced can sometimes explain more than half of the
49 315 variances of the diagnoses depending on the prevalence rates and between-variable
50 316 correlations of the input symptoms (e.g. assuming input symptoms with 0.7 or 0.9
51 317 prevalence rates for the three diagnoses). The findings show that the design of the
52 318 diagnostic criteria important for bias introduction and significant for the prevalence of the
53 319 diagnoses in populations, the relationships between the input symptoms and the diagnoses,
54 320 and the relationships between the bias variables and the diagnoses.

321 **The impact of the diagnostic criteria**

322 With the same assumptions in the prevalence rates and between-variable
323 correlations of the input symptoms, the design of the diagnostic criteria of three mental
324 illnesses can be compared to each other. The design of diagnostic criteria transform input
325 symptoms to various diagnosis prevalence rates with implicit upper limits (i.e. no more
326 prevalent than the input symptoms), unacknowledged differential weights on the input
327 symptoms (i.e. certain input symptoms explaining the diagnoses better), and the introduction
328 of biases (i.e. due to censoring, data categorization, or multiplication).

329 We were the first to notice that the prevalence rates of the three diagnoses were
330 lower than those of the input symptoms, if randomly distributed with uniform prevalence
331 rates and correlations. Given similar assumed input symptom prevalence and correlations,
332 dysthymic disorder is the most prevalent and major depressive episodes are the least. The
333 diagnosis of dysthymic disorder can be better explained by own input symptoms individually
334 or collectively. The diagnosis of major depressive episodes is the worst explained by own
335 input symptoms individually or collectively. As expected, the diagnosis of the three mental
336 illness are similar to composite measures or indices and are subject to the biases introduced
337 by data processing given all combinations of the assumed prevalence rates and between-
338 variable correlations of the input symptoms.[6] There is only one exception: dysthymic
339 disorder with the input symptoms that are randomly and independently present in 70% of the
340 population. This is because the diagnosis of dysthymic disorder is a multiplication product of
341 the major and minor criteria. Without correlations, everyone in the population is certain to
342 qualify for the minor criteria (probability of 100% because of having at least two out of six
343 item in the minor criteria: mathematically $[C(2,6) + C(3,6) + C(4,6) + C(5,6) + C(6,6)] \times$
344 $(0.7)^6 = 37 \times 0.117 = 4.35 > 100\%$). When 70% of the population are also randomly
345 assigned with the major criteria, the diagnosis of dysthymic disorder can be fully explained
346 by the major criteria alone. In fact, without correlations between input symptoms it only
347 requires each of the six items in the minor criteria to be randomly assigned to 54.8%
348 $[(1/37)^{(1/6)}]$ of the population for everyone to qualify for the minor criteria and the diagnosis
349 can be fully explained by the minor and major criteria.

350 **Distortion of the input symptoms**

351 The importance of the input symptoms has been distorted due to the functions to
352 generate the diagnoses. This has been proven in the diagnosis of frailty.[6] In other words,
353 based on the functions to generate the diagnoses, the input symptoms are differentially
354 weighted without the weights being explicitly acknowledged. The most prominent is the
355 diagnosis of dysthymic disorder, more than 90% of whose variance can be explained by its
356 major criteria assuming 0.7 or 0.9 between-variable correlations for the input symptoms in
357 Table 6. Another example is that the third item of the major criteria for the diagnosis of manic
358 episodes, "irritable mood", individually predicts the diagnosis better than any other input
359 symptoms. Assuming 0.9 correlations between input symptoms, this input symptom has
360 been put more weight than others and can explain more than 91.8% of the diagnosis
361 variance. Based on the texts in the DSM-IV-TR, we don't think this symptom should be
362 emphasized to this degree and consider the diagnostic criteria are imposing implicit and
363 unequal weights to the input symptoms, as well as introducing biases.

364 **Future directions**

365 We think it important to rethink the role and importance of the diagnostic system.
366 Current approaches are embedded with implicit assumptions of the prevalence rates of the
367 diagnoses (no higher than input symptoms), unacknowledged weights to input symptoms
368 (certain input symptoms explaining the diagnoses much better), and biases that could not be
369 explained by the input symptoms. The diagnosis of dysthymic disorder is probably
370 accidentally “designed” to be more prevalent than that of major depressive episodes or
371 manic episodes based on the diagnostic criteria assuming input symptoms with the same
372 prevalence rates. In the real world, there are other important issues related to the diagnostic
373 criteria. For example, diagnosis is not closely linked to treatment,[19, 29] diagnosis is not
374 well made particularly by non-psychiatrists,[30] and there are two diagnostic systems (DSM
375 and International Classification of Disease) that require efforts to harmonize.[31] Amid these
376 issues, we think the diagnostic criteria for mental illnesses should be reviewed and improved
377 in a way that they can be easier to understand and use without introducing biases and can
378 be closely linked to clinical decisions. We are developing methods to better detect symptom-
379 based conditions and proposing methods to search for neglected mental symptoms.

380 **Limitations**

381 The strength of this study is the use of simple assumptions in simulated populations
382 that enables the comparison of the diagnostic criteria of three mental illnesses. However, the
383 assumptions in the prevalence rates and between-variable correlations for the input
384 symptoms might not be realistic. Some of the assumptions are unlikely to hold in the real
385 world. However, this is the only option for us due to the lack of real-world data on the
386 prevalence of the input symptoms. In addition, the translation from symptoms to diagnoses
387 was assumed to be perfect based on the diagnostic criteria.

388 **Conclusion**

389 To the best of our knowledge, there is no study on the relationships between the
390 input symptoms and the diagnoses. The input symptoms were extracted from the diagnostic
391 criteria and the diagnostic criteria were transformed to mathematical equations. Without
392 mental illness data available to the public, 100,000 subjects were simulated with different
393 assumptions on the prevalence rates (0.05, 0.1, 0.3, 0.5, and 0.7) and correlations (0, 0.1,
394 0.4, 0.7, and 0.9) of the input symptoms. We found that biases were introduced to the
395 diagnoses of three mental illnesses, major depressive episodes, dysthymic disorder, and
396 manic episodes. The prevalence rates of the diagnoses were proportional to the assumed
397 prevalence rates and between-variable correlations of the input symptoms. Certain input
398 symptoms were more important than the others to explain the diagnoses. However, the input
399 symptoms could not fully explain the diagnoses, except when the input symptoms
400 independent to each other with 0.7 prevalence rates were used for the diagnosis of
401 dysthymic disorder.

1
2
3 402 **Declarations**

4
5
6 403 **Patient and Public Involvement**

7 404 This is a simulation study that does not require the input from patients or the public.
8

9 405 **Acknowledgments**

10 406 Not applicable.
11
12

13 407 **Ethical Statement**

14 408 Not applicable.
15

16 409 **Funding Statement**

17 410 The authors received no specific funding for this work.
18
19

20 411 **Consent to participate**

21 412 Not applicable
22

23 413 **Consent for publication**

24 414 Not applicable
25
26

27 415 **Data Availability**

28 416 No real-world data used. All analysis based on simulations reproducible with the files in the
29 417 Supplemental materials.
30
31

32 418 **Competing Interests**

33 419 YSC is currently employed by the Canadian Agency for Drugs and Technologies in Health.
34 420 The other authors have declared that no competing interests exist.
35
36

37 421 **Authors' Contributions**

38 422 YSC conceptualized and designed this study, managed and analyzed data and drafted the
39 423 manuscript. KFL assisted in the interpretation of the diagnostic criteria. CJW assisted in data
40 424 management and computation. HCW, HTH, LCT, YPC, YCL, and WCC participated in the design of this
41 425 study. All authors reviewed and approved the manuscript.
42
43
44 426
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

427 **References**

- 428 1. American Psychiatric Association. Diagnostic and Statistical Manual of Mental Disorders,
429 Fourth Edition, Text Revision (DSM-IV-TR®). Washington, DC: American Psychiatric Association
430 Publishing; 2010.
- 431 2. Center for Substance Abuse Treatment. Managing Depressive Symptoms in Substance Abuse
432 Clients During Early Recovery. Rockville, MD: Substance Abuse and Mental Health Services
433 Administration (US); 2008.
- 434 3. Feighner JP, Robins E, Guze SB, Woodruff RA, Jr., Winokur G, Munoz R. Diagnostic criteria for
435 use in psychiatric research. *Arch Gen Psychiatry*. 1972;26(1):57-63. Epub 1972/01/01.
- 436 4. Kendler KS, Munoz RA, Murphy G. The development of the Feighner criteria: a historical
437 perspective. *The American journal of psychiatry*. 2010;167(2):134-42. Epub 2009/12/17. doi:
438 10.1176/appi.ajp.2009.09081155.
- 439 5. Chao Y-S, Wu C-J. PP46 When Composite Measures Or Indices Fail: Data Processing Lessons.
440 *International Journal of Technology Assessment in Health Care*. 2019;34(S1):83-. Epub 01/03. doi:
441 10.1017/S0266462318002088.
- 442 6. Chao Y-S, Wu H-C, Wu C-J, Chen W-C. Index or illusion: The case of frailty indices in the
443 Health and Retirement Study. *PLOS ONE*. 2018;13(7):e0197859. doi: 10.1371/journal.pone.0197859.
- 444 7. Barnwell-Menard JL, Li Q, Cohen AA. Effects of categorization method, regression type, and
445 variable distribution on the inflation of Type-I error rate when categorizing a confounding variable.
446 *Stat Med*. 2015;34(6):936-49. Epub 2014/12/17. doi: 10.1002/sim.6387.
- 447 8. Cigolle CT, Ofstedal MB, Tian Z, Blaum CS. Comparing models of frailty: the Health and
448 Retirement Study. *J Am Geriatr Soc*. 2009;57(5):830-9. Epub 2009/05/21. doi: 10.1111/j.1532-
449 5415.2009.02225.x.
- 450 9. Brown TA, Chorpita BF, Korotitsch W, Barlow DH. Psychometric properties of the Depression
451 Anxiety Stress Scales (DASS) in clinical samples. *Behaviour research and therapy*. 1997;35(1):79-89.
- 452 10. Lim GY, Tam WW, Lu Y, Ho CS, Zhang MW, Ho RC. Prevalence of Depression in the
453 Community from 30 Countries between 1994 and 2014. *Scientific reports*. 2018;8(1):2861-. doi:
454 10.1038/s41598-018-21243-x. PubMed PMID: 29434331.
- 455 11. Smith DJ, Nicholl BI, Cullen B, Martin D, Ul-Haq Z, Evans J, et al. Prevalence and
456 Characteristics of Probable Major Depression and Bipolar Disorder within UK Biobank: Cross-
457 Sectional Study of 172,751 Participants. *PLOS ONE*. 2013;8(11):e75362. doi:
458 10.1371/journal.pone.0075362.
- 459 12. Chao Y-S, Wu C-J. PD26 Principal Component Approximation: Canadian Health Measures
460 Survey. *International Journal of Technology Assessment in Health Care*. 2019;34(S1):138-9. Epub
461 01/03. doi: 10.1017/S026646231800301X.
- 462 13. Chao Y-S, Wu C-J, Wu H-C, Chen W-C. Trend analysis for national surveys: Application to all
463 variables from the Canadian Health Measures Survey cycle 1 to 4. *PLOS ONE*. 2018;13(8):e0200127.
464 doi: 10.1371/journal.pone.0200127.
- 465 14. Chao Y-S, Wu H-C, Wu C-J, Chen W-C. Principal Component Approximation and
466 Interpretation in Health Survey and Biobank Data. *Frontiers in Digital Humanities*. 2018;5(11). doi:
467 10.3389/fdigh.2018.00011.
- 468 15. Chao YS, Wu HC, Wu CJ, Chen WC. Stages of Biological Development across Age: An Analysis
469 of Canadian Health Measure Survey 2007-2011. *Front Public Health*. 2018;5(2296-2565 (Print)). doi:
470 10.3389/fpubh.2017.00355. eCollection 2017.
- 471 16. Chao Y-S, Wu C-J. PD25 Principal Component Approximation: Medical Expenditure Panel
472 Survey. *International Journal of Technology Assessment in Health Care*. 2019;34(S1):138-. Epub
473 01/03. doi: 10.1017/S0266462318003008.
- 474 17. Chao Y-S, Wu C-J, Chen T-S. Risk adjustment and observation time: comparison between
475 cross-sectional and 2-year panel data from the Medical Expenditure Panel Survey (MEPS). *Health
476 Information Science and Systems*. 2014;2:5. doi: 10.1186/2047-2501-2-5. PubMed PMID:
477 PMC4340859.

- 1
2
3 478 18. Chao YS, Wu HT, Scutari M, Chen TS, Wu CJ, Durand M, et al. A network perspective on
4 479 patient experiences and health status: the Medical Expenditure Panel Survey 2004 to 2011. BMC
5 480 health services research. 2017;17(1472-6963 (Electronic)). doi: 10.1186/s12913-017-2496-5.
6 481
7 482 19. Bonnin JE. Treating without diagnosis: psychoanalysis in medical settings in Argentina. 2015.
8 483 20. Leisch F, Weingessel A, Hornik K. On the generation of correlated artificial binary data. 1998.
9 484 21. Leisch F, Weingessel A, Hornik K. bindata: Generation of Artificial Binary Data, 2012. URL
10 484 [http://cran/](http://cran.r-project.org/package=bindata) R-project org/package= bindata R package version 09-19.
11 485 22. Lumley T, Lumley MT. Package 'leaps'. Regression Subset Selection Thomas Lumley Based on
12 486 Fortran Code by Alan Miller Available online: [http://cran/](http://cran.r-project.org/package=leaps) R-project org/package= leaps (accessed on
13 487 18 March 2018). 2013.
14 488 23. Hastie T, Tibshirani R, Friedman J. The Elements of Statistical Learning: Data Mining,
15 489 Inference, and Prediction, Second Edition: Springer New York; 2009.
16 490 24. James G, Witten D, Hastie T, Tibshirani R. An Introduction to Statistical Learning: with
17 491 Applications in R: Springer New York; 2013.
18 492 25. Chao Y-S, Wu C-J. Principal component-based weighted indices and a framework to evaluate
19 493 indices: Results from the Medical Expenditure Panel Survey 1996 to 2011. PLoS ONE.
20 494 2017;12(9):e0183997. doi: 10.1371/journal.pone.0183997. PubMed PMID: PMC5590867.
21 495 26. Chao YS, Wu HT, Wu CJ. Feasibility of Classifying Life Stages and Searching for the
22 496 Determinants: Results from the Medical Expenditure Panel Survey 1996-2011. Front Public Health.
23 497 2017;5:247(2296-2565 (Print)). doi: 10.3389/fpubh.2017.00247. eCollection 2017.
24 498 27. R Development Core Team. R: A language and environment for statistical computing.
25 499 Vienna, Austria: R Foundation for Statistical Computing; 2016.
26 500 28. RStudio Team. RStudio: Integrated Development for R. Boston, MA: RStudio, Inc.; 2016.
27 501 29. Demyttenaere K, Bonnewyn A, Bruffaerts R, De Girolamo G, Gasquet I, Kovess V, et al.
28 502 Clinical factors influencing the prescription of antidepressants and benzodiazepines: Results from
29 503 the European study of the epidemiology of mental disorders (ESEMeD). Journal of affective
30 504 disorders. 2008;110(1-2):84-93.
31 505 30. Margolis RL. Nonpsychiatrist house staff frequently misdiagnose psychiatric disorders in
32 506 general hospital inpatients. Psychosomatics. 1994;35(5):485-91.
33 507 31. First MB. Harmonisation of ICD-11 and DSM-V: opportunities and challenges. The British
34 508 Journal of Psychiatry. 2009;195(5):382-90.

509

510

511

512 **Table 1. The assumptions and parameters in the simulations**

Assumptions		
1	Equal prevalence rates for the input symptoms of the same diagnosis; presence of input symptoms assigned randomly	
2	Same correlations between the input symptoms of the diagnoses of major depressive episodes and dysthymic disorder; same correlations between the input symptoms of manic episodes	
3	The input symptoms of manic episodes created independent of those of major depressive episodes and dysthymic disorder	
4	Diagnoses made accurately based on the diagnostic criteria and symptoms reported accurately by patients	
Parameters of input symptoms of the same diagnosis for each simulation		
1	Population sizes	10,000
2	Prevalence rates (uniform for all input symptoms in a simulation)	0.05, 0.1, 0.3, 0.5, and 0.7
3	Correlations (uniform between all input symptoms of the same diagnosis in a simulation)	0, 0.1, 0.4, 0.7, and 0.9
4	Number of simulations for each combination of the assumed prevalence rates and between-variable correlations of the input symptoms	100

513

514

515

516

517

518

519 **Table 2. The input symptoms, intermediate variables, and bias variables for the diagnosis of major depressive episodes.**

Classification of symptoms	Criterion variable	Domains in the major or minor criteria	Domain variables	Symptoms	Symptom variables	Equations to derive diagnosis or domain variables	Approximation by linear regression	Mechanisms related to introducing biases
Major depressive episode (variable = mde)						$mde = mde_ma1 \times mde_ma2 \times (mde_mi3 + mde_mi4 + mde_mi5 + mde_mi6 + mde_mi7 + mde_mi8 + mde_mi9 + mde_bias1) + (1 - mde_ma1 \times mde_ma2) \times (mde_ma1 \times mde_ma2) \times (mde_mi3 + mde_mi4 + mde_mi5 + mde_mi6 + mde_mi7 + mde_mi8 + mde_mi9 + mde_bias2)$	$mde = intercept + coef1 \times mde_ma1 + coef2 \times mde_ma2 + coef3 \times mde_mi3 + coef4 \times mde_mi4 + coef5 \times mde_mi5 + coef6 \times mde_mi6 + coef7 \times mde_mi7 + coef8 \times mde_mi8 + coef9 \times mde_mi9 + coef10 \times mde_bias$	1) Multiplication to create the situations when one or two symptoms in the major criteria confirmed and the bias (mde_bias) calculated by extracting the information of the diagnosis not explained by the input symptoms and two bias variables generated by censoring (mde_bias1 and mde_bias2) 2) Categorizing of the sum of the input symptoms in the minor criteria at the threshold of three or four (mde_bias1 and mde_bias2)
Major criteria, essential for diagnosis		Depressed mood or a loss of interest or pleasure in daily activities for more than two weeks.						
		Depressed mood for more than two weeks.	mde_ma1					
		Loss of interest or pleasure in daily activities for more than two weeks.	mde_ma2					
Minor criteria (at least 5 of the symptoms including the two in major criteria)	mde_mi							
		Significant unintentional weight loss or gain	mde_mi3			$mde_mi3 = mde_mi3_1 + mde_mi3_2 + mde_mi3_bias$		Censoring of the sum of multiple input variables
				Significant unintentional weight gain	mde_mi3_1			
				Significant unintentional weight loss	mde_mi3_2			
				Information of the domain not explained by the input variables	mde_mi3_bias			
		Insomnia or sleeping too much	mde_mi4			$mde_mi4 = mde_mi4_1 + mde_mi4_2 + mde_mi4_bias$		Censoring of the sum of multiple input variables
				Insomnia	mde_mi4_1			
				Sleeping too much	mde_mi4_2			
				Information of the domain not explained by the input variables	mde_mi4_bias			
		Agitation or psychomotor retardation noticed by others	mde_mi5			$mde_mi5 = mde_mi5_1 + mde_mi5_2 + mde_mi5_bias$		Censoring of the sum of multiple input variables
				Agitation	mde_mi5_1			

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

			Psychomotor retardation noticed by others	mde_mi5_2	
			Information of the domain not explained by the input variables	mde_mi5_bias	
	Fatigue or loss of energy	mde_mi6		mde_mi6 = mde_mi6_1 + mde_mi6_2 + mde_mi6_bias	Censoring of the sum of multiple input variables
			Fatigue	mde_mi6_1	
			Loss of energy	mde_mi6_2	
			Information of the domain not explained by the input variables	mde_mi6_bias	
	Feelings of worthlessness or excessive guilt	mde_mi7		mde_mi7 = mde_mi7_1 + mde_mi7_2 + mde_mi7_bias	Censoring of the sum of multiple input variables
			Feelings of worthlessness	mde_mi7_1	
			Feelings of excessive guilt	mde_mi7_2	
			Information of the domain not explained by the input variables	mde_mi7_bias	
	Diminished ability to think or concentrate, or indecisiveness+	mde_mi8		mde_mi8 = mde_mi8_1 + mde_mi8_2 + mde_mi8_bias	Censoring of the sum of multiple input variables
			Diminished ability to think or concentrate	mde_mi8_1	
			Indecisiveness	mde_mi8_2	
			Information of the domain not explained by the input variables	mde_mi8_bias	
	Recurrent thoughts of death	mde_mi9			
	Information due to categorization (choosing three domains in minor criteria)	mde_bias1			Bias introduced to categorize the sum of the number of confirmed symptoms in the minor criteria
	Information due to categorization (choosing four domains in minor criteria)	mde_bias2			Bias introduced to categorize the sum of the number of confirmed symptoms in the minor criteria
	Information of diagnosis not explained by the domains	mde_bias			Information of the diagnosis not explained by the input variables and two bias variables generated due to data categorization

520

521

523 Table 3. The input symptoms, intermediate variables, and bias variables for the diagnosis of dysthymic disorder.

Classification of symptoms	Criterion variable	Major or minor criteria (domains)	Intermediate variables	Symptoms	Symptom variables	Equations to generate diagnosis or domain variables	Approximation	Mechanisms related to introducing biases
Dysthymia (variable = dys)						$dys = dys_ma \times dys_mi$	$dys = intercept + coef1 \times dys_ma + coef2 \times dys_mi + coef3 \times dys_bias$	Multiplication to create the situations where both the major and minor criteria met (union of two binomial variables, $mde_ma \times mde_mi$) and the bias variable (dys_bias) equivalent to the residual of the diagnosis not explained by the input symptoms and the bias variables due to censoring and categorization
Major criteria, essential for diagnosis		Depressed mood most of the day for more days than not, for at least 2 years	dys_ma					
Minor criteria (at least 2 items)			dys_mi			$dys_mi = dys_mi1 + dys_mi2 + dys_mi3 + dys_mi4 + dys_mi5 + dys_mi6 + dys_mi_bias$		Categorizing of the sum of multiple input variables
		Poor appetite or overeating	dys_mi1	Poor appetite Overeating Information of the domain not explained by the input variables	dys_mi1_1 dys_mi1_2 dys_mi1_bias	$dys_mi1 = dys_mi1_1 + dys_mi1_2 + dys_mi1_bias$		Censoring of the sum of multiple input variables
		Insomnia or sleeping too much*	dys_mi2/mde_mi4	Insomnia Sleeping too much Information of the domain not explained by the input variables	mde_mi4_1 mde_mi4_2 mde_mi4_bias	$dys_mi2 = mde_mi4 = mde_mi4_1 + mde_mi4_2 + mde_mi4_bias$		Censoring of the sum of multiple input variables
		Low energy or fatigue*	dys_mi3/mde_mi6	Fatigue Loss of energy (low energy) Information of the domain not explained by the input variables	mde_mi6_1 mde_mi6_2 mde_mi6_bias	$dys_mi3 = mde_mi6 = mde_mi6_1 + mde_mi6_2 + mde_mi6_bias$		Censoring of the sum of multiple input variables
		Low self-esteem Poor concentration or difficulty making decisions*	dys_mi4 dys_mi5/mde_mi8	Diminished ability to think or concentrate (Poor concentration) difficulty making decisions (indecisiveness) Information of the domain not explained by the input variables	mde_mi8_1 mde_mi8_2 mde_mi8_bias	$dys_mi5 = mde_mi8 = mde_mi8_1 + mde_mi8_2 + mde_mi8_bias$		Censoring of the sum of multiple input variables
		Feelings of hopelessness	dys_mi6					

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

	Information of minor criteria not explained by input variables	dys_mi_bias	Bias introduced by categorizing the number of input symptoms confirmed in the minor criteria
Information of diagnosis not explained by major or minor criteria	dys_bias		Information of the diagnosis not explained by the input symptoms and the bias variables generated due to data categorization (dys_mi_bias)

524 *The same input symptoms used for the diagnosis of major depressive episodes.

525

For peer review only

527 Table 4. The input symptoms, intermediate variables, and bias variables for the diagnosis of manic episodes.

Classification of symptoms	Criterion variable	Major or minor criteria (domains)	Domain variables	Symptoms	Symptom variables	Equations	Approximation	Mechanisms related to introducing biases
Manic episode (variable = manic)						$\text{manic} = (1 - \text{man_ma1} \times \text{man_ma2}) \times (\text{man_ma1} + \text{man_ma2}) \times \text{man_ma3} \times (\text{man_mi1} + \text{man_mi2} + \text{man_mi3} + \text{man_mi4} + \text{man_mi5} + \text{man_mi6} + \text{man_mi7} + \text{man_bias1}) + [1 - (1 - \text{man_ma1} \times \text{man_ma2})(\text{man_ma1} + \text{man_ma2})] \times \text{man_ma3} \times (\text{man_mi1} + \text{man_mi2} + \text{man_mi3} + \text{man_mi4} + \text{man_mi5} + \text{man_mi6} + \text{man_mi7} + \text{man_bias2})$	$\text{manic} = \text{intercept} + \text{coef1} \times \text{man_ma1} + \text{coef2} \times \text{man_ma2} + \text{coef3} \times \text{man_ma3} + \text{coef4} \times \text{man_mi1} + \text{coef5} \times \text{man_mi2} + \text{coef6} \times \text{man_mi3} + \text{coef7} \times \text{man_mi4} + \text{coef8} \times \text{man_mi5} + \text{coef9} \times \text{man_mi6} + \text{coef10} \times \text{man_mi7} + \text{coef11} \times \text{man_bias}$	<ol style="list-style-type: none"> 1) Multiplication to create the situations where one of the symptom in the major criteria met (union of three binomial variables, such as $\text{man_ma1} + \text{man_ma2}$ and $\text{man_ma1} \times \text{man_ma2}$), \n 2) multiplication for the condition of presenting irritable mood (... x man_ma3), and 3) the bias variable (man_bias) equivalent to the residual of the diagnosis not explained by the input symptoms and the bias variables due to censoring; 4) the bias variables introduced by categorizing the number of input symptoms confirmed in the minor criteria (man_bias1 and man_bias2)
			A distinct period of abnormally and persistently elevated, expansive, or irritable mood, lasting at least 1 week (or any duration if hospitalization is necessary)		Elevated mood, lasting at least 1 week	man_ma1		
					Expansive mood, lasting at least 1 week	man_ma2		
					Irritable mood, lasting at least 1 week	man_ma3		
Minor criteria (3 or more of the following symptoms have persisted; 4 if the mood is only irritable)		Increased self-esteem or grandiosity	man_mi1	Increased self-esteem	man_mi1_1	$\text{man_mi1} = \text{man_mi1_1} + \text{man_mi1_2} + \text{man_mi1_bias}$		Censoring of the sum of multiple input variables
				Grandiosity	man_mi1_2			
				Information of the domain not explained by the input variables	man_mi1_bias			

		Decreased need for sleep (e.g., feels rested after only 3 hours of sleep)	man_mi2			
		More talkative than usual or pressure to keep talking	man_mi3	More talkative than usual Pressure to keep talking Information of the domain not explained by the input variables	man_mi3_1 man_mi3_2 man_mi3_bias	man_mi3 = man_mi3_1 + man_mi3_2 + man_mi3_bias Censoring of the sum of multiple input variables
		Flight of ideas or subjective experience that thoughts are racing	man_mi4	Flight of ideas Subjective experience that thoughts are racing Information of the domain not explained by the input variables	man_mi4_1 man_mi4_2 man_mi4_bias	man_mi4 = man_mi4_1 + man_mi4_2 + man_mi4_bias Censoring of the sum of multiple input variables
		Distractibility (i.e., attention too easily drawn to unimportant or irrelevant external stimuli)	man_mi5			
		Increase in goal-directed activity (either socially, at work or school, or sexually) or psychomotor agitation	man_mi6	Increase in goal-directed activity Psychomotor agitation Information of the domain not explained by the input variables	man_mi6_1 man_mi6_2 man_mi6_bias	man_mi6 = man_mi6_1 + man_mi6_2 + man_mi6_bias Censoring of the sum of multiple input variables
		Excessive involvement in pleasurable activities that have a high potential for painful consequences (e.g., engaging in unrestrained buying sprees, sexual indiscretions, or foolish business investments)"	man_mi7			
Information of diagnosis due to categorization (choosing at least three symptoms)	man_bias1					Bias introduced by categorizing the number of input symptoms confirmed in the minor criteria
Information of diagnosis due to categorization (choosing at least four symptoms)	man_bias2					Bias introduced by categorizing the number of input symptoms confirmed in the minor criteria
Information of diagnosis not explained by symptoms	man_bias					Information of the diagnosis not explained by the input symptoms and the bias variables generated due to data categorization, man_bias1 and man_bias2

529 **Table 5. The derived prevalence rates of the diagnoses of major depressive episodes, dysthymic disorder, and manic**
 530 **episodes based on the assumed prevalence rates and between-variable correlations of the input symptoms**

Assumed correlations between input symptoms	Assumed prevalence of input symptoms	Major depressive episodes	Dysthymic disorder	Manic episodes
0	0.05	0 (95% CI = 0 to 0)	0.004 (95% CI = 0.004 to 0.004)	0 (95% CI = 0 to 0)
0	0.1	0.001 (95% CI = 0.001 to 0.001)	0.025 (95% CI = 0.025 to 0.025)	0.002 (95% CI = 0.002 to 0.002)
0	0.3	0.067 (95% CI = 0.067 to 0.067)	0.249 (95% CI = 0.249 to 0.249)	0.136 (95% CI = 0.135 to 0.136)
0	0.5	0.245 (95% CI = 0.244 to 0.245)	0.493 (95% CI = 0.493 to 0.493)	0.436 (95% CI = 0.436 to 0.436)
0	0.7	0.49 (95% CI = 0.49 to 0.49)	0.7 (95% CI = 0.7 to 0.7)	0.692 (95% CI = 0.692 to 0.693)
0.1	0.05	0.004 (95% CI = 0.004 to 0.004)	0.018 (95% CI = 0.018 to 0.018)	0.007 (95% CI = 0.007 to 0.007)
0.1	0.1	0.011 (95% CI = 0.011 to 0.011)	0.049 (95% CI = 0.049 to 0.049)	0.022 (95% CI = 0.021 to 0.022)
0.1	0.3	0.094 (95% CI = 0.094 to 0.094)	0.25 (95% CI = 0.25 to 0.25)	0.172 (95% CI = 0.171 to 0.172)
0.1	0.5	0.267 (95% CI = 0.267 to 0.268)	0.482 (95% CI = 0.482 to 0.482)	0.425 (95% CI = 0.425 to 0.425)
0.1	0.7	0.51 (95% CI = 0.509 to 0.51)	0.697 (95% CI = 0.697 to 0.697)	0.679 (95% CI = 0.679 to 0.679)
0.4	0.05	0.019 (95% CI = 0.019 to 0.019)	0.037 (95% CI = 0.037 to 0.037)	0.029 (95% CI = 0.029 to 0.029)
0.4	0.1	0.042 (95% CI = 0.042 to 0.042)	0.078 (95% CI = 0.078 to 0.078)	0.062 (95% CI = 0.062 to 0.062)
0.4	0.3	0.166 (95% CI = 0.166 to 0.167)	0.267 (95% CI = 0.267 to 0.267)	0.231 (95% CI = 0.231 to 0.231)
0.4	0.5	0.344 (95% CI = 0.344 to 0.344)	0.476 (95% CI = 0.476 to 0.476)	0.44 (95% CI = 0.44 to 0.441)
0.4	0.7	0.57 (95% CI = 0.57 to 0.57)	0.689 (95% CI = 0.688 to 0.689)	0.666 (95% CI = 0.666 to 0.666)
0.7	0.05	0.035 (95% CI = 0.035 to 0.035)	0.046 (95% CI = 0.046 to 0.046)	0.042 (95% CI = 0.042 to 0.042)
0.7	0.1	0.071 (95% CI = 0.071 to 0.071)	0.092 (95% CI = 0.092 to 0.092)	0.085 (95% CI = 0.085 to 0.085)
0.7	0.3	0.233 (95% CI = 0.233 to 0.234)	0.285 (95% CI = 0.285 to 0.285)	0.27 (95% CI = 0.27 to 0.27)
0.7	0.5	0.422 (95% CI = 0.421 to 0.422)	0.486 (95% CI = 0.485 to 0.486)	0.469 (95% CI = 0.468 to 0.469)
0.7	0.7	0.635 (95% CI = 0.635 to 0.635)	0.69 (95% CI = 0.69 to 0.691)	0.678 (95% CI = 0.677 to 0.678)
0.9	0.05	0.042 (95% CI = 0.042 to 0.042)	0.048 (95% CI = 0.048 to 0.048)	0.046 (95% CI = 0.046 to 0.046)
0.9	0.1	0.085 (95% CI = 0.085 to 0.085)	0.096 (95% CI = 0.096 to 0.097)	0.093 (95% CI = 0.093 to 0.093)
0.9	0.3	0.268 (95% CI = 0.268 to 0.268)	0.293 (95% CI = 0.293 to 0.293)	0.286 (95% CI = 0.286 to 0.287)
0.9	0.5	0.463 (95% CI = 0.463 to 0.463)	0.493 (95% CI = 0.492 to 0.493)	0.485 (95% CI = 0.485 to 0.486)
0.9	0.7	0.669 (95% CI = 0.669 to 0.669)	0.695 (95% CI = 0.694 to 0.695)	0.688 (95% CI = 0.688 to 0.688)

531

532

534 **Table 6. The individual input symptoms that best explained the diagnoses: major depressive episodes, dysthymic**
 535 **disorder, and manic episodes**

Assumed correlations between input symptoms	Assumed prevalence of input symptoms	Major depressive episodes	Dysthymic disorder	Manic episodes
0	0.05	mde_ma1	dys_ma	man_ma3
0	0.05	0.001 (95% CI = 0.001 to 0.001)	0.076 (95% CI = 0.075 to 0.077)	0.002 (95% CI = 0.002 to 0.002)
0	0.1	mde_ma1	dys_ma	man_ma3
0	0.1	0.01 (95% CI = 0.01 to 0.01)	0.228 (95% CI = 0.227 to 0.229)	0.021 (95% CI = 0.02 to 0.021)
0	0.3	mde_ma1	dys_ma	man_ma3
0	0.3	0.167 (95% CI = 0.167 to 0.167)	0.774 (95% CI = 0.773 to 0.774)	0.366 (95% CI = 0.366 to 0.367)
0	0.5	mde_ma2	dys_ma	man_ma3
0	0.5	0.324 (95% CI = 0.324 to 0.325)	0.971 (95% CI = 0.971 to 0.971)	0.773 (95% CI = 0.772 to 0.773)
0	0.7	mde_ma2	dys_ma	man_ma3
0	0.7	0.412 (95% CI = 0.412 to 0.412)	0.999 (95% CI = 0.999 to 0.999)	0.964 (95% CI = 0.964 to 0.964)
0.1	0.05	mde_ma2	dys_ma	man_ma3
0.1	0.05	0.07 (95% CI = 0.07 to 0.071)	0.353 (95% CI = 0.352 to 0.355)	0.136 (95% CI = 0.135 to 0.137)
0.1	0.1	mde_ma1	dys_ma	man_ma3
0.1	0.1	0.101 (95% CI = 0.1 to 0.101)	0.462 (95% CI = 0.461 to 0.463)	0.199 (95% CI = 0.198 to 0.199)
0.1	0.3	mde_ma2	dys_ma	man_ma3
0.1	0.3	0.242 (95% CI = 0.242 to 0.243)	0.777 (95% CI = 0.777 to 0.778)	0.483 (95% CI = 0.483 to 0.484)
0.1	0.5	mde_ma2	dys_ma	man_ma3
0.1	0.5	0.365 (95% CI = 0.365 to 0.366)	0.932 (95% CI = 0.931 to 0.932)	0.74 (95% CI = 0.74 to 0.741)
0.1	0.7	mde_ma2	dys_ma	man_ma3
0.1	0.7	0.445 (95% CI = 0.445 to 0.446)	0.986 (95% CI = 0.986 to 0.986)	0.906 (95% CI = 0.906 to 0.907)
0.4	0.05	mde_ma1	dys_ma	man_ma3
0.4	0.05	0.375 (95% CI = 0.373 to 0.376)	0.731 (95% CI = 0.729 to 0.732)	0.561 (95% CI = 0.559 to 0.562)
0.4	0.1	mde_ma1	dys_ma	man_ma3
0.4	0.1	0.395 (95% CI = 0.394 to 0.396)	0.763 (95% CI = 0.762 to 0.764)	0.595 (95% CI = 0.594 to 0.596)
0.4	0.3	mde_ma1	dys_ma	man_ma3
0.4	0.3	0.465 (95% CI = 0.465 to 0.466)	0.851 (95% CI = 0.85 to 0.851)	0.701 (95% CI = 0.701 to 0.702)
0.4	0.5	mde_ma2	dys_ma	man_ma3
0.4	0.5	0.525 (95% CI = 0.524 to 0.525)	0.908 (95% CI = 0.908 to 0.908)	0.787 (95% CI = 0.786 to 0.787)
0.4	0.7	mde_ma2	dys_ma	man_ma3
0.4	0.7	0.568 (95% CI = 0.568 to 0.569)	0.946 (95% CI = 0.946 to 0.947)	0.855 (95% CI = 0.854 to 0.855)
0.7	0.05	mde_ma2	dys_ma	man_ma3
0.7	0.05	0.688 (95% CI = 0.687 to 0.69)	0.909 (95% CI = 0.908 to 0.909)	0.831 (95% CI = 0.83 to 0.832)
0.7	0.1	mde_ma1	dys_ma	man_ma3
0.7	0.1	0.688 (95% CI = 0.687 to 0.689)	0.912 (95% CI = 0.911 to 0.913)	0.836 (95% CI = 0.835 to 0.836)
0.7	0.3	mde_ma2	dys_ma	man_ma3
0.7	0.3	0.71 (95% CI = 0.709 to 0.711)	0.93 (95% CI = 0.93 to 0.93)	0.862 (95% CI = 0.861 to 0.862)
0.7	0.5	mde_ma2	dys_ma	man_ma3
0.7	0.5	0.729 (95% CI = 0.728 to 0.729)	0.944 (95% CI = 0.943 to 0.944)	0.882 (95% CI = 0.882 to 0.883)
0.7	0.7	mde_ma1	dys_ma	man_ma3
0.7	0.7	0.745 (95% CI = 0.744 to 0.745)	0.954 (95% CI = 0.954 to 0.955)	0.9 (95% CI = 0.9 to 0.9)
0.9	0.05	mde_ma1	dys_ma	man_ma3
0.9	0.05	0.828 (95% CI = 0.827 to 0.829)	0.958 (95% CI = 0.957 to 0.958)	0.918 (95% CI = 0.917 to 0.919)
0.9	0.1	mde_ma2	dys_ma	man_ma3
0.9	0.1	0.838 (95% CI = 0.838 to 0.839)	0.961 (95% CI = 0.961 to 0.961)	0.925 (95% CI = 0.924 to 0.925)
0.9	0.3	mde_ma2	dys_ma	man_ma3
0.9	0.3	0.856 (95% CI = 0.856 to 0.857)	0.969 (95% CI = 0.968 to 0.969)	0.937 (95% CI = 0.936 to 0.937)
0.9	0.5	mde_ma2	dys_ma	man_ma3
0.9	0.5	0.862 (95% CI = 0.862 to 0.863)	0.972 (95% CI = 0.972 to 0.972)	0.942 (95% CI = 0.942 to 0.943)
0.9	0.7	mde_ma2	dys_ma	man_ma3
0.9	0.7	0.865 (95% CI = 0.865 to 0.866)	0.974 (95% CI = 0.974 to 0.974)	0.946 (95% CI = 0.946 to 0.946)

536

537

538

539

540

541

542

543

544

545

546

547 Table 7. The individual bias variables that best explained the diagnoses: major depressive episodes,
548 dysthymic disorder, and manic episodes

Assumed correlations between input symptoms	Assumed prevalence of input symptoms	Major depressive episodes	Dysthymic disorder	Manic episodes
0	0.05	mde_bias2	dys_bias	man_bias2
0	0.05	0 (95% CI = 0 to 0)	0.028 (95% CI = 0.028 to 0.028)	0.001 (95% CI = 0.001 to 0.001)
0	0.1	mde_bias2	dys_bias	man_bias2
0	0.1	0.004 (95% CI = 0.004 to 0.004)	0.053 (95% CI = 0.053 to 0.054)	0.011 (95% CI = 0.011 to 0.011)
0	0.3	mde_bias2	dys_bias	man_bias1
0	0.3	0.015 (95% CI = 0.015 to 0.015)	0.045 (95% CI = 0.045 to 0.045)	0.089 (95% CI = 0.089 to 0.09)
0	0.5	mde_bias	dys_bias	man_bias1
0	0.5	0.013 (95% CI = 0.013 to 0.014)	0.007 (95% CI = 0.007 to 0.007)	0.035 (95% CI = 0.034 to 0.035)
0	0.7	mde_bias	dys_bias	man_bias1
0	0.7	0.01 (95% CI = 0.01 to 0.01)	0 (95% CI = 0 to 0)	0.002 (95% CI = 0.002 to 0.002)
0.1	0.05	mde_bias2	dys_bias	man_bias1
0.1	0.05	0.037 (95% CI = 0.037 to 0.037)	0.113 (95% CI = 0.113 to 0.114)	0.083 (95% CI = 0.083 to 0.084)
0.1	0.1	mde_bias2	dys_bias	man_bias1
0.1	0.1	0.047 (95% CI = 0.047 to 0.048)	0.122 (95% CI = 0.121 to 0.122)	0.116 (95% CI = 0.115 to 0.116)
0.1	0.3	mde_bias2	dys_mi_bias	man_bias1
0.1	0.3	0.077 (95% CI = 0.077 to 0.077)	0.105 (95% CI = 0.105 to 0.106)	0.198 (95% CI = 0.197 to 0.198)
0.1	0.5	mde_bias2	dys_mi_bias	man_bias1
0.1	0.5	0.079 (95% CI = 0.079 to 0.08)	0.073 (95% CI = 0.073 to 0.073)	0.166 (95% CI = 0.166 to 0.167)
0.1	0.7	mde_bias2	dys_mi_bias	man_bias1
0.1	0.7	0.065 (95% CI = 0.065 to 0.065)	0.047 (95% CI = 0.046 to 0.047)	0.094 (95% CI = 0.093 to 0.094)
0.4	0.05	mde_bias1	dys_mi_bias	man_bias1
0.4	0.05	0.294 (95% CI = 0.293 to 0.295)	0.415 (95% CI = 0.413 to 0.416)	0.432 (95% CI = 0.431 to 0.433)
0.4	0.1	mde_bias1	dys_mi_bias	man_bias1
0.4	0.1	0.304 (95% CI = 0.303 to 0.304)	0.419 (95% CI = 0.418 to 0.42)	0.445 (95% CI = 0.444 to 0.445)
0.4	0.3	mde_bias1	dys_mi_bias	man_bias1
0.4	0.3	0.335 (95% CI = 0.334 to 0.335)	0.411 (95% CI = 0.411 to 0.412)	0.473 (95% CI = 0.472 to 0.473)
0.4	0.5	mde_bias1	dys_mi_bias	man_bias1
0.4	0.5	0.354 (95% CI = 0.354 to 0.355)	0.395 (95% CI = 0.395 to 0.396)	0.475 (95% CI = 0.474 to 0.475)
0.4	0.7	mde_bias1	dys_mi_bias	man_bias1
0.4	0.7	0.356 (95% CI = 0.355 to 0.356)	0.367 (95% CI = 0.366 to 0.367)	0.451 (95% CI = 0.45 to 0.451)
0.7	0.05	mde_bias1	dys_mi_bias	man_bias1
0.7	0.05	0.616 (95% CI = 0.615 to 0.617)	0.705 (95% CI = 0.704 to 0.706)	0.723 (95% CI = 0.722 to 0.724)
0.7	0.1	mde_bias1	dys_mi_bias	man_bias1
0.7	0.1	0.611 (95% CI = 0.611 to 0.612)	0.699 (95% CI = 0.698 to 0.699)	0.72 (95% CI = 0.72 to 0.721)
0.7	0.3	mde_bias1	dys_mi_bias	man_bias1
0.7	0.3	0.623 (95% CI = 0.623 to 0.624)	0.699 (95% CI = 0.699 to 0.7)	0.728 (95% CI = 0.728 to 0.729)
0.7	0.5	mde_bias1	dys_mi_bias	man_bias1
0.7	0.5	0.632 (95% CI = 0.632 to 0.633)	0.696 (95% CI = 0.696 to 0.697)	0.731 (95% CI = 0.731 to 0.732)
0.7	0.7	mde_bias1	dys_mi_bias	man_bias1
0.7	0.7	0.639 (95% CI = 0.638 to 0.639)	0.693 (95% CI = 0.692 to 0.693)	0.732 (95% CI = 0.731 to 0.732)
0.9	0.05	mde_bias1	dys_mi_bias	man_bias1
0.9	0.05	0.777 (95% CI = 0.776 to 0.778)	0.835 (95% CI = 0.834 to 0.835)	0.847 (95% CI = 0.847 to 0.848)
0.9	0.1	mde_bias1	dys_mi_bias	man_bias1
0.9	0.1	0.788 (95% CI = 0.788 to 0.789)	0.842 (95% CI = 0.841 to 0.843)	0.855 (95% CI = 0.854 to 0.855)
0.9	0.3	mde_bias1	dys_mi_bias	man_bias1
0.9	0.3	0.807 (95% CI = 0.806 to 0.807)	0.854 (95% CI = 0.853 to 0.854)	0.867 (95% CI = 0.867 to 0.868)
0.9	0.5	mde_bias1	dys_mi_bias	man_bias1
0.9	0.5	0.811 (95% CI = 0.811 to 0.811)	0.855 (95% CI = 0.855 to 0.856)	0.87 (95% CI = 0.87 to 0.871)
0.9	0.7	mde_bias1	dys_mi_bias	man_bias1
0.9	0.7	0.812 (95% CI = 0.811 to 0.812)	0.853 (95% CI = 0.853 to 0.853)	0.869 (95% CI = 0.869 to 0.87)

549

550

551

553 Table 8. Approximating the diagnoses using input symptoms and derived adjusted R-squared

Assumed correlations between input symptoms	Assumed prevalence of input symptoms	Major depressive episodes	Dysthymic disorder	Manic episodes
0	0.05	0.003 (95% CI = 0.002 to 0.003)	0.122 (95% CI = 0.121 to 0.123)	0.004 (95% CI = 0.004 to 0.005)
0	0.1	0.024 (95% CI = 0.023 to 0.024)	0.305 (95% CI = 0.304 to 0.306)	0.039 (95% CI = 0.038 to 0.039)
0	0.3	0.348 (95% CI = 0.348 to 0.349)	0.842 (95% CI = 0.841 to 0.842)	0.483 (95% CI = 0.482 to 0.483)
0	0.5	0.649 (95% CI = 0.649 to 0.649)	0.986 (95% CI = 0.986 to 0.986)	0.817 (95% CI = 0.817 to 0.817)
0	0.7	0.823 (95% CI = 0.823 to 0.823)	1 (95% CI = 1 to 1)	0.967 (95% CI = 0.967 to 0.967)
0.1	0.05	0.143 (95% CI = 0.141 to 0.144)	0.435 (95% CI = 0.433 to 0.436)	0.212 (95% CI = 0.211 to 0.213)
0.1	0.1	0.198 (95% CI = 0.197 to 0.199)	0.539 (95% CI = 0.538 to 0.54)	0.29 (95% CI = 0.289 to 0.291)
0.1	0.3	0.45 (95% CI = 0.45 to 0.451)	0.826 (95% CI = 0.826 to 0.827)	0.588 (95% CI = 0.588 to 0.589)
0.1	0.5	0.663 (95% CI = 0.663 to 0.664)	0.952 (95% CI = 0.952 to 0.952)	0.799 (95% CI = 0.799 to 0.799)
0.1	0.7	0.809 (95% CI = 0.809 to 0.809)	0.991 (95% CI = 0.991 to 0.991)	0.922 (95% CI = 0.922 to 0.922)
0.4	0.05	0.587 (95% CI = 0.585 to 0.588)	0.782 (95% CI = 0.781 to 0.783)	0.675 (95% CI = 0.674 to 0.676)
0.4	0.1	0.607 (95% CI = 0.606 to 0.608)	0.807 (95% CI = 0.807 to 0.808)	0.698 (95% CI = 0.697 to 0.698)
0.4	0.3	0.688 (95% CI = 0.688 to 0.689)	0.878 (95% CI = 0.877 to 0.878)	0.775 (95% CI = 0.774 to 0.775)
0.4	0.5	0.761 (95% CI = 0.761 to 0.762)	0.925 (95% CI = 0.924 to 0.925)	0.838 (95% CI = 0.838 to 0.838)
0.4	0.7	0.821 (95% CI = 0.821 to 0.822)	0.956 (95% CI = 0.956 to 0.956)	0.887 (95% CI = 0.887 to 0.888)
0.7	0.05	0.813 (95% CI = 0.812 to 0.814)	0.925 (95% CI = 0.925 to 0.926)	0.877 (95% CI = 0.877 to 0.878)
0.7	0.1	0.826 (95% CI = 0.826 to 0.827)	0.928 (95% CI = 0.927 to 0.928)	0.881 (95% CI = 0.881 to 0.882)
0.7	0.3	0.86 (95% CI = 0.86 to 0.86)	0.942 (95% CI = 0.942 to 0.942)	0.9 (95% CI = 0.9 to 0.9)
0.7	0.5	0.88 (95% CI = 0.88 to 0.88)	0.953 (95% CI = 0.953 to 0.953)	0.913 (95% CI = 0.913 to 0.913)
0.7	0.7	0.895 (95% CI = 0.895 to 0.895)	0.962 (95% CI = 0.962 to 0.962)	0.925 (95% CI = 0.925 to 0.925)
0.9	0.05	0.903 (95% CI = 0.903 to 0.904)	0.965 (95% CI = 0.965 to 0.966)	0.941 (95% CI = 0.94 to 0.941)
0.9	0.1	0.91 (95% CI = 0.91 to 0.911)	0.968 (95% CI = 0.968 to 0.968)	0.945 (95% CI = 0.945 to 0.945)
0.9	0.3	0.923 (95% CI = 0.923 to 0.923)	0.974 (95% CI = 0.974 to 0.974)	0.954 (95% CI = 0.953 to 0.954)
0.9	0.5	0.928 (95% CI = 0.928 to 0.928)	0.976 (95% CI = 0.976 to 0.977)	0.958 (95% CI = 0.957 to 0.958)
0.9	0.7	0.932 (95% CI = 0.932 to 0.932)	0.978 (95% CI = 0.978 to 0.978)	0.96 (95% CI = 0.96 to 0.96)

554

555

557 Table 9. Approximating the diagnoses using bias variables and derived R-squared

Assumed correlations between input symptoms	Assumed prevalence of input symptoms	Major depressive episodes	Dysthymic disorder	Manic episodes
0	0.05	0.003 (95% CI = 0.002 to 0.003)	0.029 (95% CI = 0.029 to 0.03)	0.004 (95% CI = 0.004 to 0.004)
0	0.1	0.013 (95% CI = 0.012 to 0.013)	0.056 (95% CI = 0.056 to 0.056)	0.017 (95% CI = 0.017 to 0.017)
0	0.3	0.083 (95% CI = 0.083 to 0.083)	0.047 (95% CI = 0.047 to 0.047)	0.098 (95% CI = 0.098 to 0.099)
0	0.5	0.111 (95% CI = 0.111 to 0.112)	0.007 (95% CI = 0.007 to 0.007)	0.039 (95% CI = 0.038 to 0.039)
0	0.7	0.095 (95% CI = 0.095 to 0.095)	0 (95% CI = 0 to 0)	0.012 (95% CI = 0.012 to 0.013)
0.1	0.05	0.083 (95% CI = 0.082 to 0.084)	0.145 (95% CI = 0.144 to 0.146)	0.126 (95% CI = 0.125 to 0.127)
0.1	0.1	0.096 (95% CI = 0.095 to 0.097)	0.156 (95% CI = 0.155 to 0.156)	0.154 (95% CI = 0.153 to 0.154)
0.1	0.3	0.145 (95% CI = 0.144 to 0.145)	0.139 (95% CI = 0.138 to 0.139)	0.216 (95% CI = 0.216 to 0.216)
0.1	0.5	0.172 (95% CI = 0.172 to 0.173)	0.097 (95% CI = 0.097 to 0.097)	0.182 (95% CI = 0.181 to 0.182)
0.1	0.7	0.175 (95% CI = 0.175 to 0.175)	0.065 (95% CI = 0.064 to 0.065)	0.115 (95% CI = 0.115 to 0.116)
0.4	0.05	0.421 (95% CI = 0.419 to 0.423)	0.455 (95% CI = 0.453 to 0.456)	0.505 (95% CI = 0.504 to 0.506)
0.4	0.1	0.422 (95% CI = 0.421 to 0.423)	0.454 (95% CI = 0.453 to 0.455)	0.507 (95% CI = 0.506 to 0.508)
0.4	0.3	0.435 (95% CI = 0.434 to 0.435)	0.442 (95% CI = 0.442 to 0.443)	0.512 (95% CI = 0.512 to 0.513)
0.4	0.5	0.452 (95% CI = 0.452 to 0.453)	0.427 (95% CI = 0.427 to 0.427)	0.506 (95% CI = 0.505 to 0.506)
0.4	0.7	0.46 (95% CI = 0.459 to 0.46)	0.403 (95% CI = 0.402 to 0.403)	0.481 (95% CI = 0.481 to 0.482)
0.7	0.05	0.728 (95% CI = 0.727 to 0.729)	0.729 (95% CI = 0.728 to 0.731)	0.764 (95% CI = 0.763 to 0.765)
0.7	0.1	0.722 (95% CI = 0.721 to 0.723)	0.723 (95% CI = 0.722 to 0.724)	0.76 (95% CI = 0.759 to 0.761)
0.7	0.3	0.726 (95% CI = 0.726 to 0.727)	0.722 (95% CI = 0.722 to 0.723)	0.761 (95% CI = 0.761 to 0.762)
0.7	0.5	0.732 (95% CI = 0.731 to 0.732)	0.72 (95% CI = 0.719 to 0.72)	0.76 (95% CI = 0.76 to 0.761)
0.7	0.7	0.737 (95% CI = 0.736 to 0.737)	0.717 (95% CI = 0.716 to 0.717)	0.758 (95% CI = 0.758 to 0.759)
0.9	0.05	0.852 (95% CI = 0.851 to 0.853)	0.85 (95% CI = 0.849 to 0.851)	0.871 (95% CI = 0.871 to 0.872)
0.9	0.1	0.86 (95% CI = 0.859 to 0.861)	0.857 (95% CI = 0.856 to 0.857)	0.876 (95% CI = 0.876 to 0.877)
0.9	0.3	0.872 (95% CI = 0.871 to 0.872)	0.867 (95% CI = 0.867 to 0.868)	0.886 (95% CI = 0.886 to 0.886)
0.9	0.5	0.874 (95% CI = 0.874 to 0.875)	0.869 (95% CI = 0.868 to 0.869)	0.888 (95% CI = 0.887 to 0.888)
0.9	0.7	0.874 (95% CI = 0.874 to 0.875)	0.867 (95% CI = 0.866 to 0.867)	0.886 (95% CI = 0.886 to 0.886)

1
2
3 559 **Figure 1. The assumed and derived prevalence rates for major depressive episodes, dysthymic disorder, and manic**
4 560 **episodes.**

5
6 561

7
8 562 Note: each of the combinations of assumed prevalence rates and between-variable correlations of
9 563 the input symptoms for major depressive episodes, dysthymic disorder, and manic episodes was
10 564 represented by one circle, but the circles overlapped in the graph.

11
12
13 565
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

1
2
3 567 Figure 2. The prevalence rates of an intermediate variable for the diagnosis of major depressive
4 568 episodes.

5
6
7 569

8
9 570 Note: The intermediate variable is “significant unintentional weight loss or gain” and the input
10 571 symptoms are “significant unintentional weight loss” and “significant unintentional weight gain”.
11 572 The black line represents the situation where the prevalence rates of the input symptoms are the
12 573 same as that of the intermediate variable. Lines above the black lines have prevalence rates larger
13 574 than those of the input symptoms.

14
15
16 575
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

1
2
3 577 Figure 3. The prevalence rates of dysthymic disorder.
4
5
6 578

7 579 Note: Dysthymic disorder is diagnosed when both the major (depressed mood most of the day for
8 580 more days than not, for at least 2 years) and minor criteria (at least two of the six items) are
9 581 confirmed. The black line represents the situation where the prevalence rates of the input
10 582 symptoms are the same as those of the intermediate variable. Lines below the black lines have
11 583 prevalence rates lower than those of the input symptoms.
12
13
14

15 584

16
17 585

18
19 586
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

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

587 **Figure 4. The prevalence rates of major depressive episodes.**

588

589 Note: Major depressive episodes are diagnosed when both major and minor criteria are confirmed.
590 The black line represents the situation where the prevalence rates of the input symptoms are the
591 same as that of the intermediate variable. Lines below the black lines have prevalence rates lower
592 than those of the input symptoms.

593

594

595

For peer review only

1
2
3 596 **Figure 5. The prevalence rates of manic episodes**
4

5 597

6
7 598 Note: Manic episodes are diagnosed when the symptoms present as described in the diagnostic
8 599 manual. The black line represents the situation where the prevalence rates of the input symptoms
9 600 are the same as those of the input symptoms. Lines below the black lines have prevalence rates
10 601 lower than those of the input symptoms.
11
12

13 602

14
15 603

16
17 604
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

1
2
3 605 **Figure 6. The approximation of the diagnosis of dysthymic disorder by the input symptoms, the bias variables, and both,**
4 606 **measured by R-squared**

5
6 607

7
8 608 Note: the diagnosis of dysthymic disorder is approximated by the input symptoms, the bias
9 609 variables, and both using forward-stepwise regression. The selection of the variables was
10 610 determined by adjusted R-squared. See Table 4 for the details in the input symptoms and the bias
11 611 variables. The assumed correlations between the input symptoms are 0.4 and the assumed
12 612 prevalence rates of the input symptoms are 0.7.
13 613

14
15 614

16
17 615

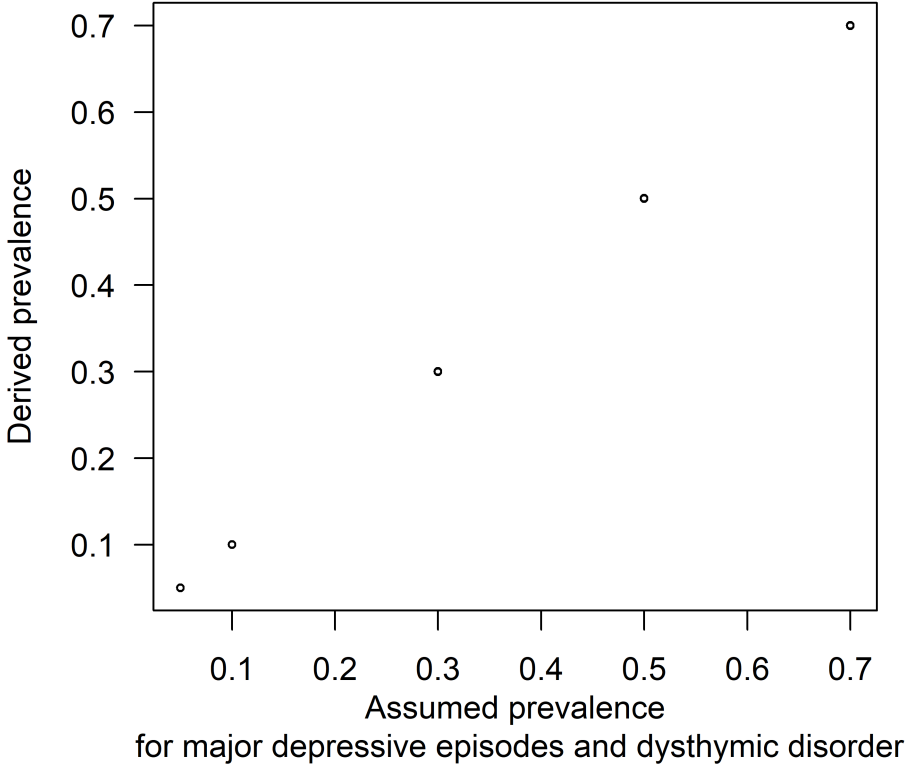
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

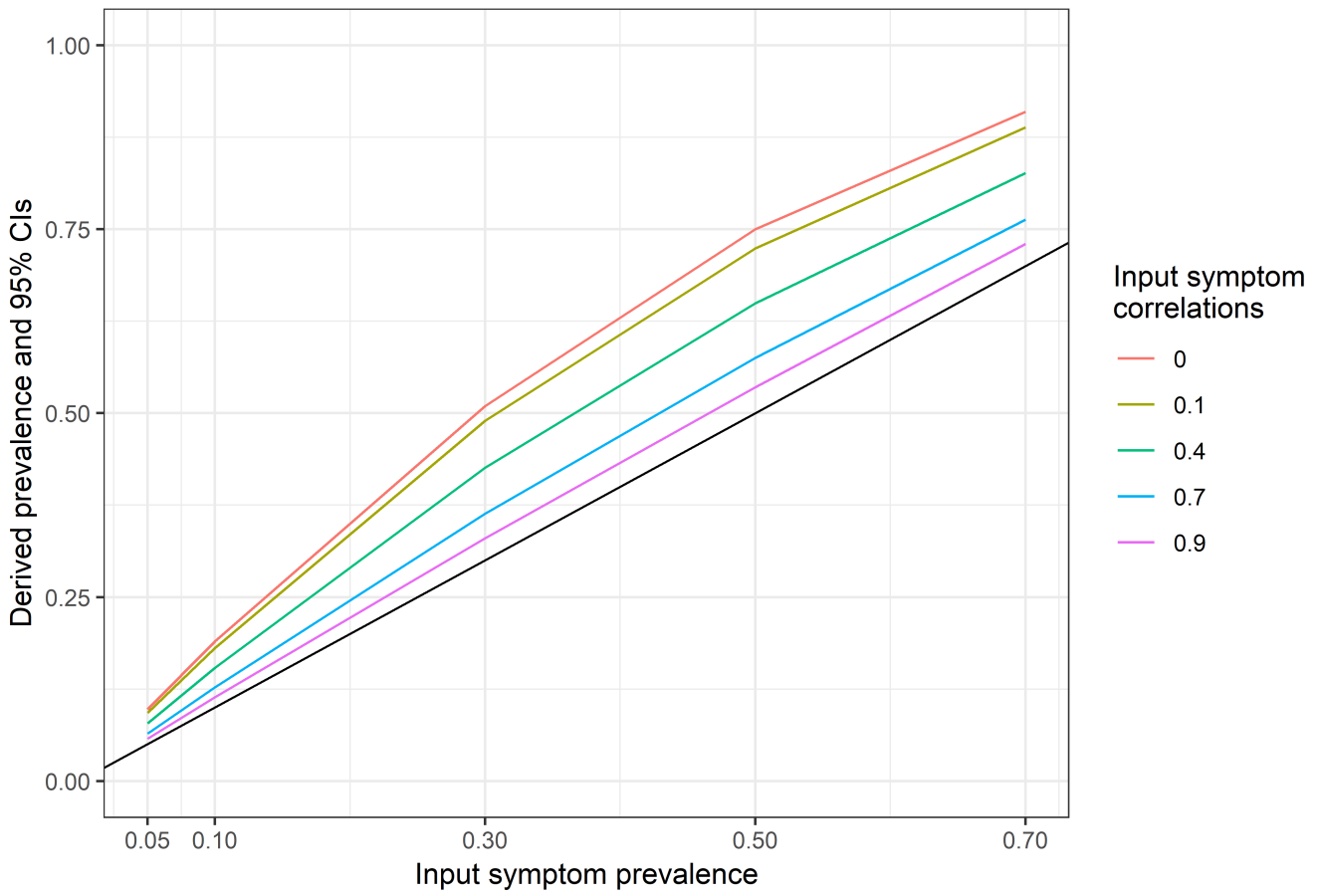
1	
2	
3	
4	616 Supplemental materials
5	617 S 1. Characteristics of the input symptoms for simulations
6	
7	618 S 2. R codes to be used with S1 to simulate populations
8	
9	619 S 3. Correlations between the symptoms
10	
11	620
12	
13	621
14	
15	622
16	
17	
18	
19	
20	
21	
22	
23	
24	
25	
26	
27	
28	
29	
30	
31	
32	
33	
34	
35	
36	
37	
38	
39	
40	
41	
42	
43	
44	
45	
46	
47	
48	
49	
50	
51	
52	
53	
54	
55	
56	
57	
58	
59	
60	

For peer review only

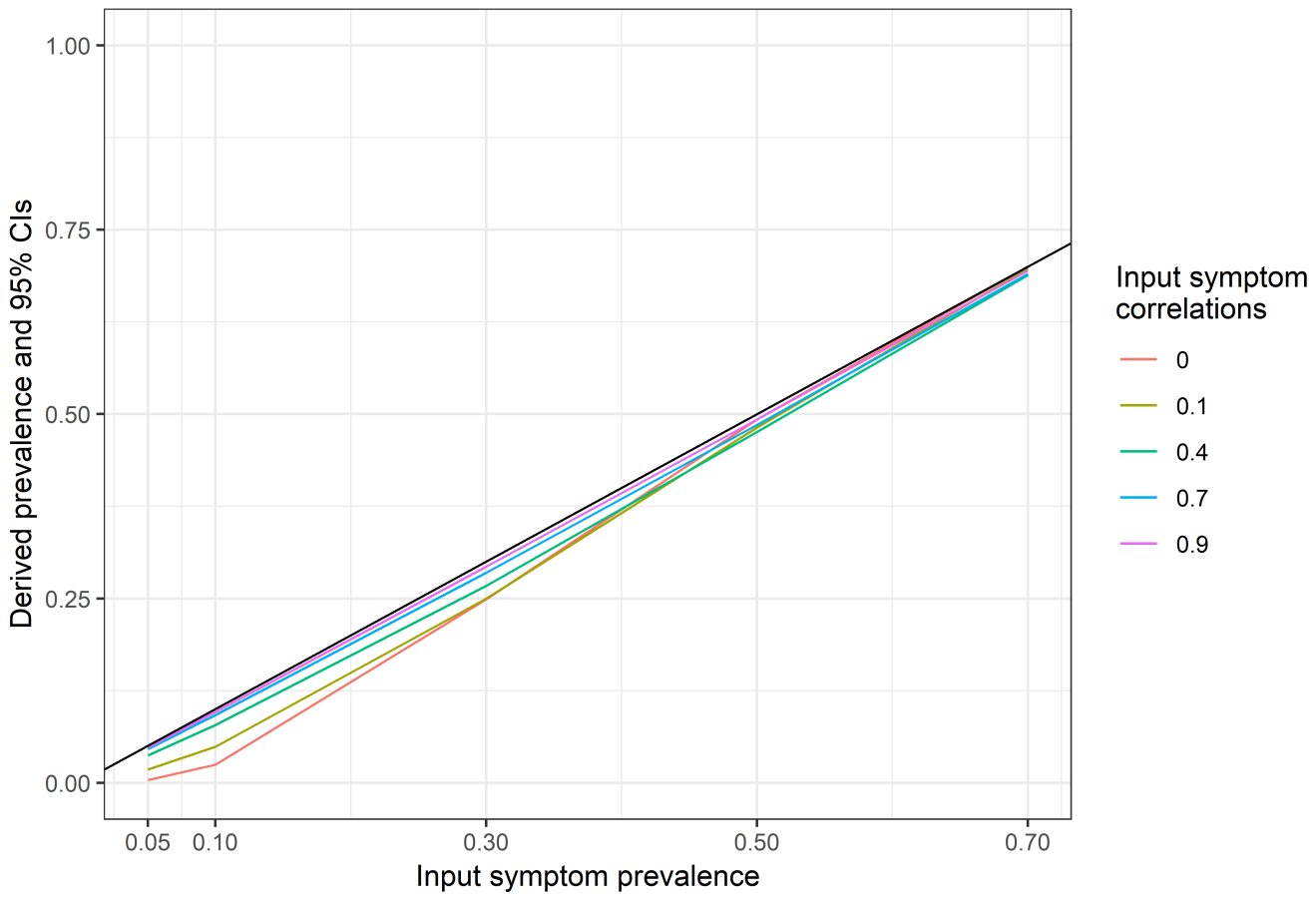
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



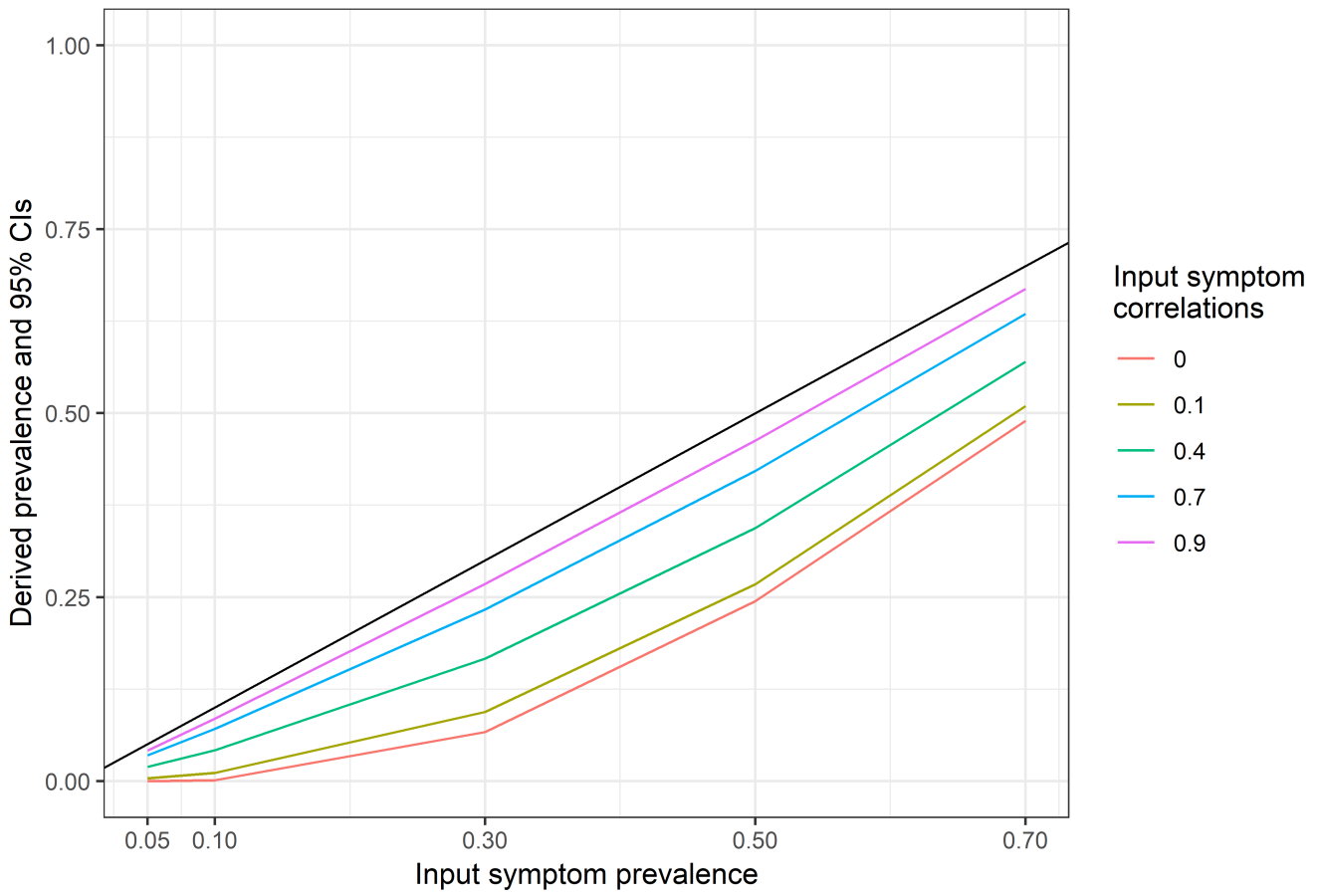
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



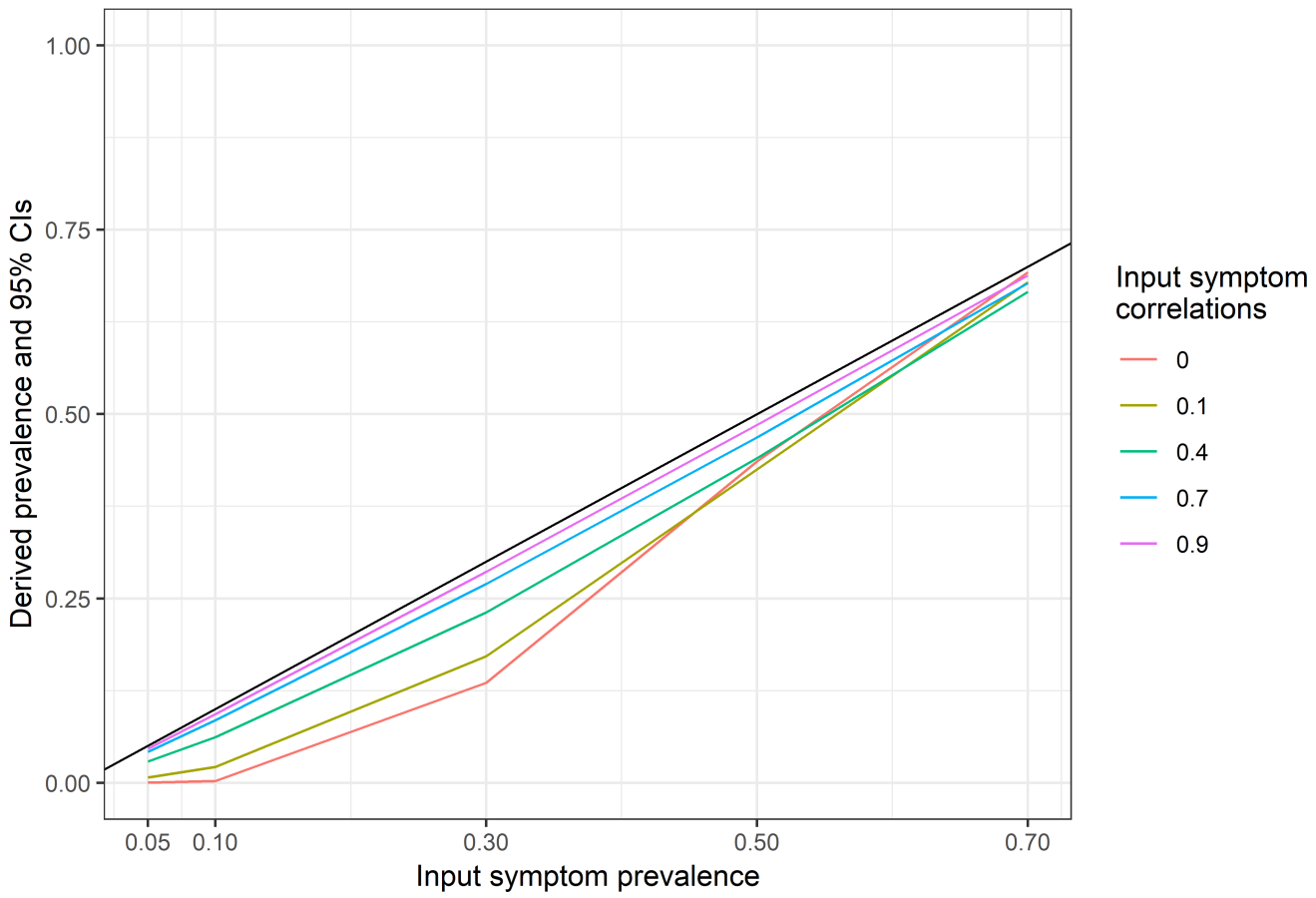
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

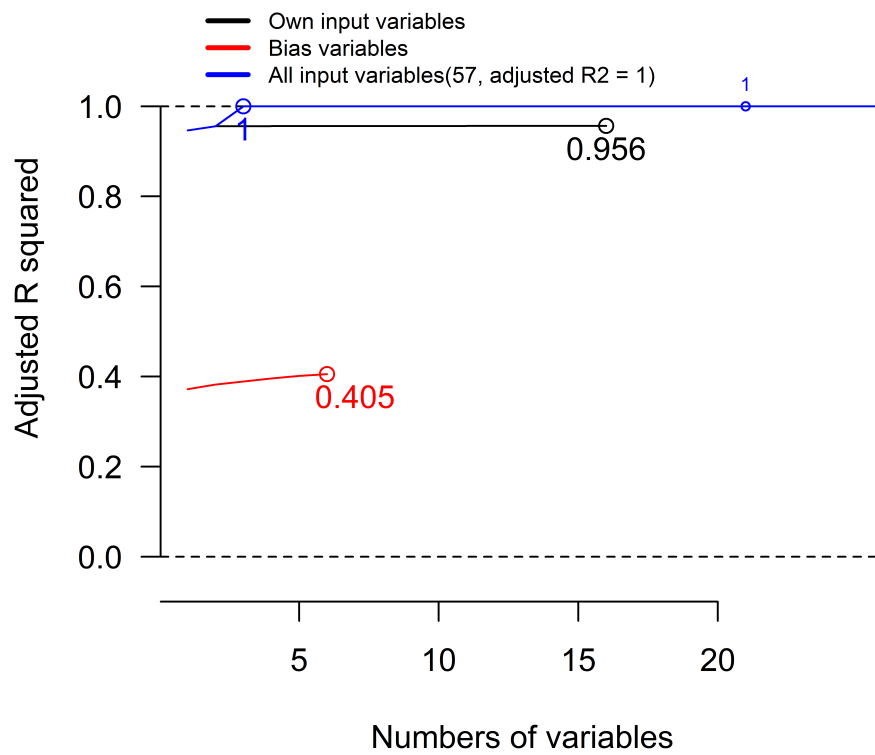


1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



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





1
2
3
4 1 A simulation study to demonstrate the
5
6
7 2 biases in the diagnoses of mental
8
9
10 3 illnesses: major depressive episodes,
11
12 4 dysthymia, and manic episodes
13
14

15 5 **Yi-Sheng Chao^{1*}, Kuan-Fu Lin,² Chao-Jung Wu³, Hsing-Chien Wu⁴, Hui-Ting**
16 6 **Hsu⁵, Lien-Cheng Tsao⁵, Yen-Po Cheng⁵, Yi-Chun Lai⁶, Wei-Chih Chen^{7,8}**

17
18 7 *¹Independent researcher, Montréal, H2X 0A8 Canada, ²National Taiwan*
19 8 *University Hospital Yun-Lin Branch, Yunlin County, 640 Taiwan, ³Département*
20 9 *d'informatique Université du Québec à Montréal, Montréal H3B 1B4 Canada,*
21 10 *⁴Taipei Hospital Ministry of Health and Welfare New Taipei city, 242 Taiwan,*
22 11 *⁵Changhua Christian Hospital, Changhua County 526, Taiwan, ⁶National Yang-*
23 12 *Ming University Hospital, Yilan 260 Taiwan, ⁷Department of Chest Medicine,*
24 13 *Taipei Veterans General Hospital, Taipei 112, Taiwan, ⁸Institute of Emergency*
25 14 *and Critical Care Medicine, National Yang-Ming University, Taipei 112, Taiwan*
26 15 **chaoyisheng@post.harvard.edu*

27
28 16 **Keywords:** Frailty; bias; forward-stepwise regression; the Health and
29 17 Retirement Study; index mining
30 18
31
32 19

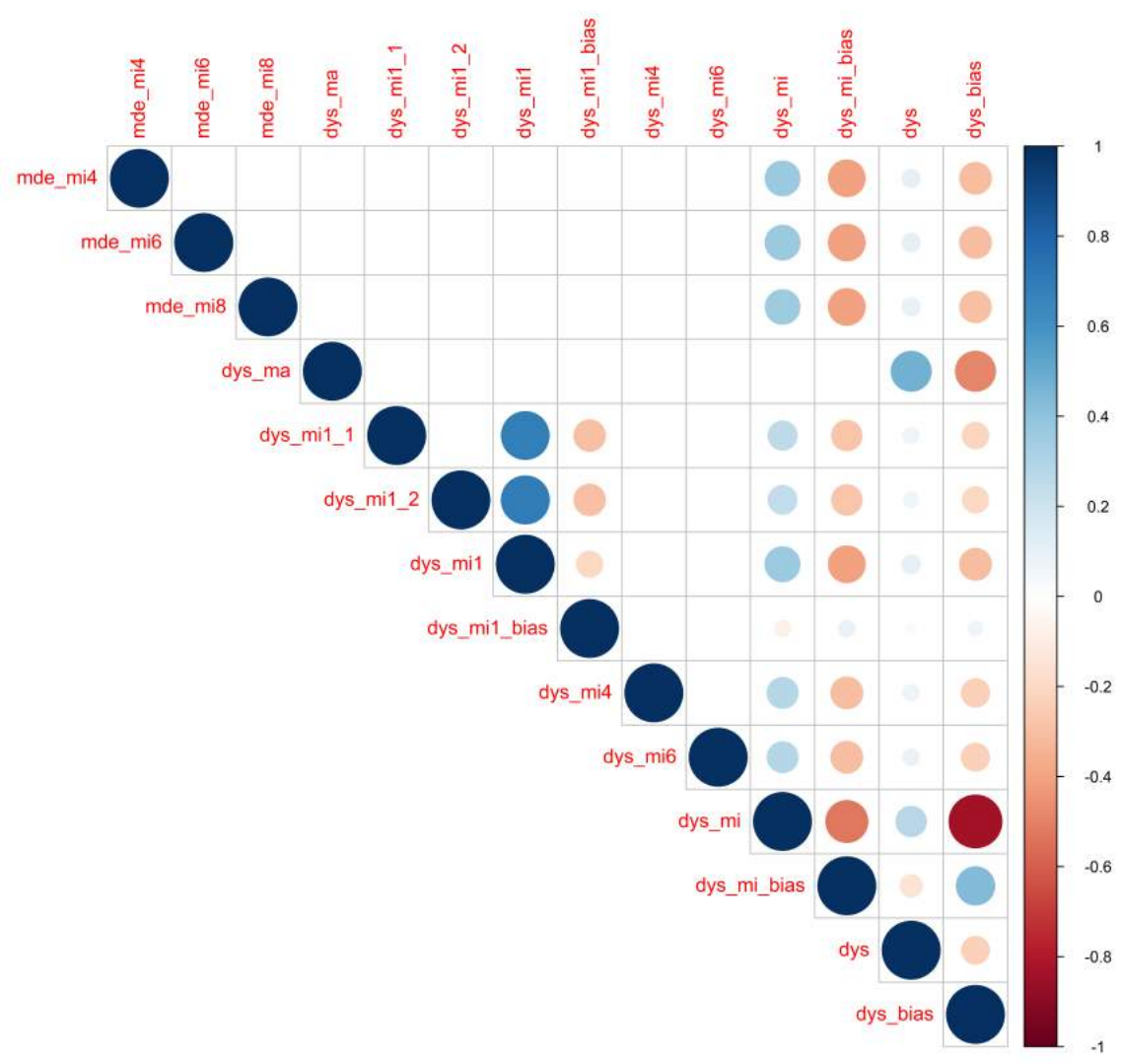

```

1
2
3
4 title: "2019_09_06 simulated mental illnesses"
5 author: "Yi-Sheng Chao"
6 date: "November 22, 2018"
7 output: pdf_document
8 editor_options:
9   chunk_output_type: inline
10
11
12
13 ##Adding correlations to the random variables
14
15 ```{r}
16 library(bindata)
17
18 library(openxlsx)
19 resu = read.xlsx("A simulation study to demonstrate the biases in three
20 diagnoses of mental illnesses.xlsx", sheet = "Prob 1")
21 names(resu)
22 unique(resu$variable)
23 memory.limit(size = 10^13)
24 ssize = 10^5
25 times = 10^2
26
27 prevalence = c(0.05, 0.1, 0.3, 0.5, 0.7)
28 rho = c(0, 0.1, 0.4, 0.7, 0.9)#correlation coefficients of the input
29 symptoms
30
31 collect = c("mean", "max",
32 "min", "derivedprevalence", "coef", "coefse", "p", "intercept",
33 "interceptp", "r2", "subcoef", "subcoefse", "subp", "subintercept",
34 "subinterceptp", "subr2", "appbyownr2", "appbybiasr2", "appbyallr2",
35 "appbyownvar", "appbybiasvar", "appbyallvar", "appbyownn", "appbybiasn",
36 "appbyalln")
37
38
39 set.seed(1)
40
41
42 ##Create a simulated data set to extract variables
43 for(preval in 1:length(prevalence)){
44   for(rh in 1:length(rho)){
45
46
47     library(openxlsx)
48     resu = read.xlsx("A simulation study to demonstrate the biases in three
49 diagnoses of mental illnesses.xlsx", sheet = "Prob 1")
50
51     # foreach(c = 1:times) %dopar% {
52     for(c in 1:times){
53
54       library(bindata)
55       bindata = as.data.frame(rmvbin(ssize, rep(prevalence[preval], 40),
56 bincorr=(1 - rho[rh])*diag(40) + rho[rh]))
57       bindata2 = as.data.frame(rmvbin(ssize, rep(prevalence[preval], 20),
58 bincorr=(1 - rho[rh])*diag(20) + rho[rh]))
59
60       ##demographic characteristics

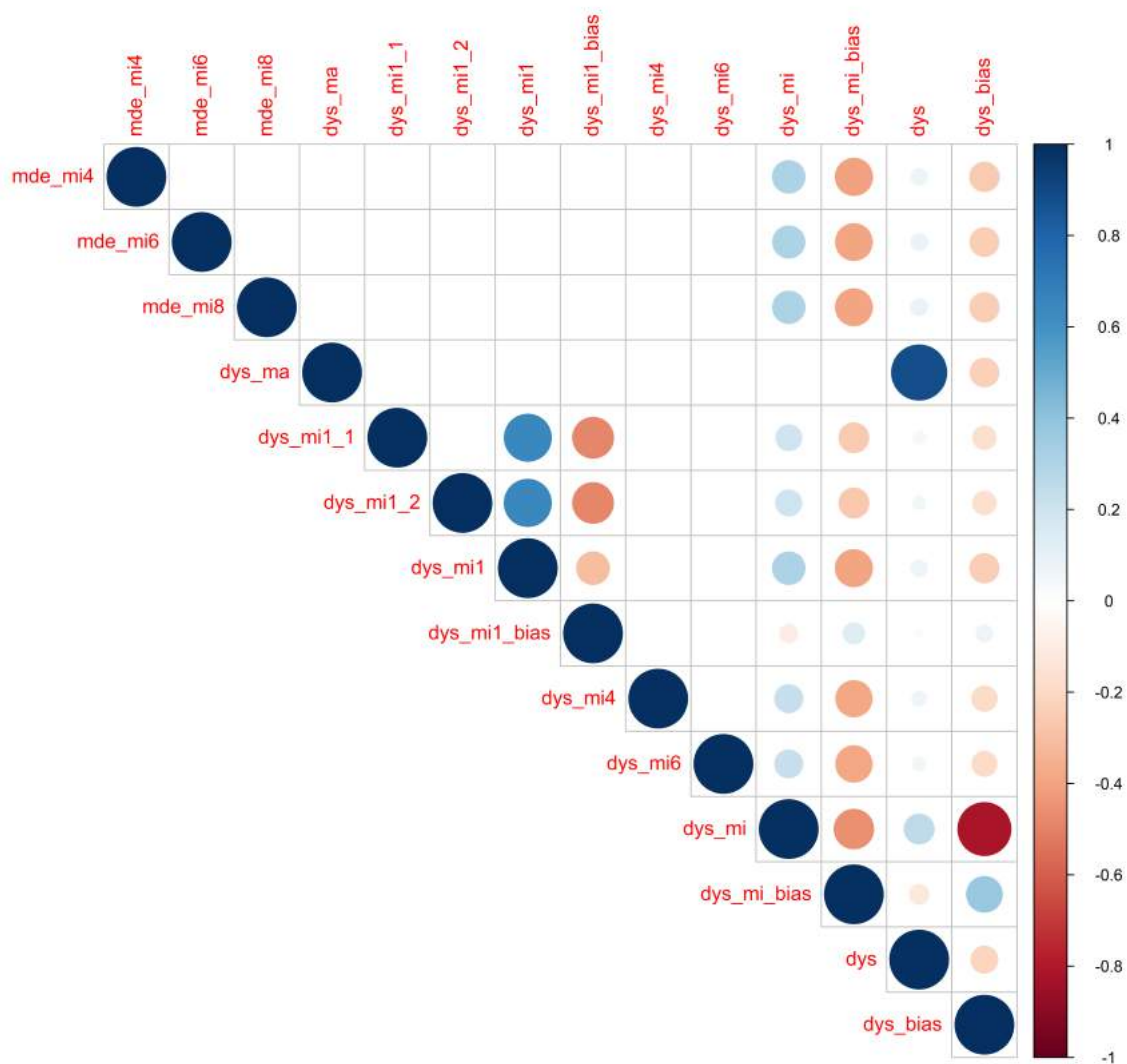
```

```
1
2
3   sim = data.frame(1:ssize)
4   names(sim) = "id"
5   sim$female = rbinom(n = ssize, size = 1, prob = 0.51)
6   sim$age = sample(30:60, ssize, replace = TRUE)
7   sim$edu = rnorm(ssize, mean = 12, sd = 5)
8   sim$edu[which(sim$edu <= 0)] = 0
9   sim$id = NULL
10
11
12   sim$mde_ma1 = bindata[,1]
13   sim$mde_ma2 = bindata[,2]
14
15   sim$mde_mi3_1 = bindata[,3]
16   sim$mde_mi3_2 = bindata[,4]
17   sim$mde_mi3 = 1*((sim$mde_mi3_1 + sim$mde_mi3_2) > 0)
18   sim$mde_mi3_bias = sim$mde_mi3 - sim$mde_mi3_1 - sim$mde_mi3_2
19
20   sim$mde_mi4_1 = bindata[,5]
21   sim$mde_mi4_2 = bindata[,6]
22   sim$mde_mi4 = 1*((sim$mde_mi4_1 + sim$mde_mi4_2) > 0)
23   sim$mde_mi4_bias = sim$mde_mi4 - sim$mde_mi4_1 - sim$mde_mi4_2
24
25   sim$mde_mi5_1 = bindata[,7]
26   sim$mde_mi5_2 = bindata[,8]
27   sim$mde_mi5 = 1*((sim$mde_mi5_1 + sim$mde_mi5_2) > 0)
28   sim$mde_mi5_bias = sim$mde_mi5 - sim$mde_mi5_1 - sim$mde_mi5_2
29
30   sim$mde_mi6_1 = bindata[,9]
31   sim$mde_mi6_2 = bindata[,10]
32   sim$mde_mi6 = 1*((sim$mde_mi6_1 + sim$mde_mi6_2) > 0)
33   sim$mde_mi6_bias = sim$mde_mi6 - sim$mde_mi6_1 - sim$mde_mi6_2
34
35   sim$mde_mi7_1 = bindata[,11]
36   sim$mde_mi7_2 = bindata[,12]
37   sim$mde_mi7 = 1*((sim$mde_mi7_1 + sim$mde_mi7_2) > 0)
38   sim$mde_mi7_bias = sim$mde_mi7 - sim$mde_mi7_1 - sim$mde_mi7_2
39
40   sim$mde_mi8_1 = bindata[,13]
41   sim$mde_mi8_2 = bindata[,14]
42   sim$mde_mi8 = 1*((sim$mde_mi8_1 + sim$mde_mi8_2) > 0)
43   sim$mde_mi8_bias = sim$mde_mi8 - sim$mde_mi8_1 - sim$mde_mi8_2
44
45   sim$mde_mi9 = bindata[,15]
46
47   sim$mde_bias1 = 1 * ((sim$mde_mi3 + sim$mde_mi4 + sim$mde_mi5 +
48   sim$mde_mi6 + sim$mde_mi7 + sim$mde_mi8 + sim$mde_mi9)>2) - (sim$mde_mi3
49   + sim$mde_mi4 + sim$mde_mi5 + sim$mde_mi6 + sim$mde_mi7 + sim$mde_mi8 +
50   sim$mde_mi9)
51   sim$mde_bias2 = 1 * ((sim$mde_mi3 + sim$mde_mi4 + sim$mde_mi5 +
52   sim$mde_mi6 + sim$mde_mi7 + sim$mde_mi8 + sim$mde_mi9)>3) - (sim$mde_mi3
53   + sim$mde_mi4 + sim$mde_mi5 + sim$mde_mi6 + sim$mde_mi7 + sim$mde_mi8 +
54   sim$mde_mi9)
55
56   sim$mde = sim$mde_ma1 * sim$mde_ma2 * (sim$mde_mi3 + sim$mde_mi4 +
57   sim$mde_mi5 + sim$mde_mi6 + sim$mde_mi7 + sim$mde_mi8 + sim$mde_mi9 +
58   sim$mde_bias1) + (1- sim$mde_ma1 * sim$mde_ma2) * (sim$mde_ma1 *
59   sim$mde_ma2) * (sim$mde_mi3 + sim$mde_mi4 + sim$mde_mi5 + sim$mde_mi6 +
60   sim$mde_mi7 + sim$mde_mi8 + sim$mde_mi9 + sim$mde_bias2)
```

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

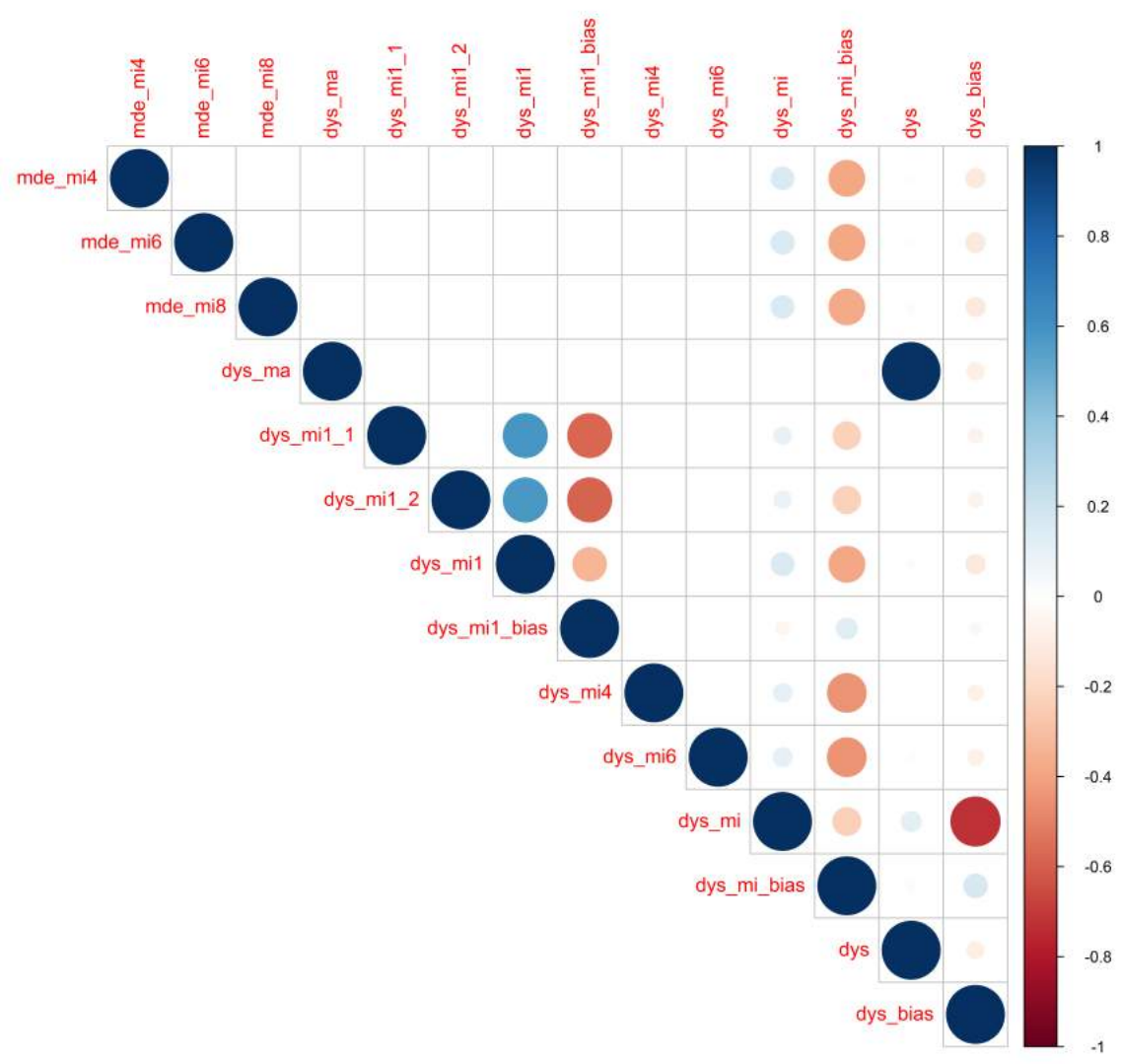


only

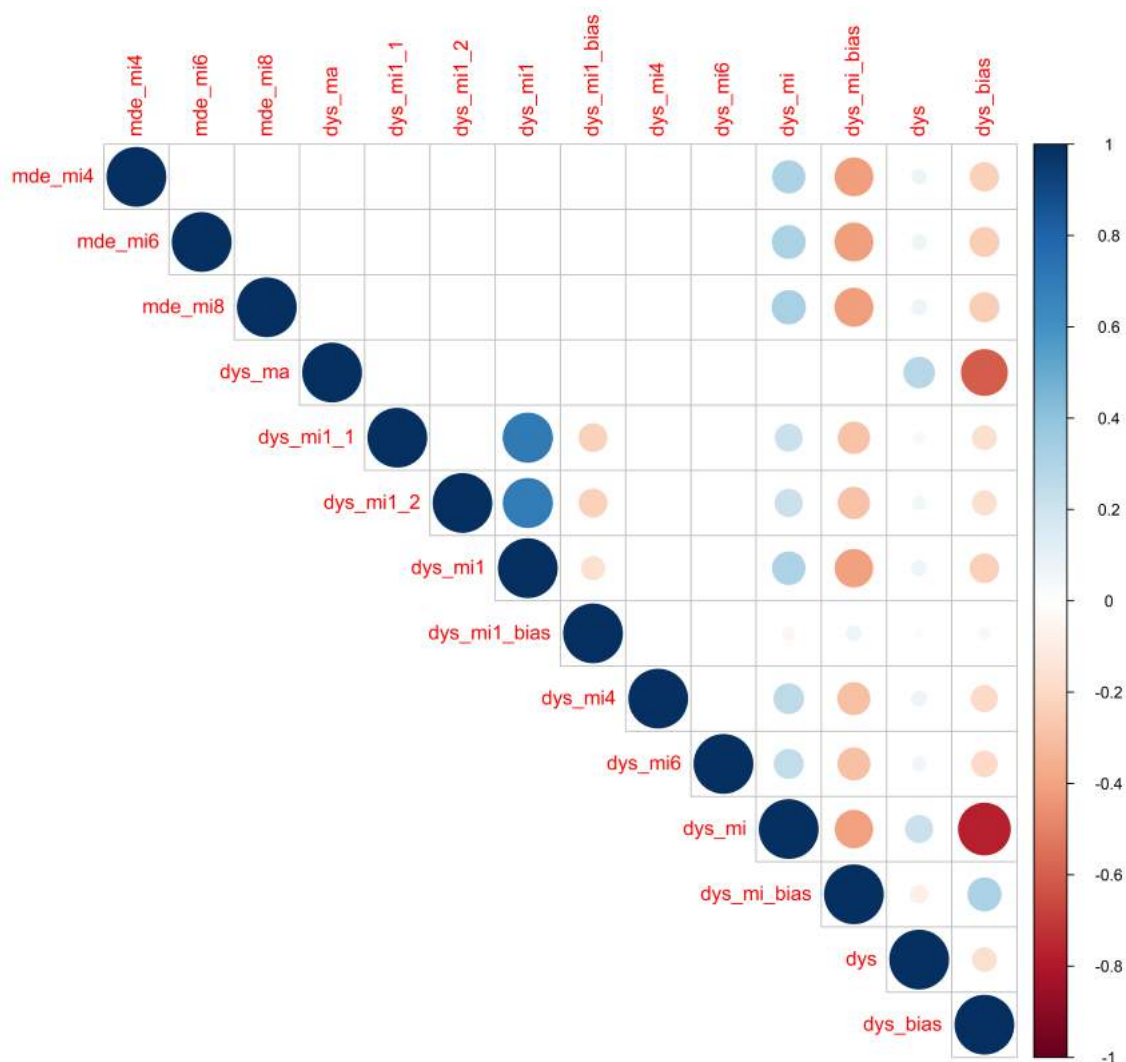


only

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

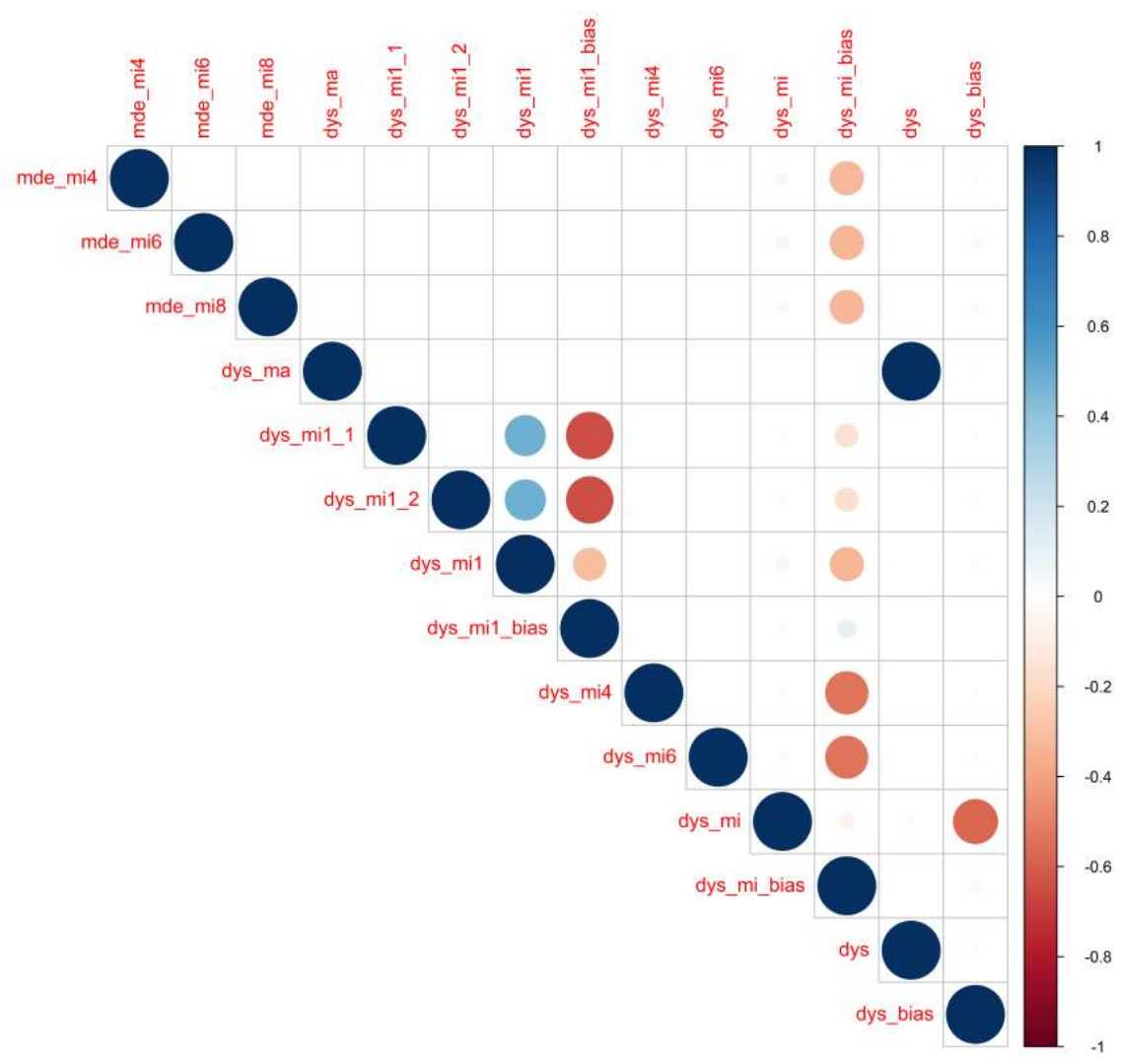


only

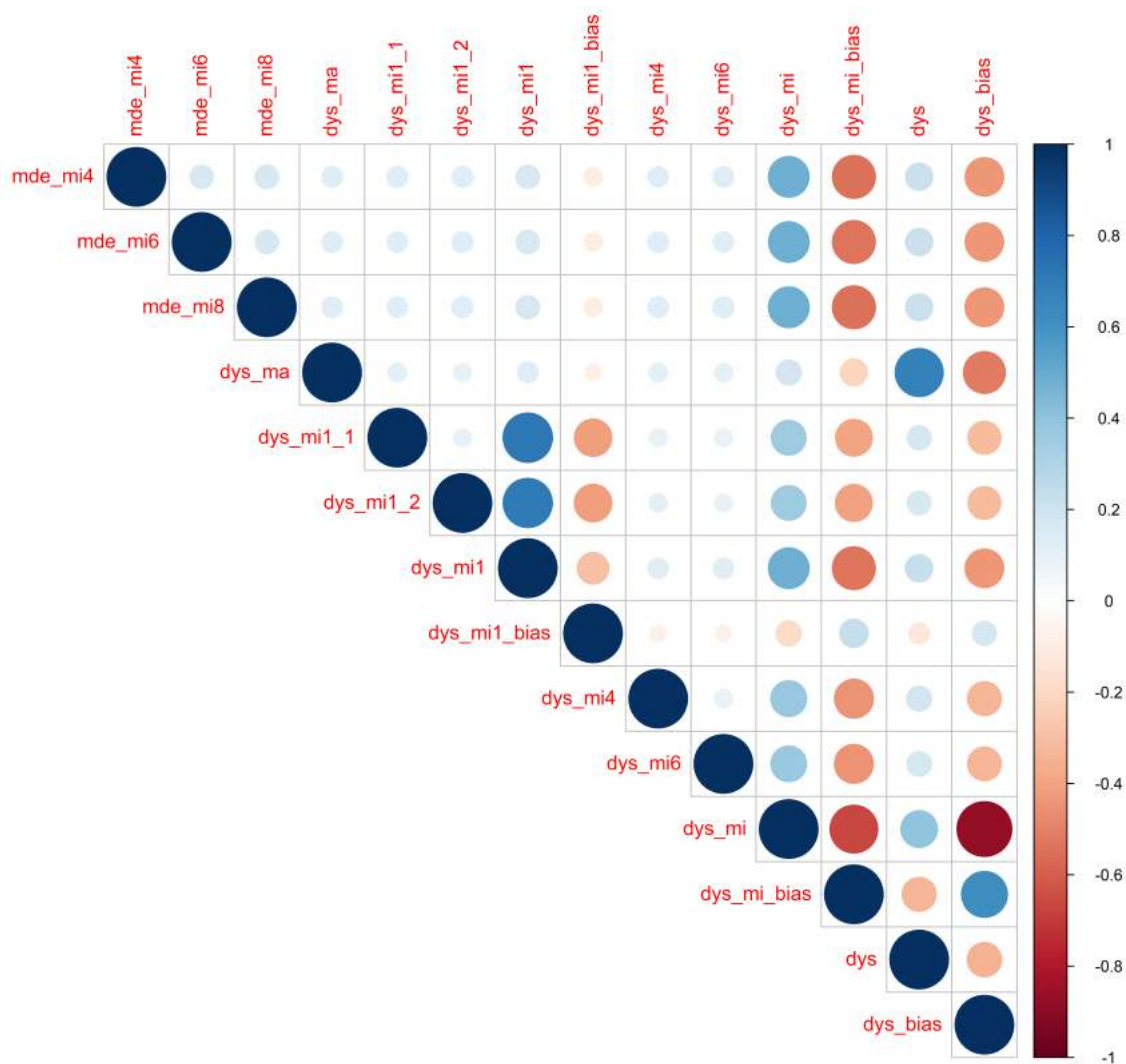


only

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

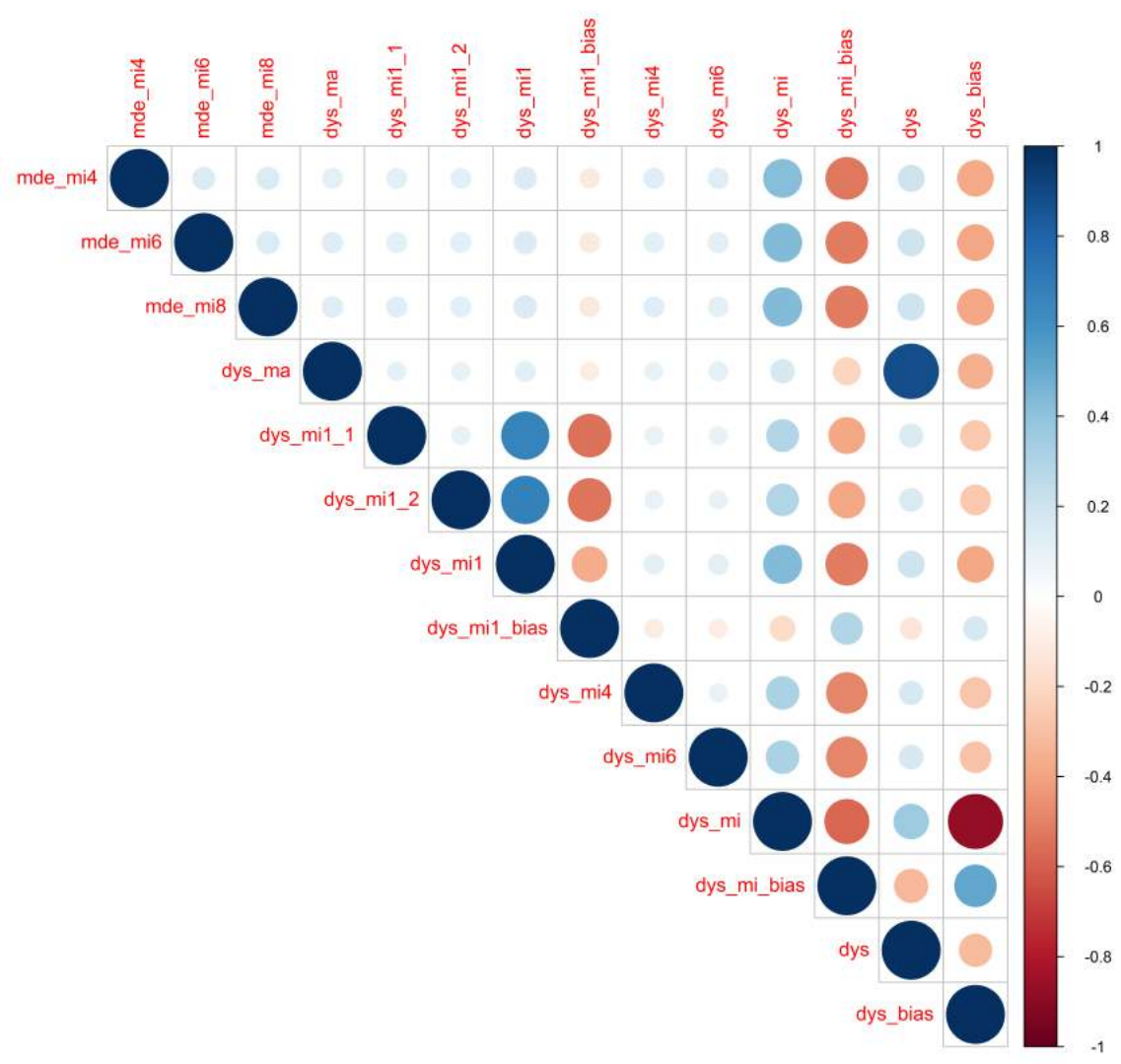


only

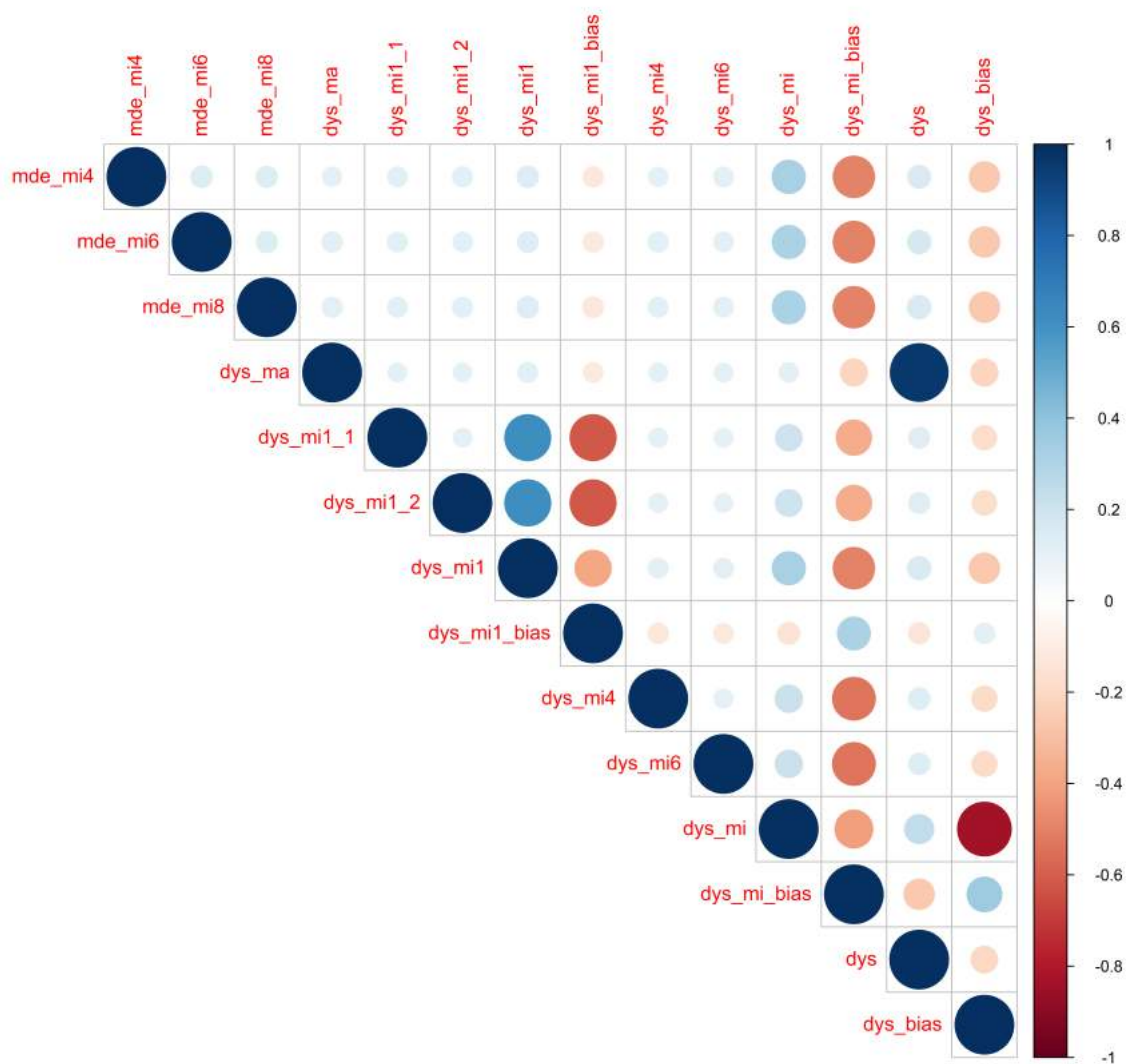


only

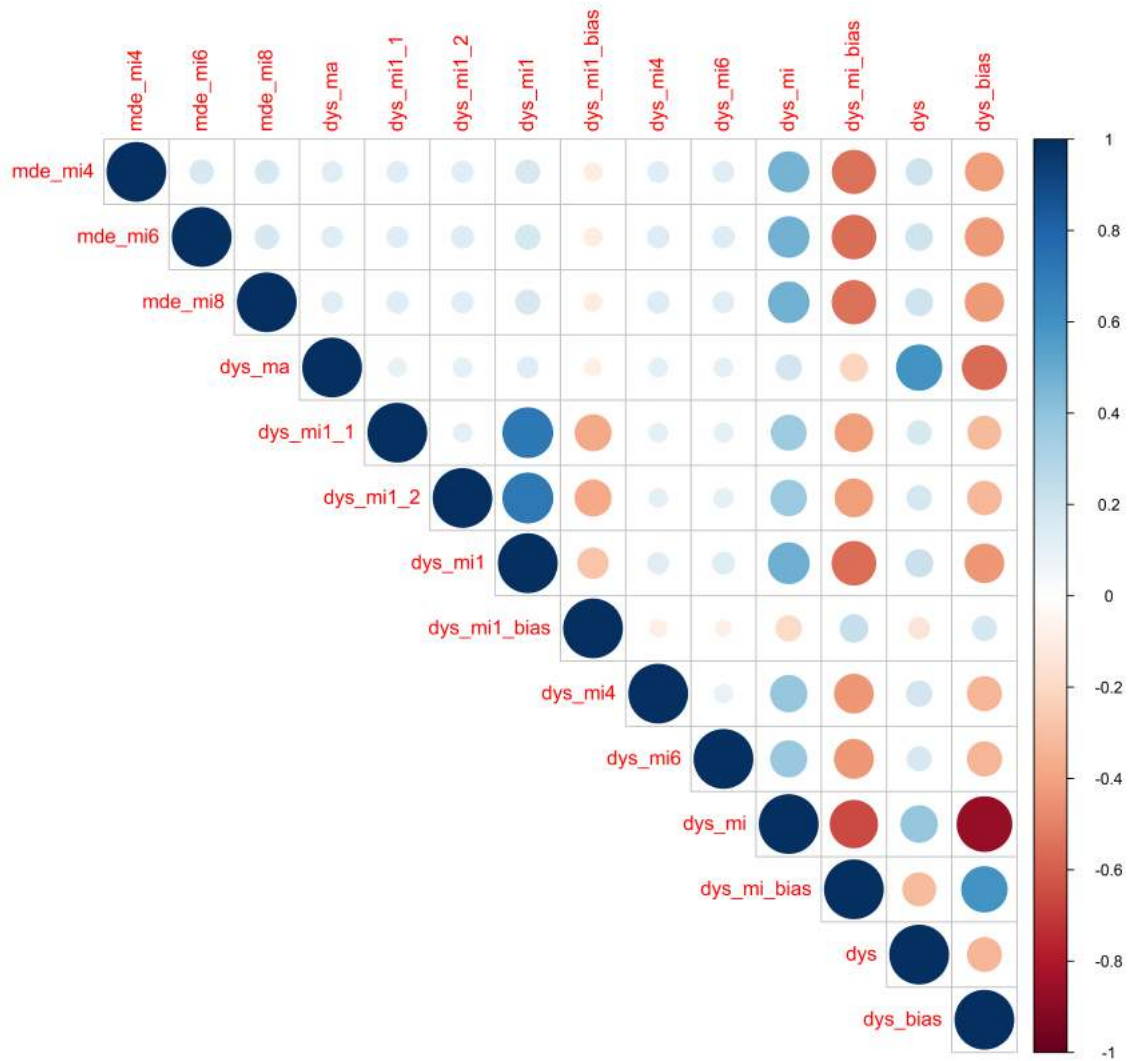
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



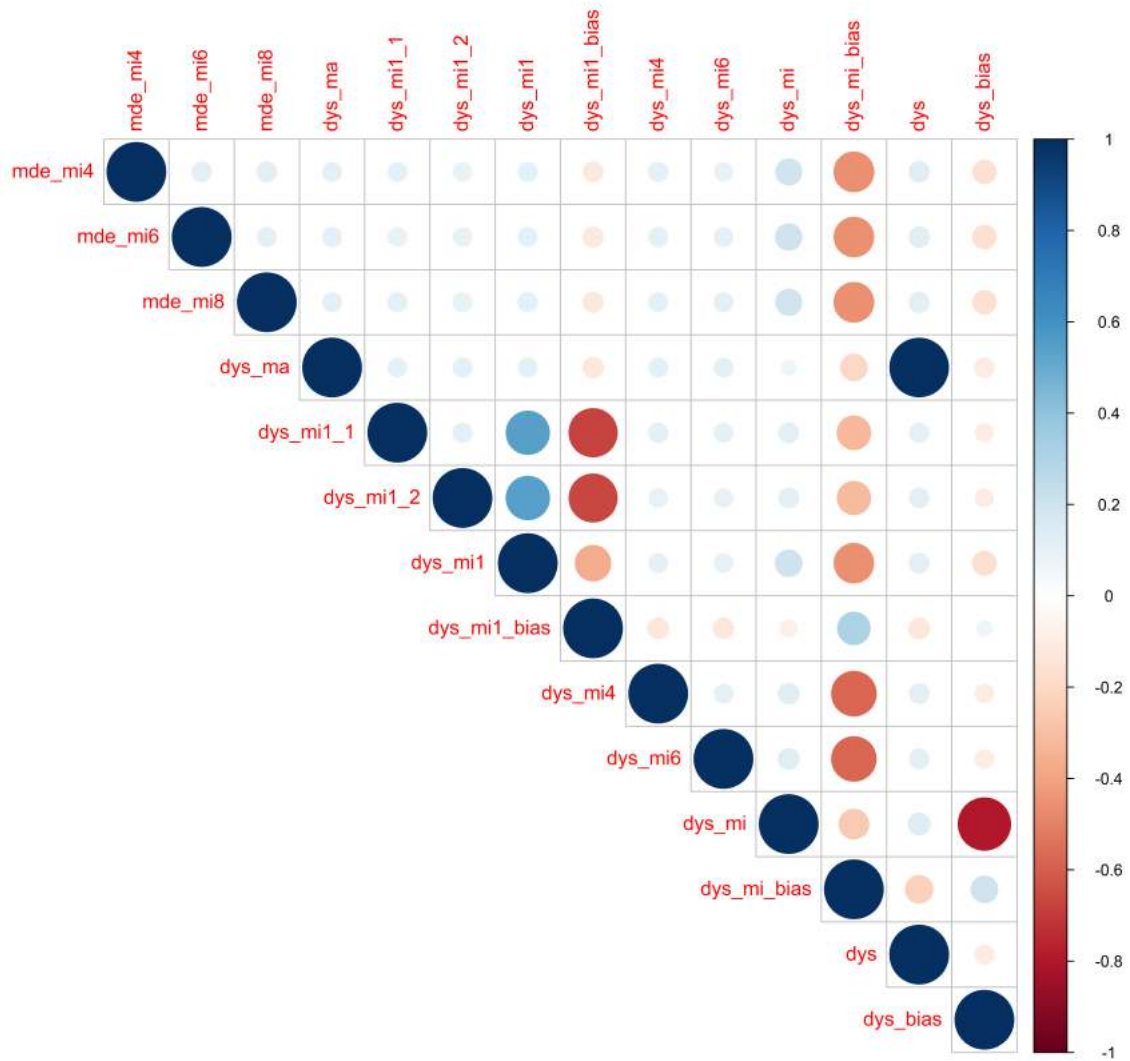
only



only

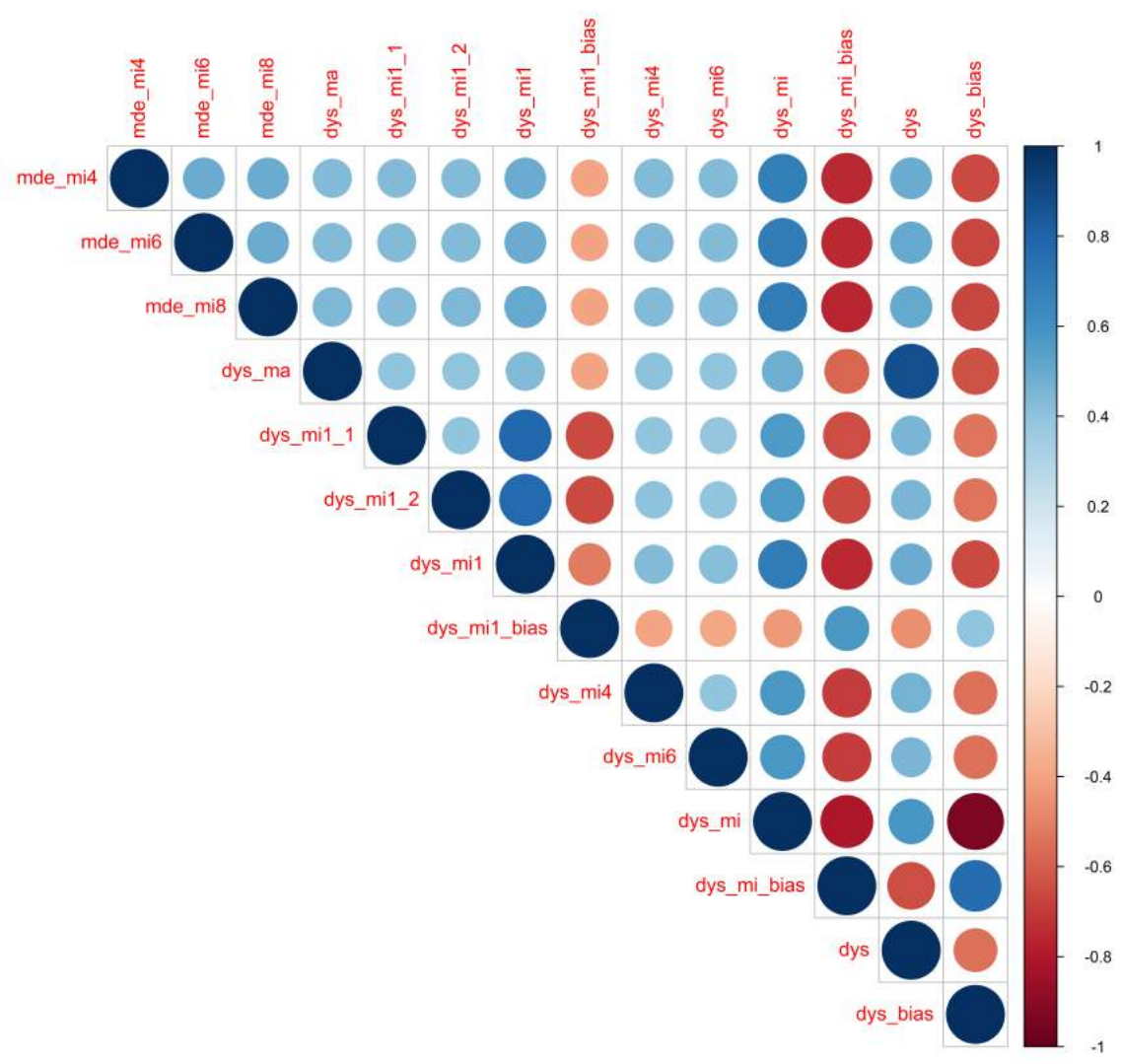


only

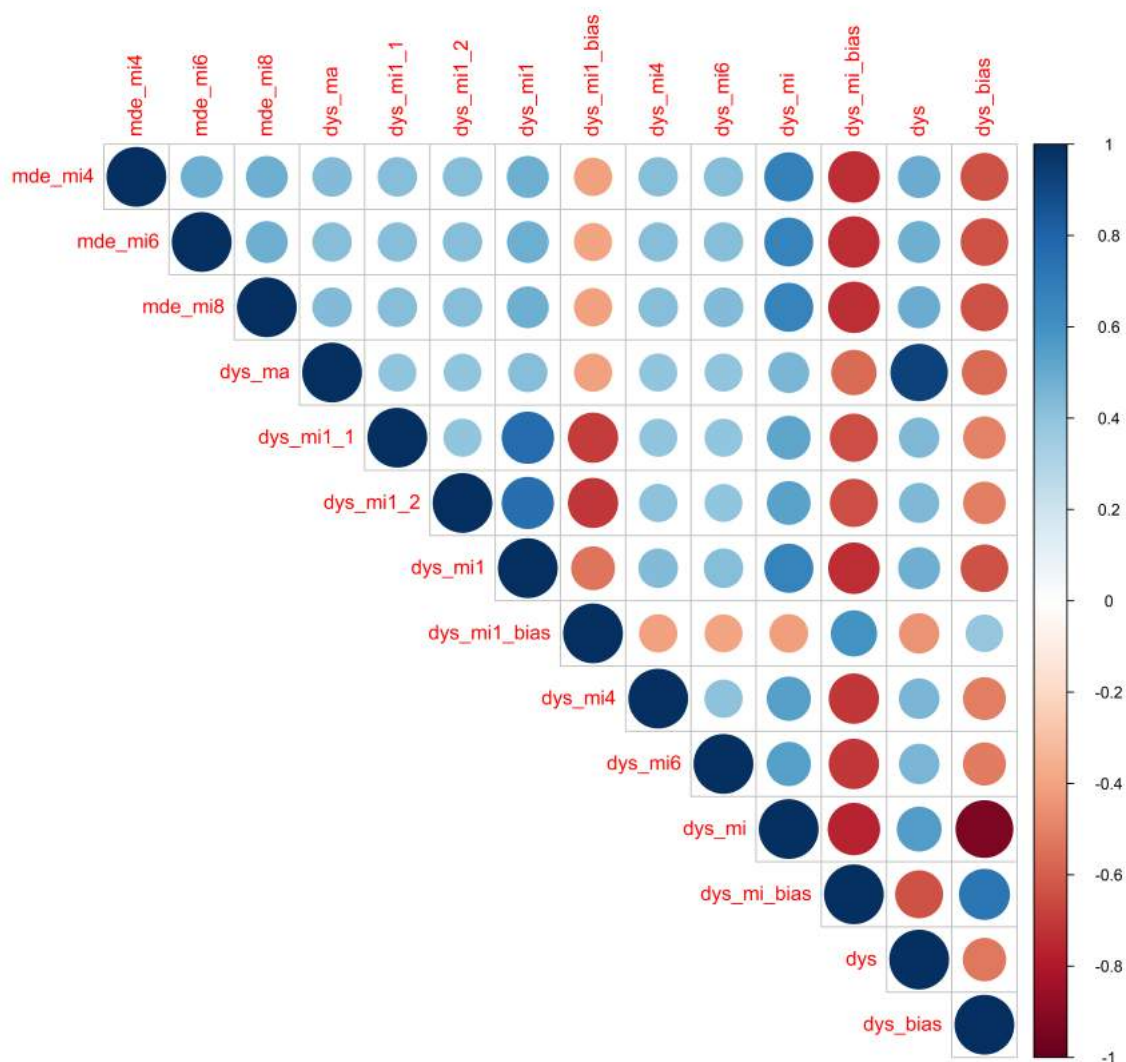


only

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

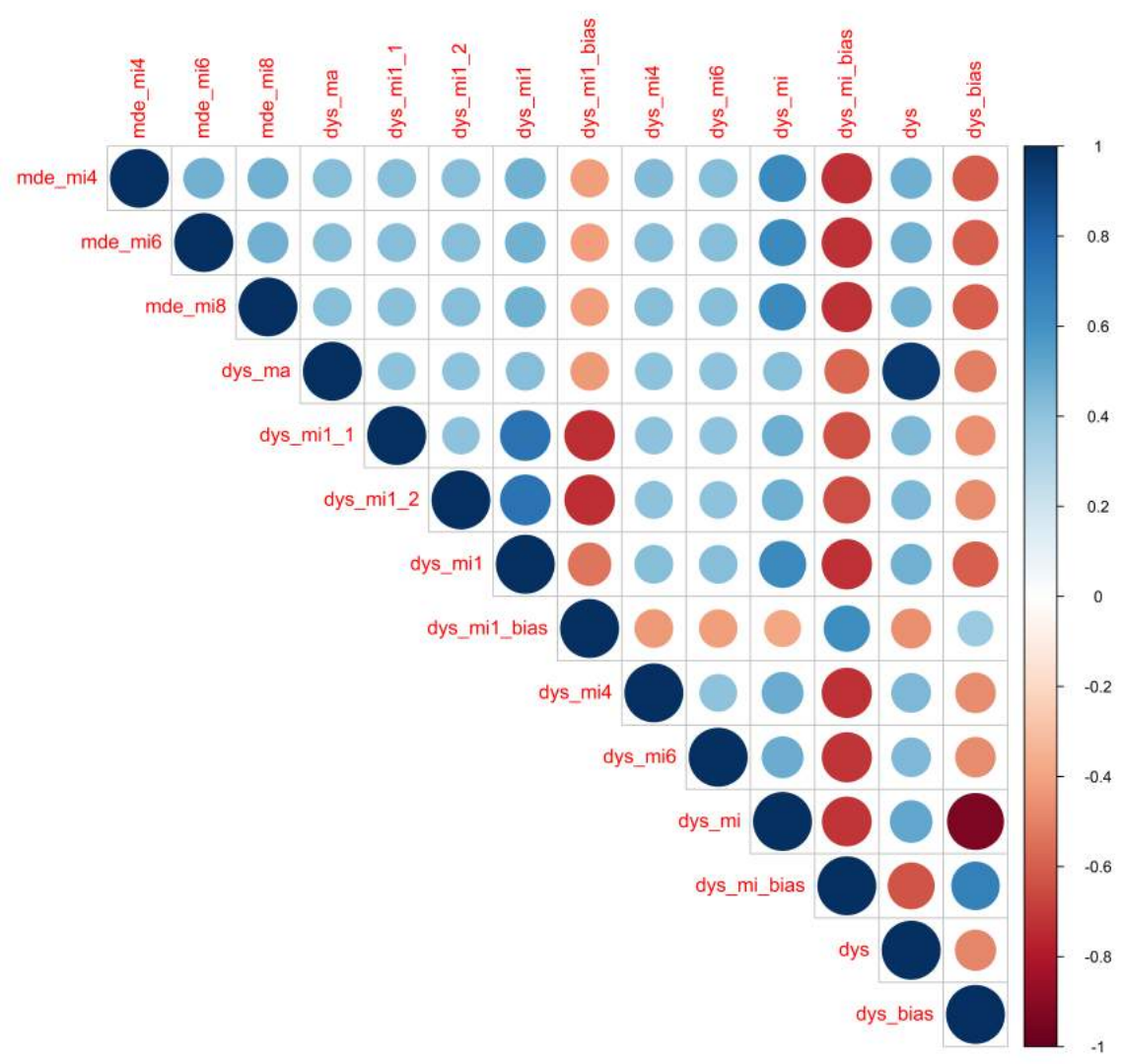


only

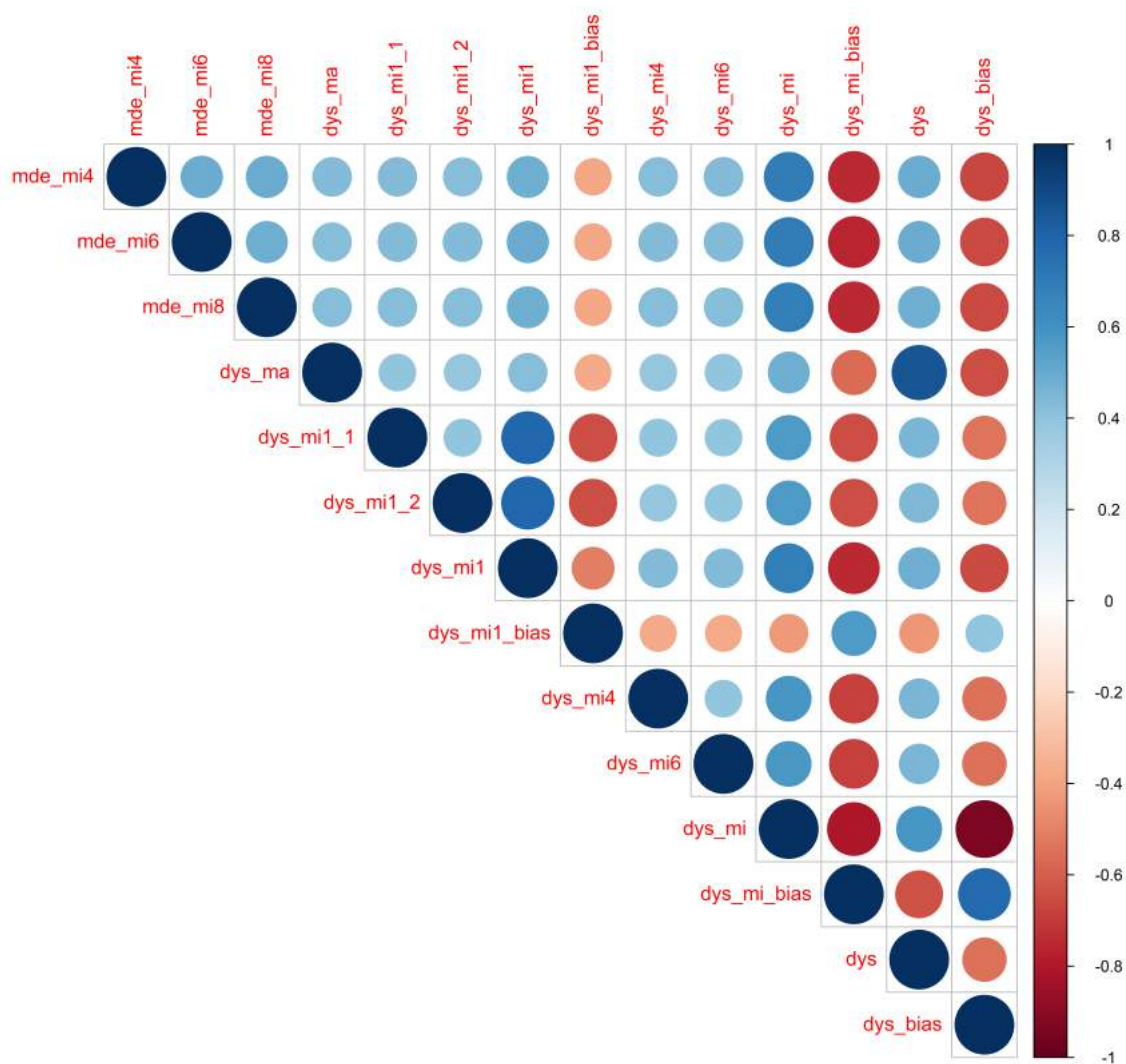


only

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

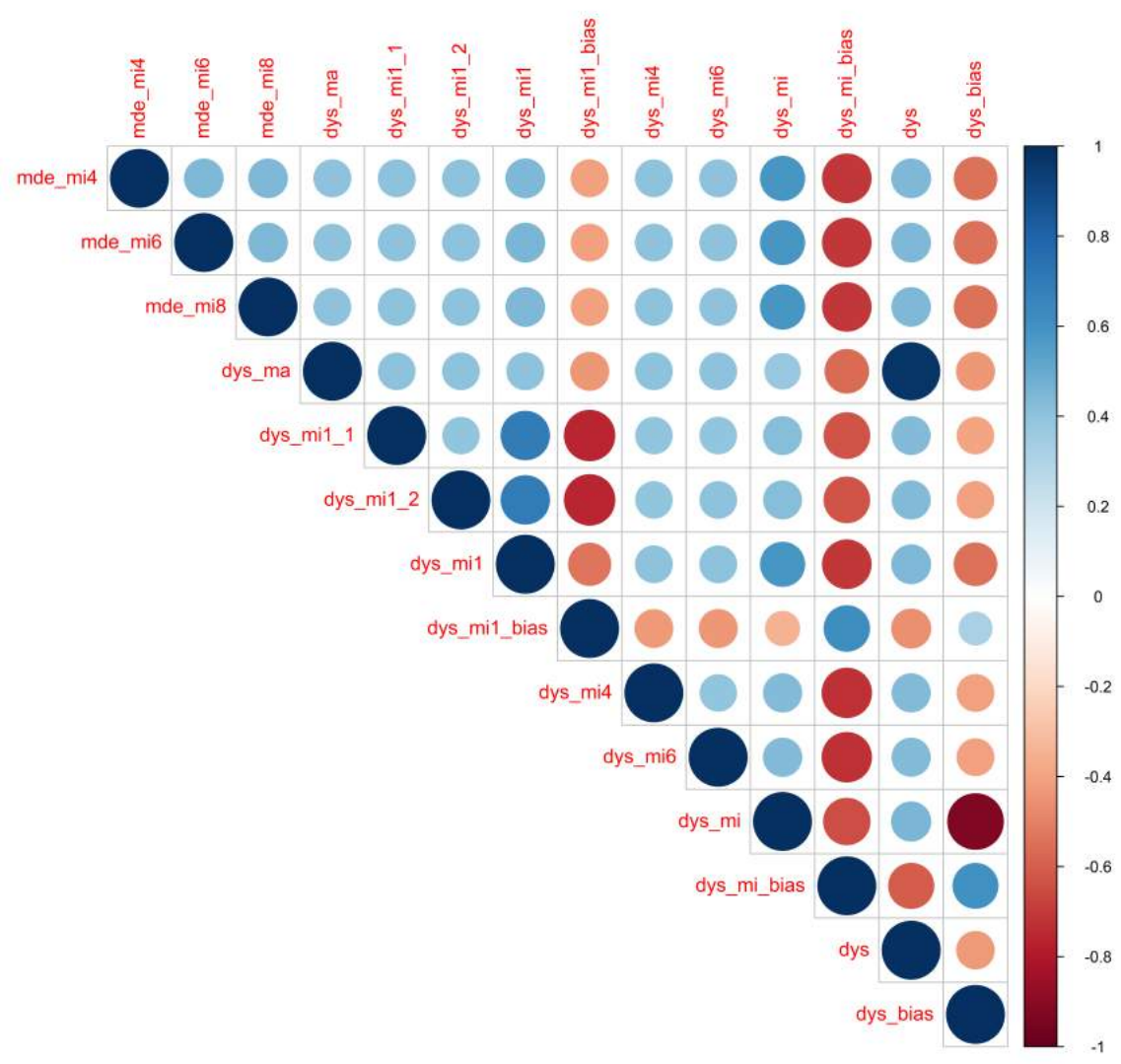


only

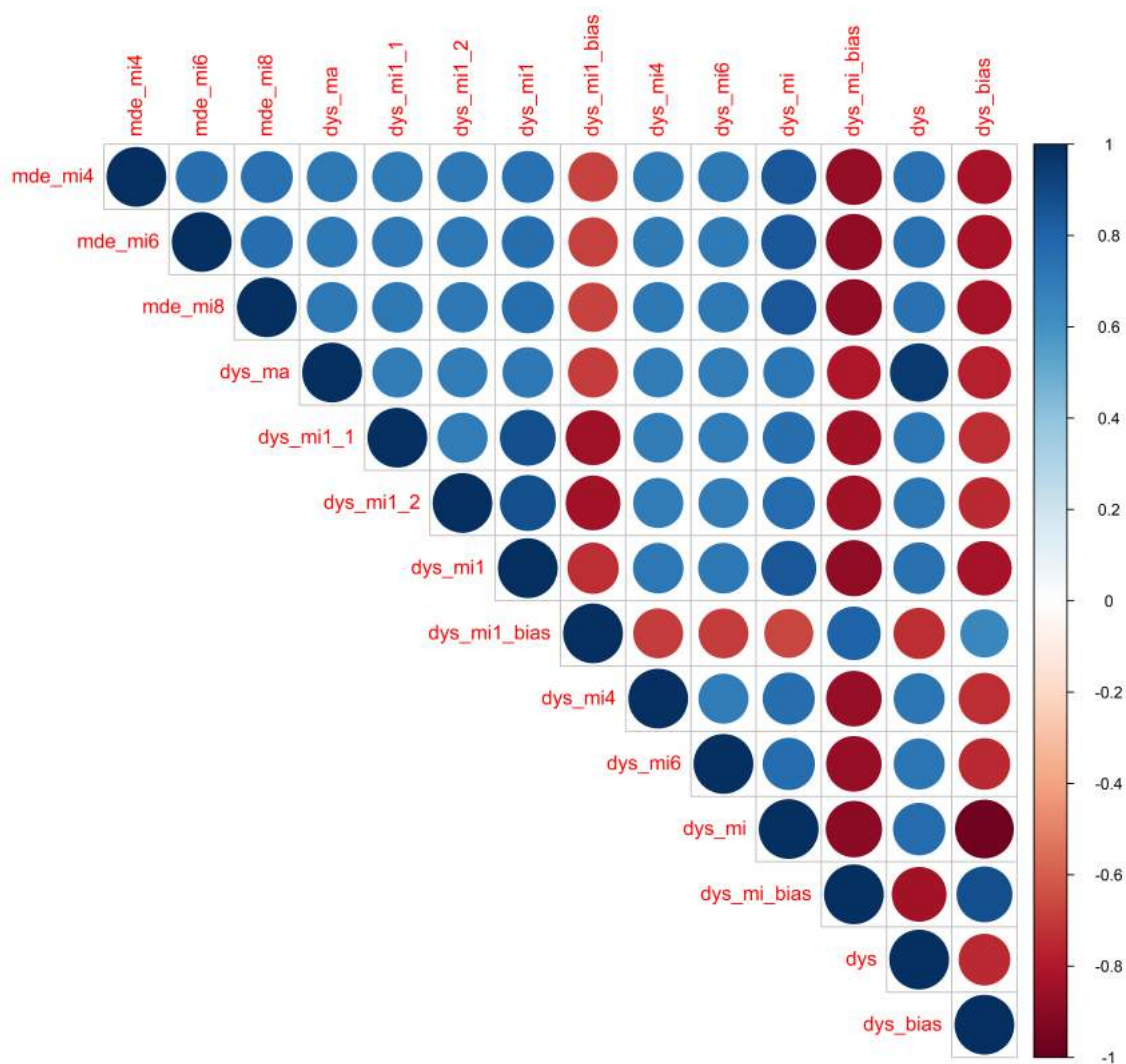


only

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

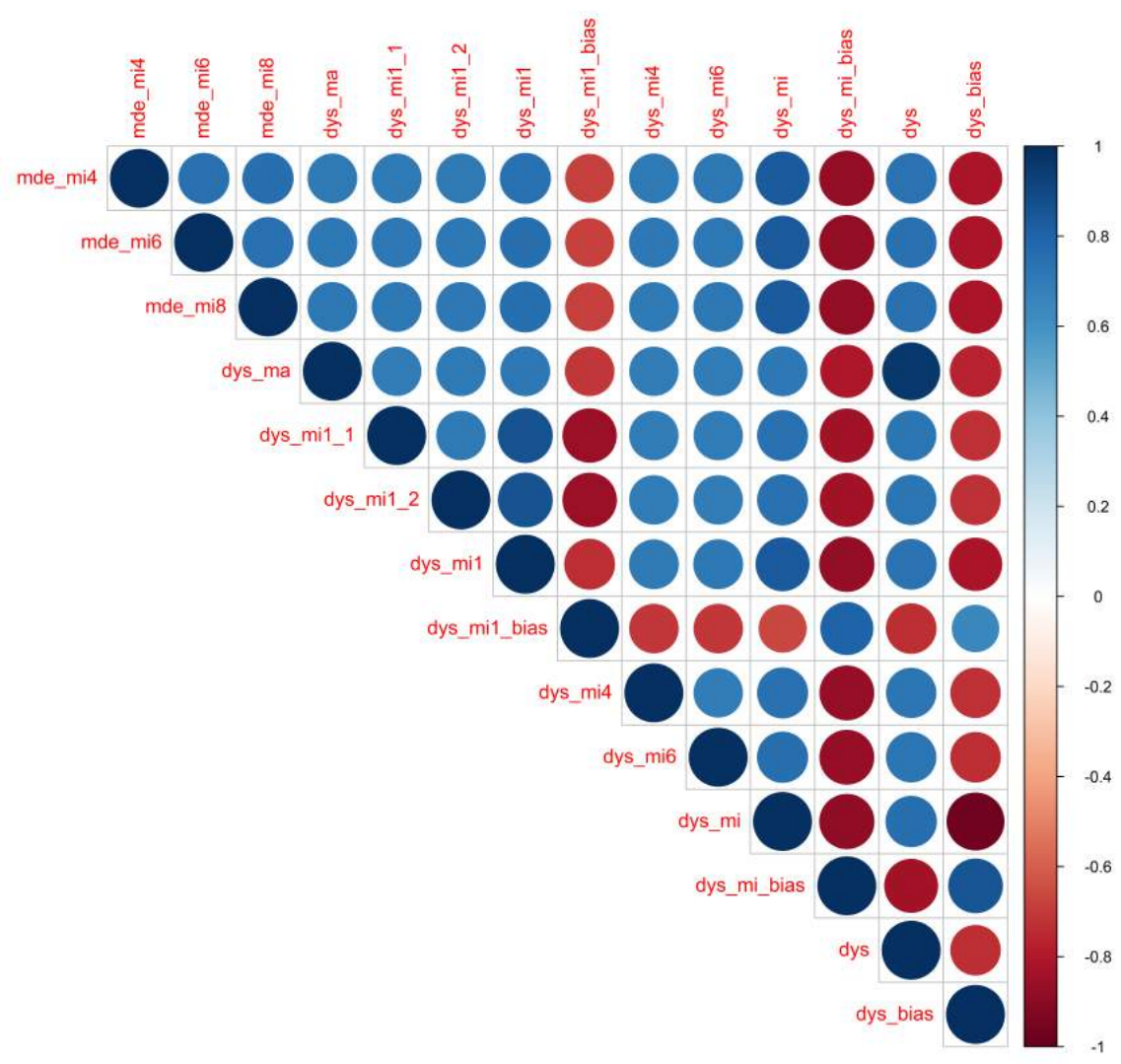


only

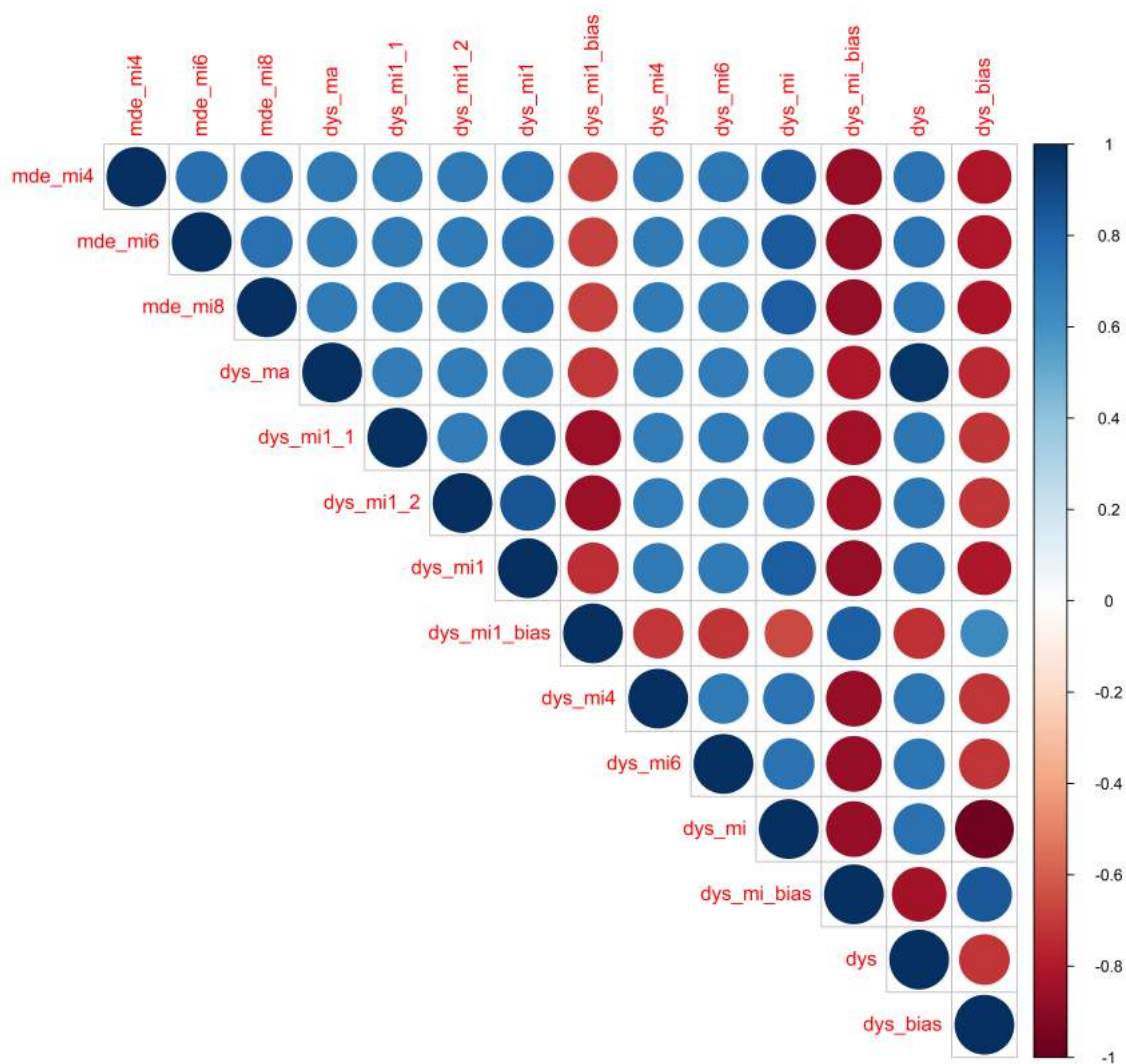


only

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

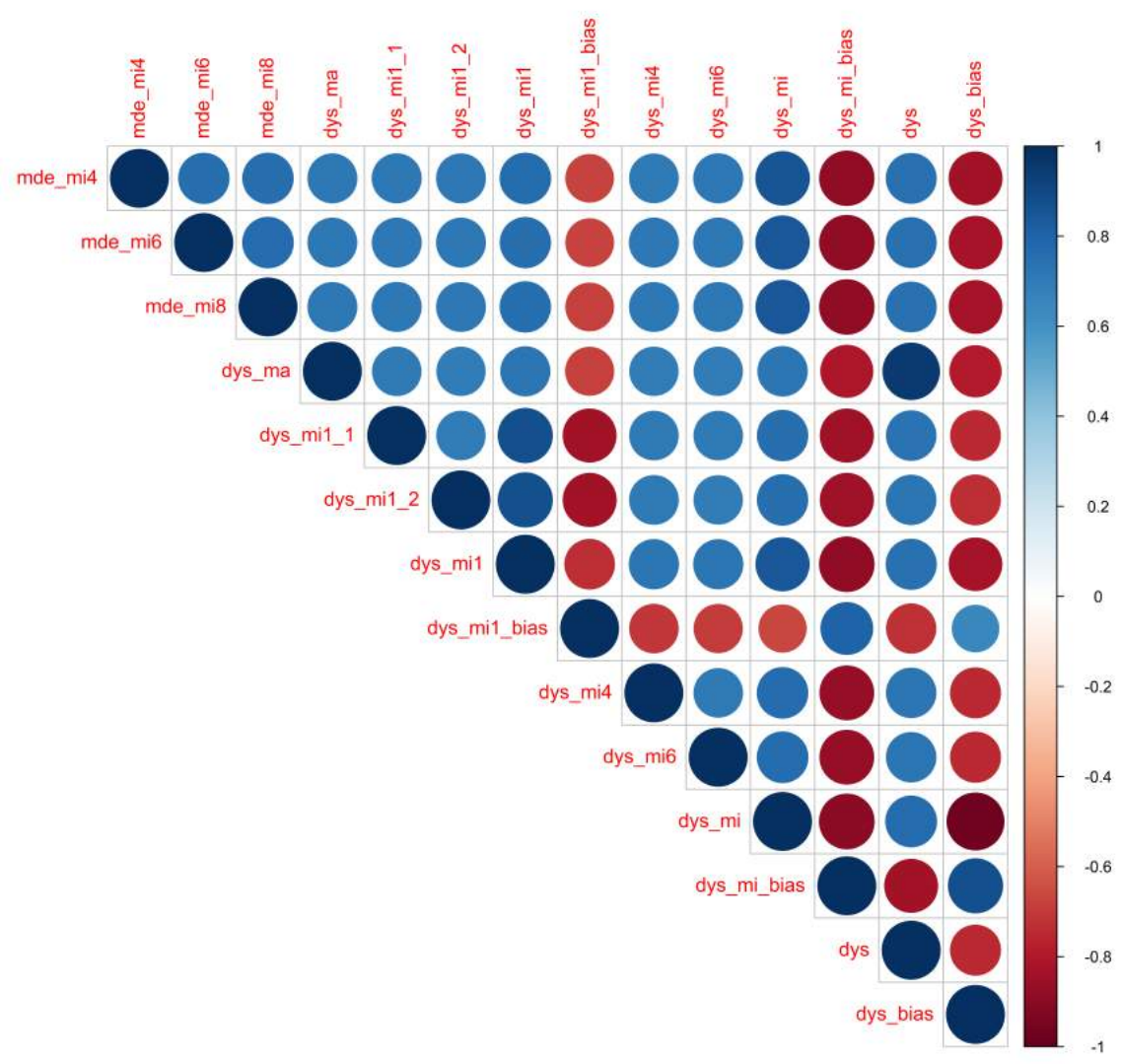


only

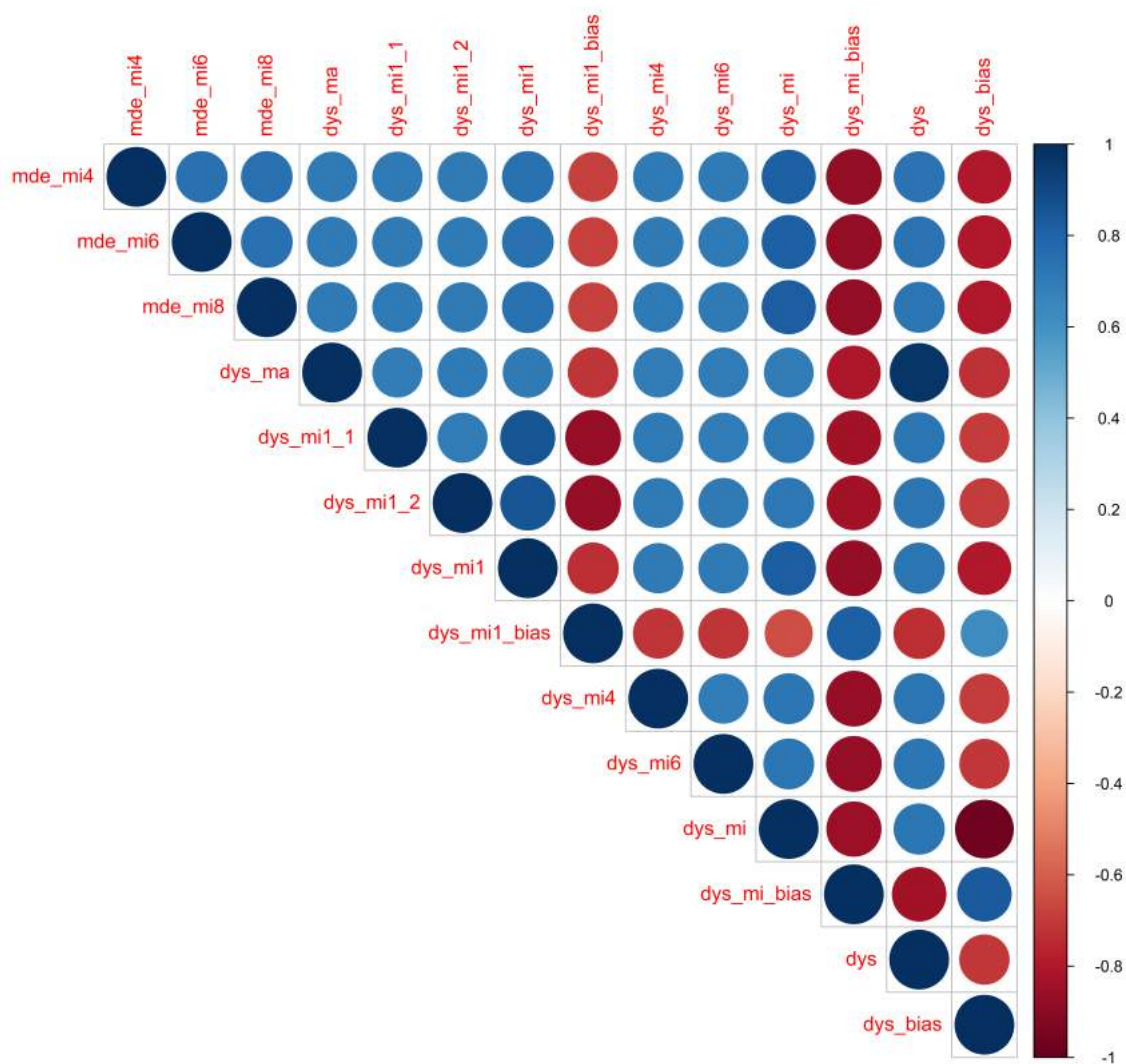


only

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

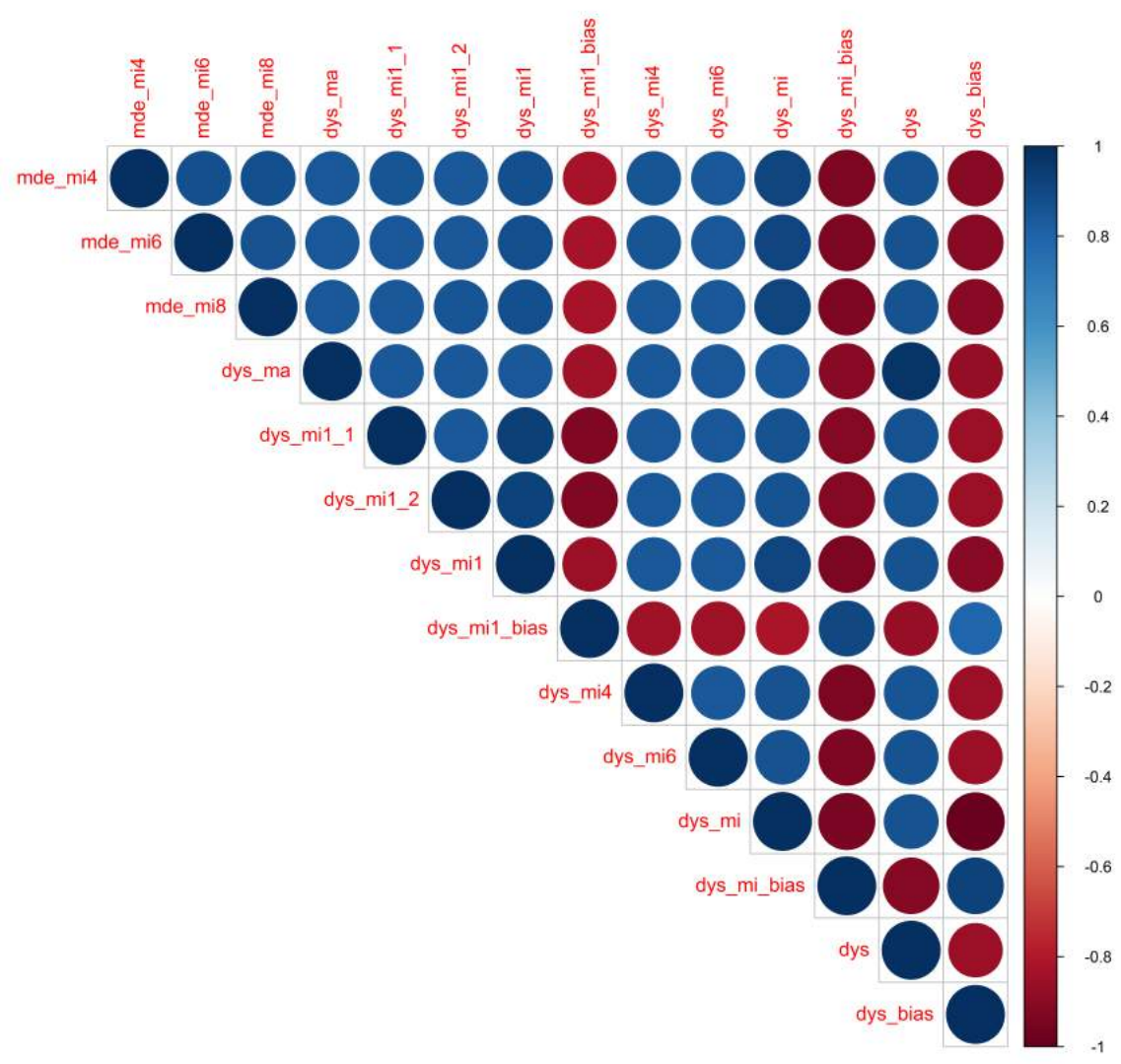


only

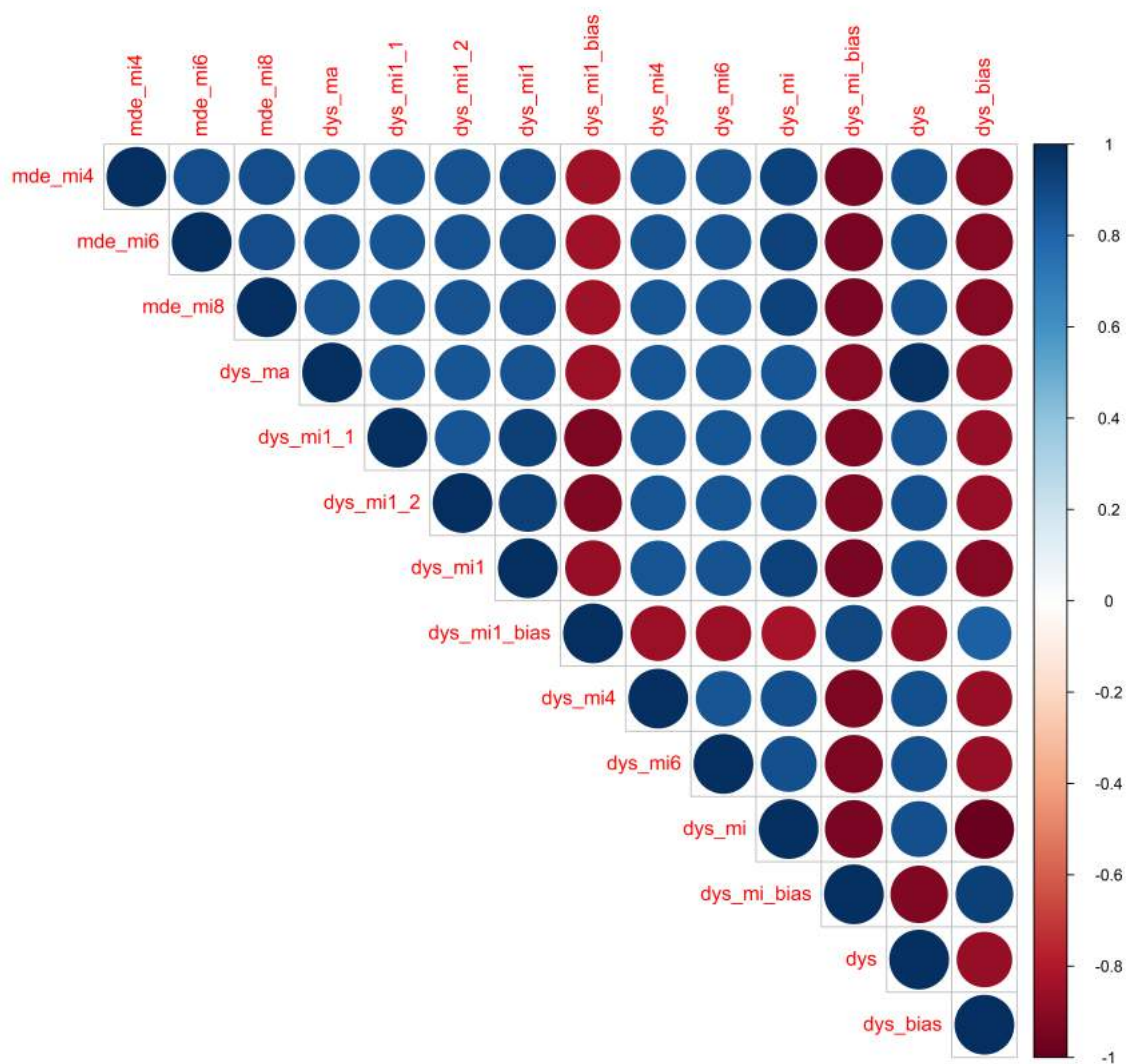


only

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

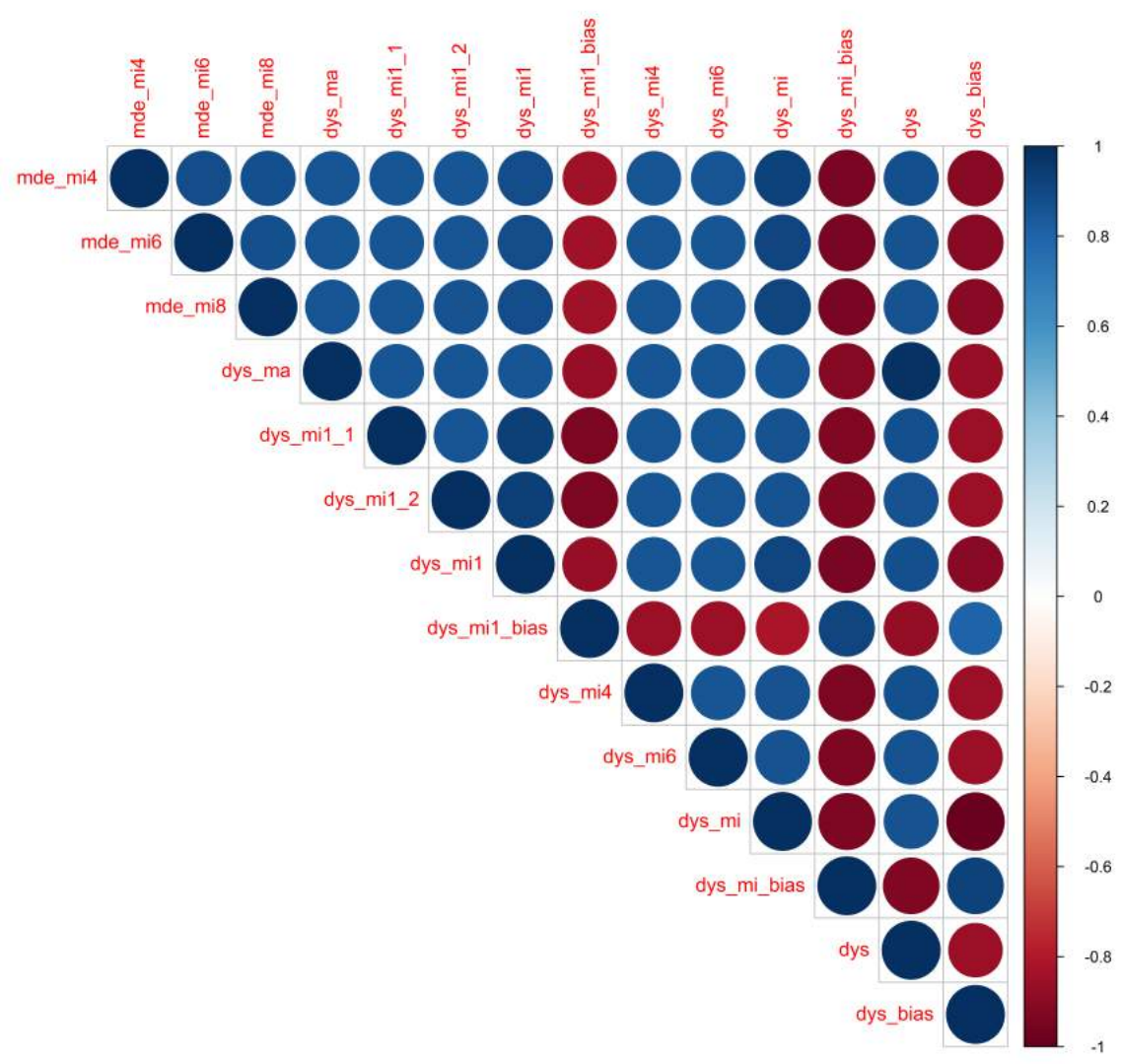


only

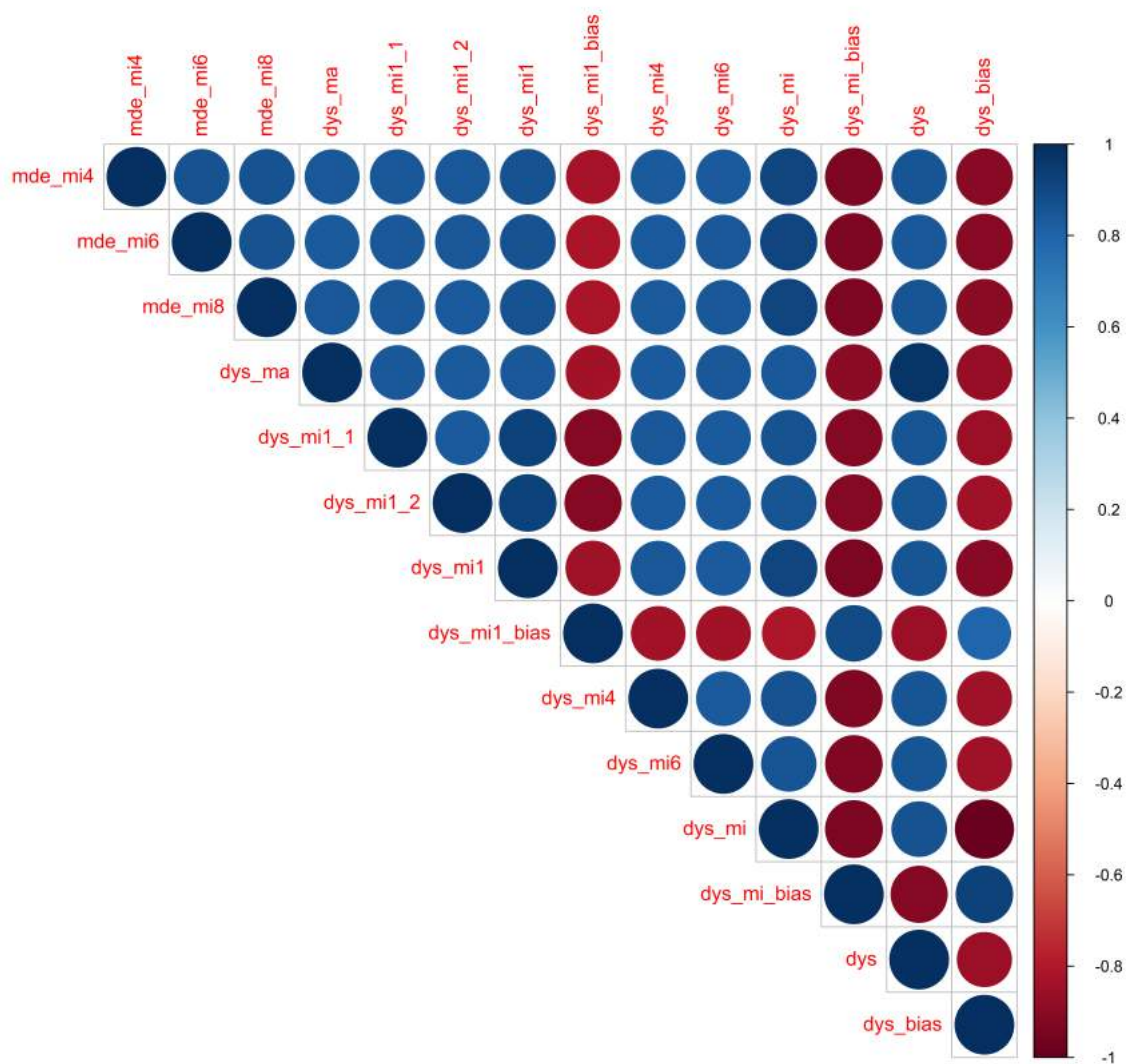


only

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

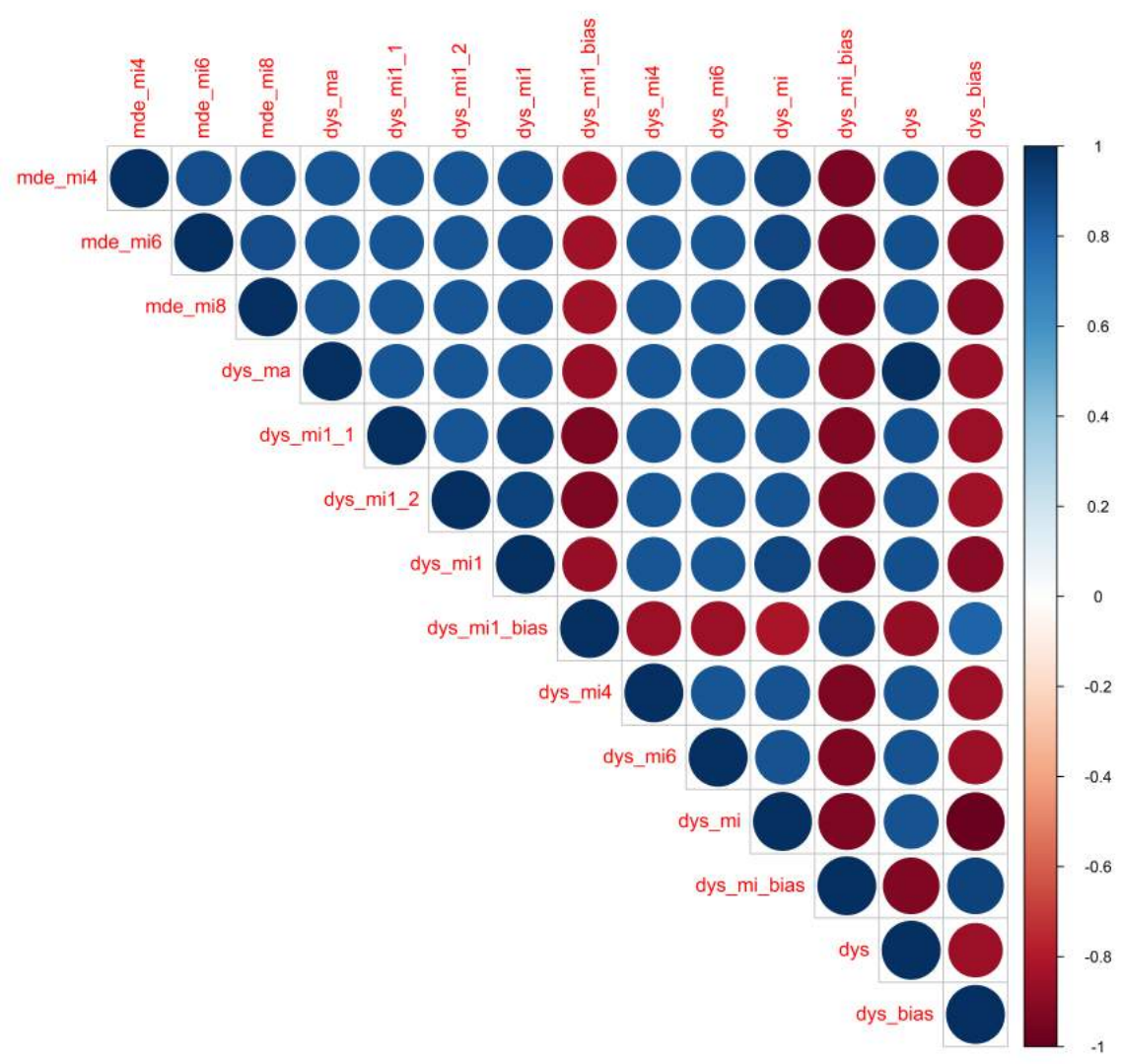


only



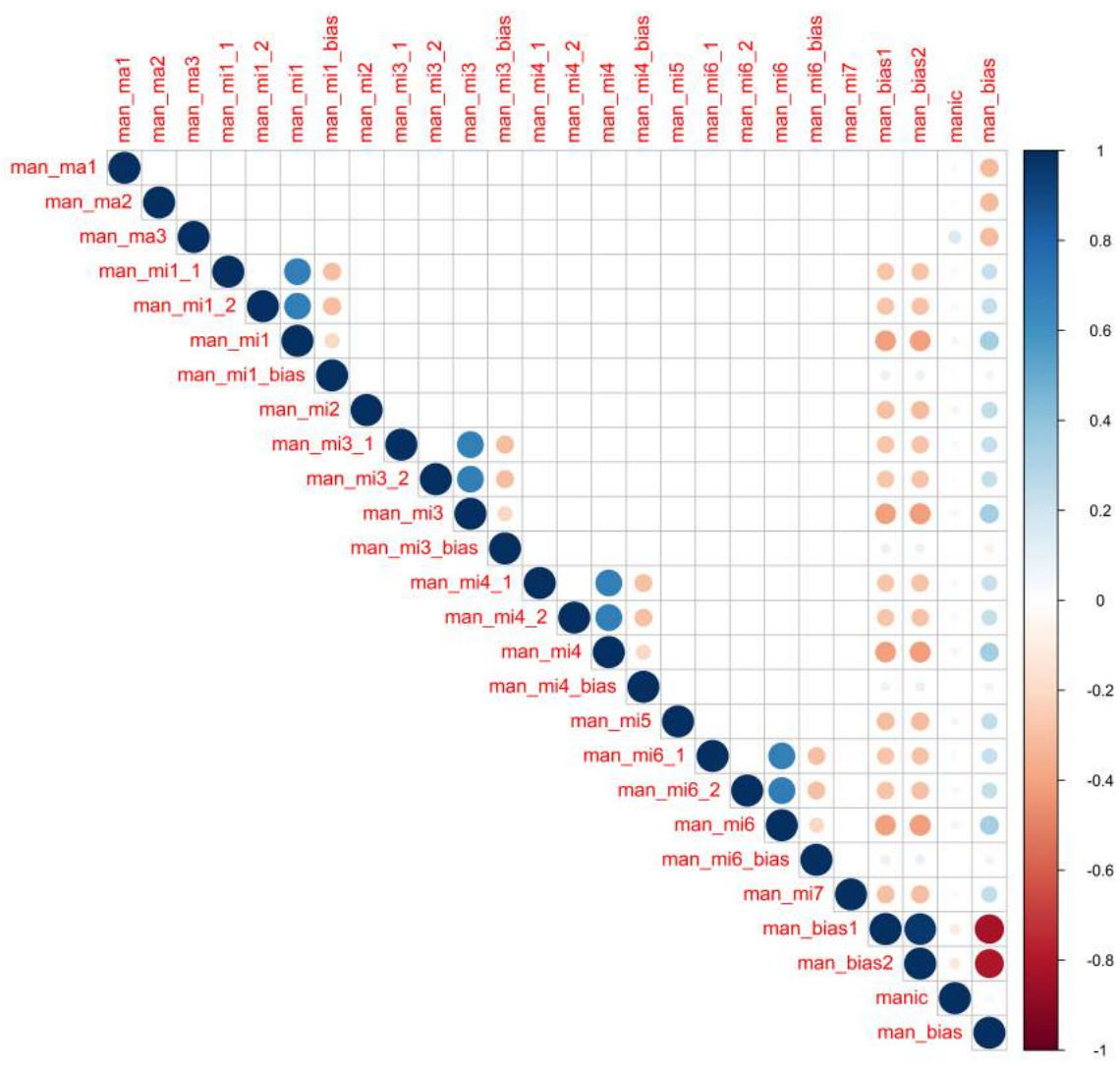
only

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



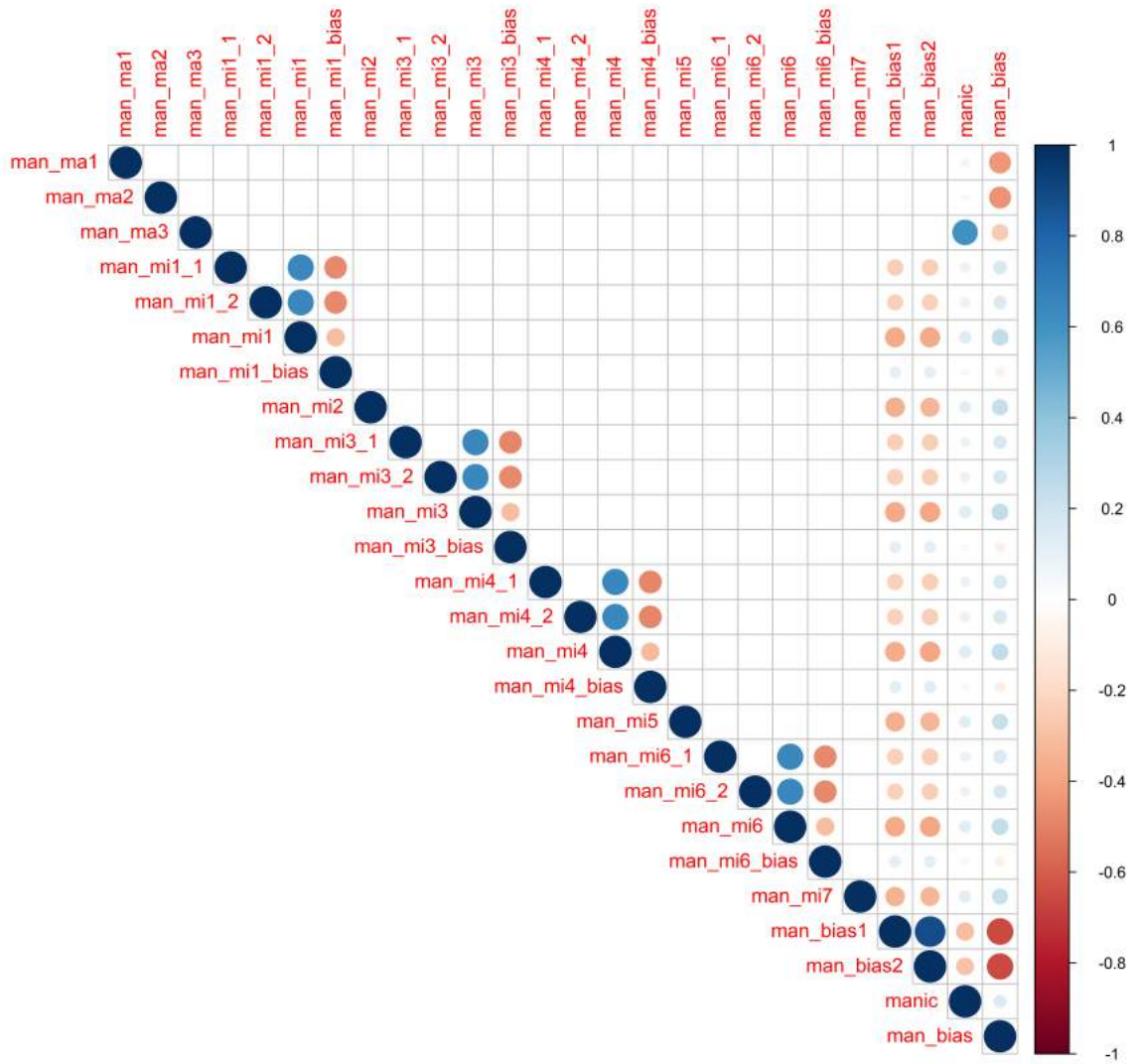
only

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

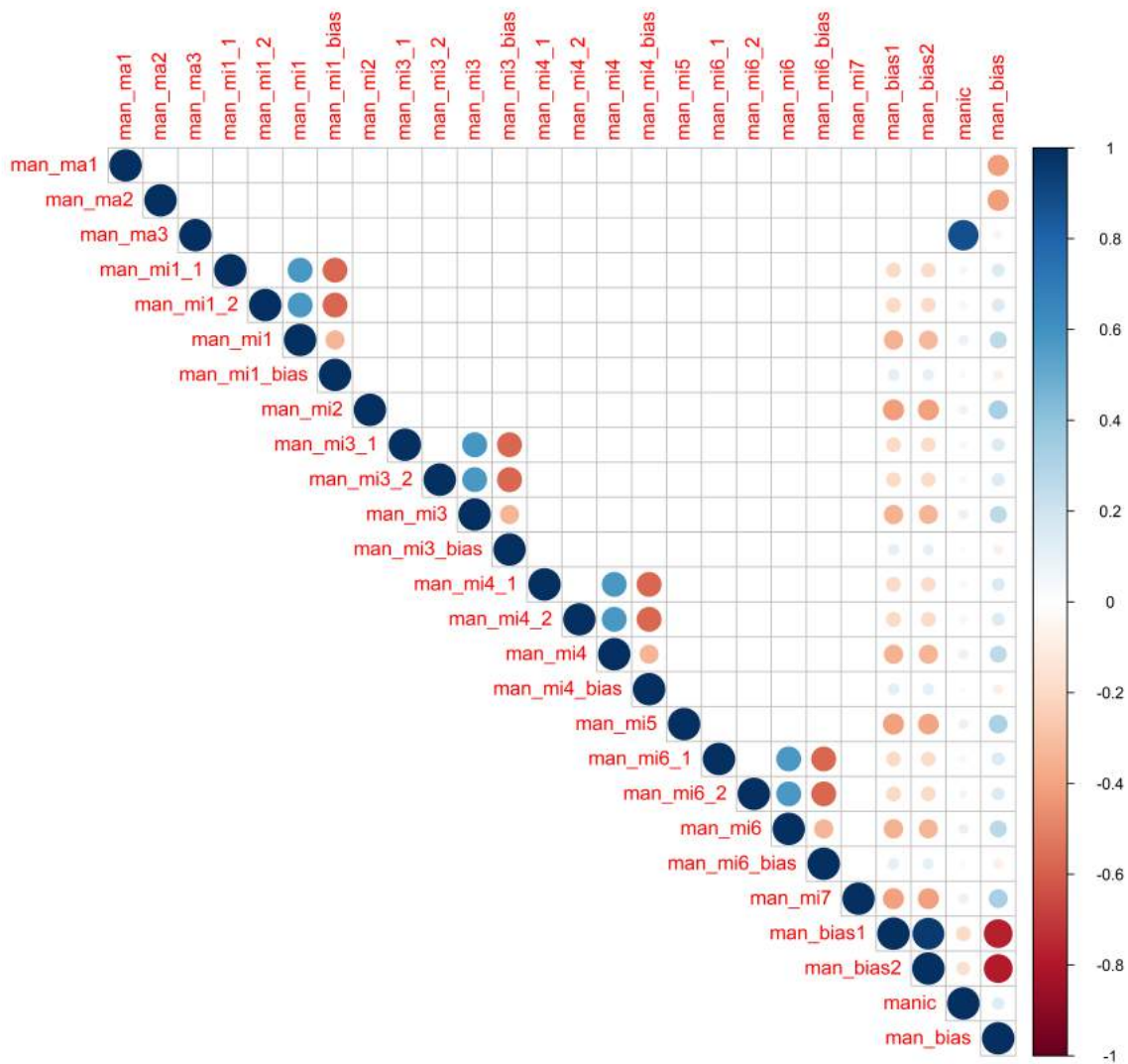


only

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

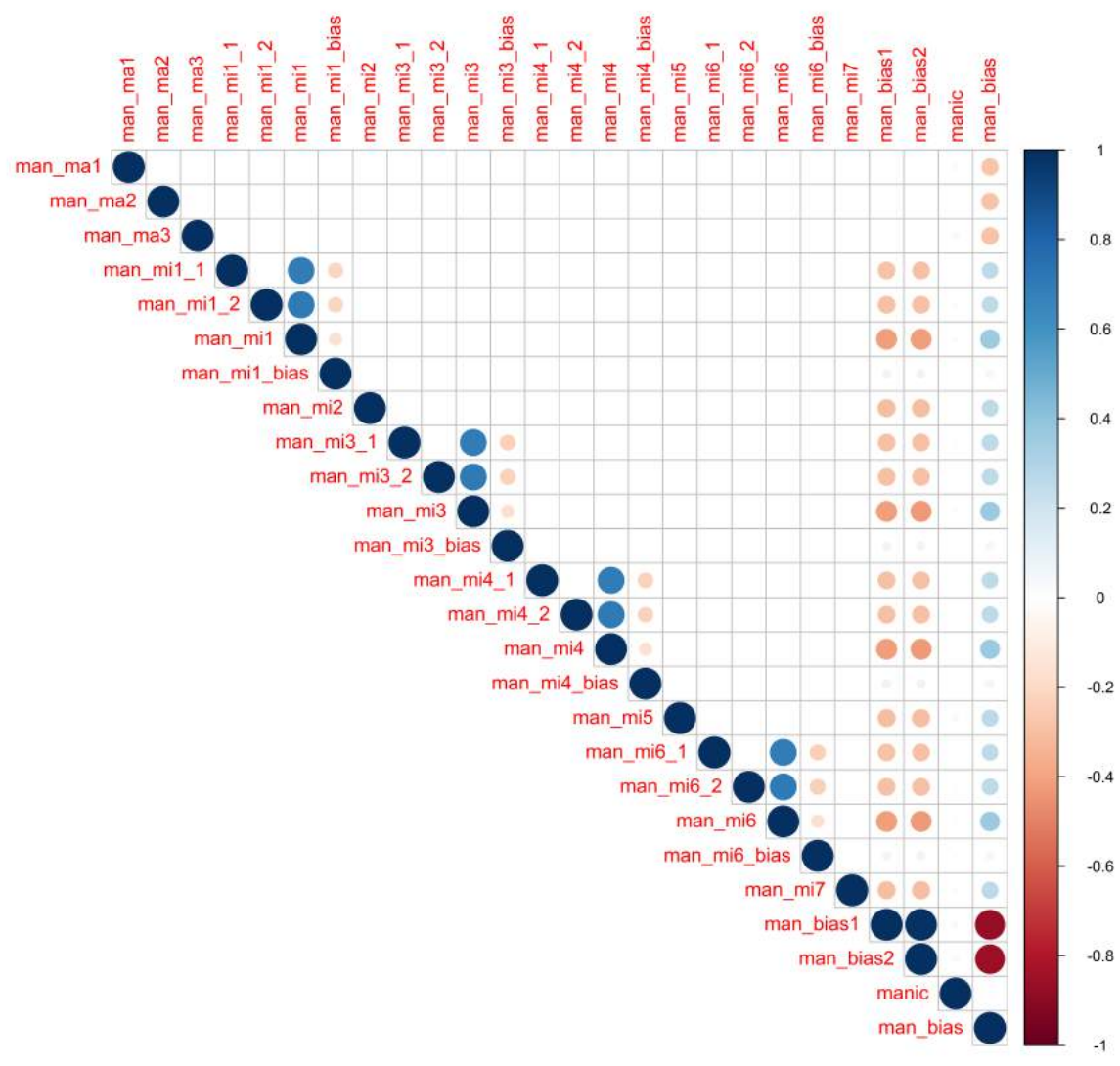


only



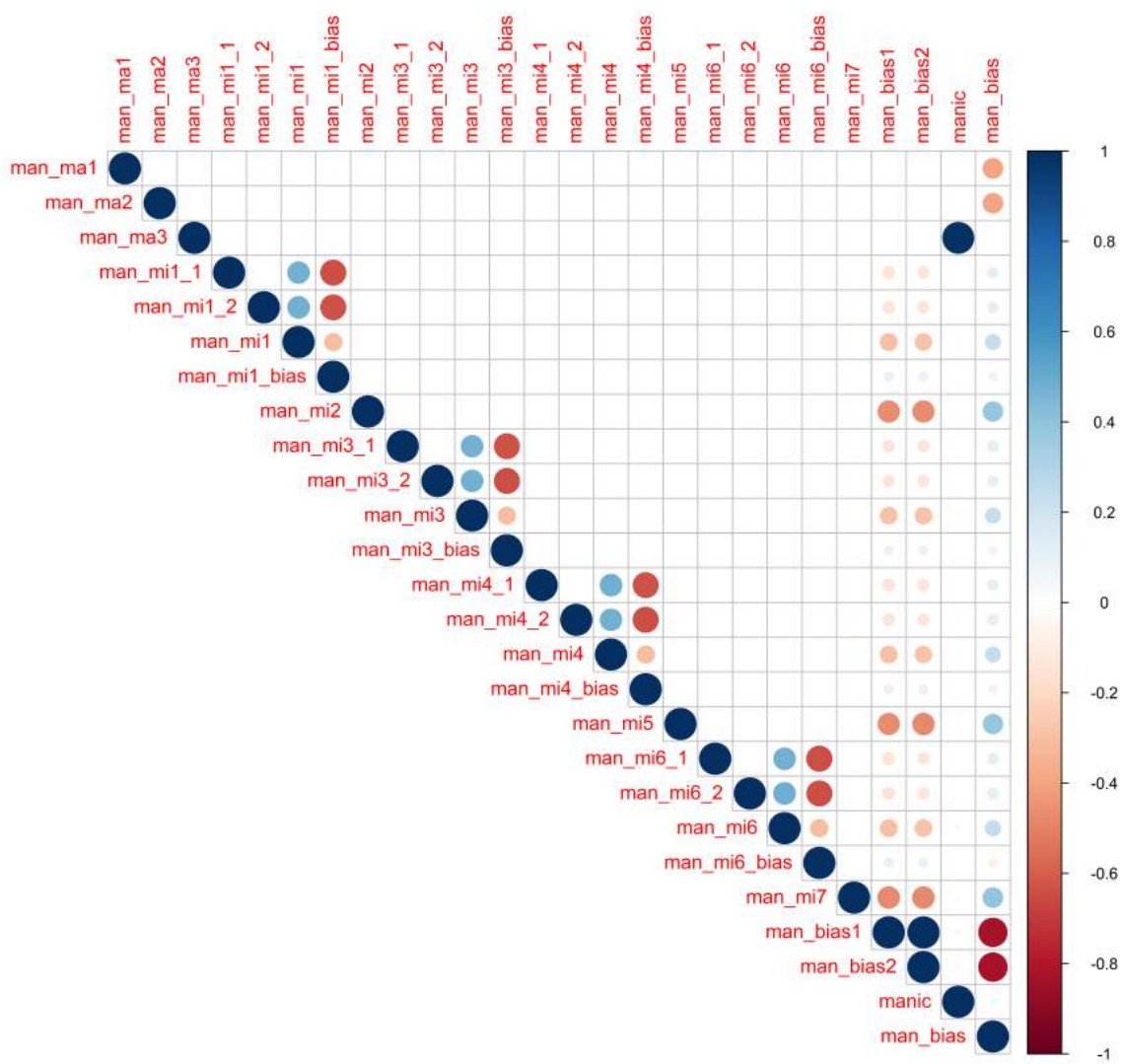
only

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

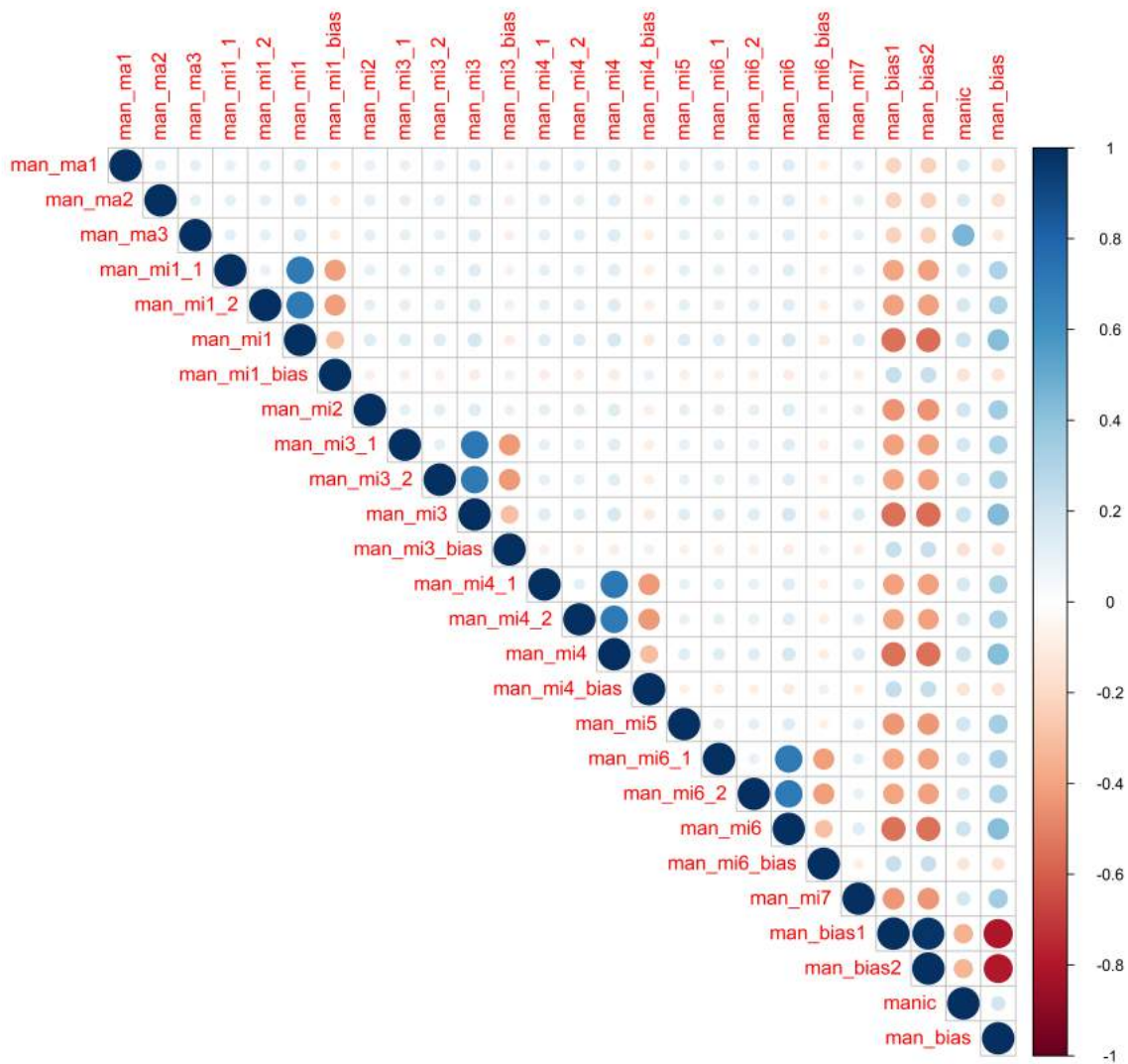


only

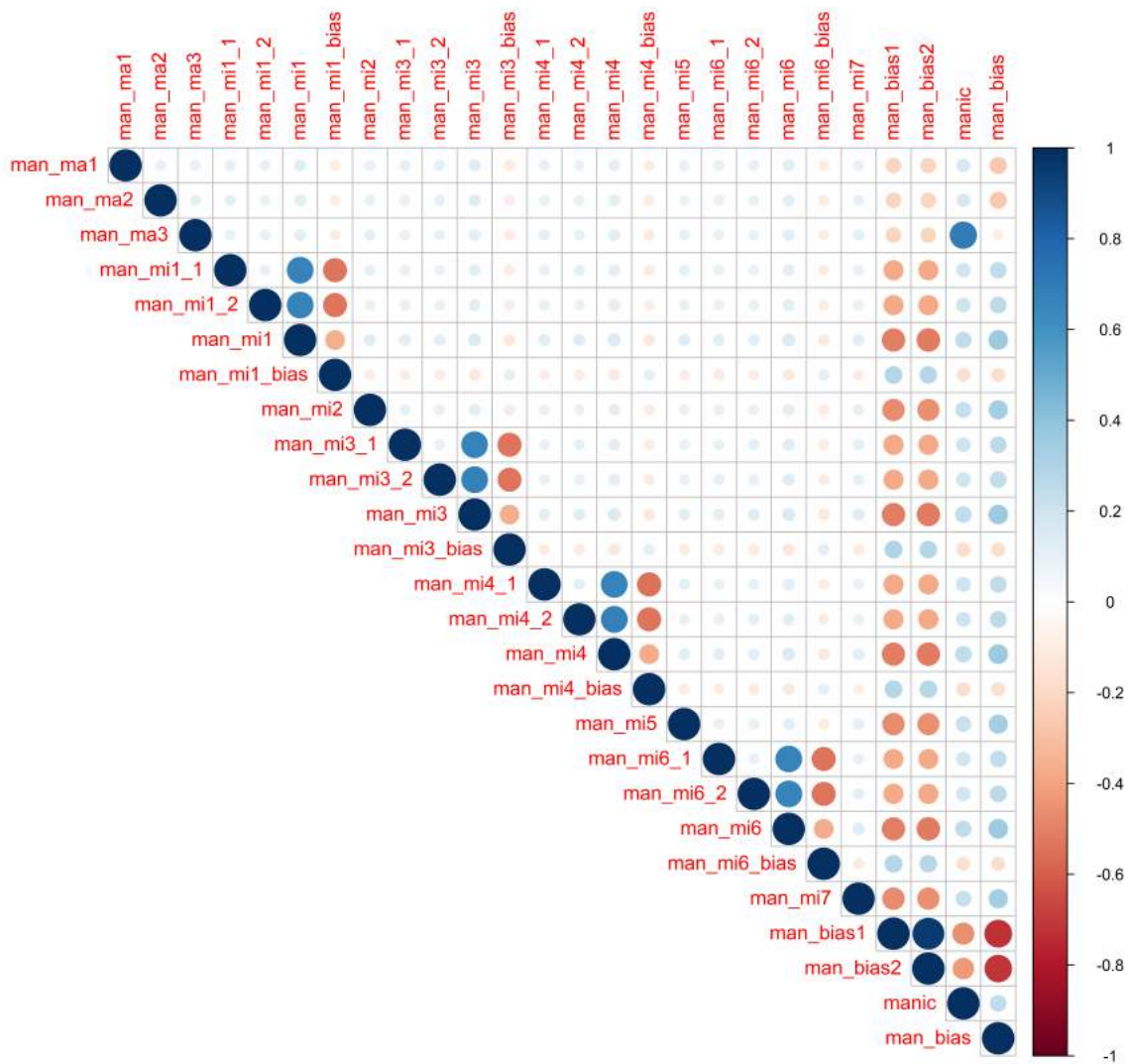
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



only

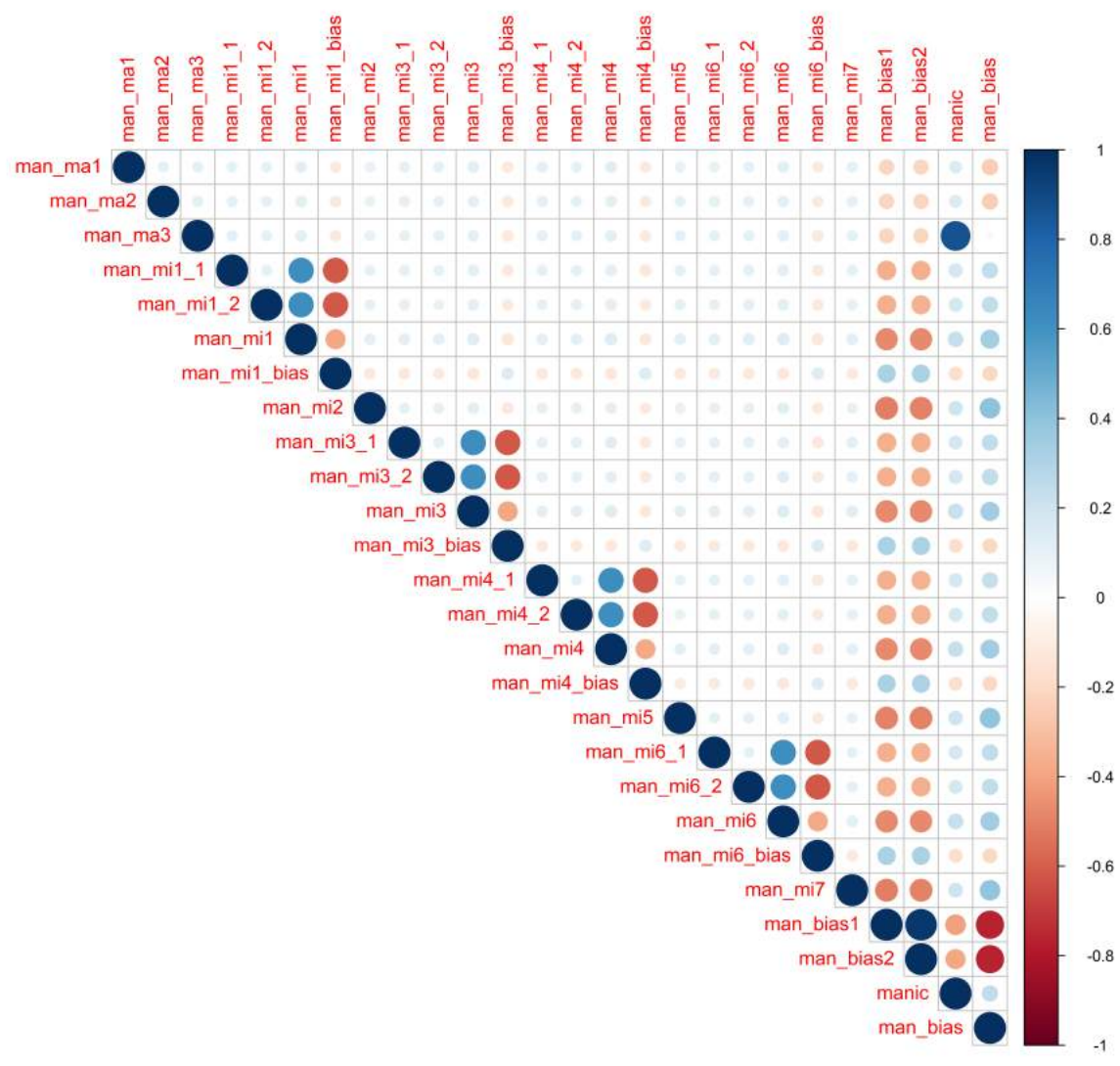


only

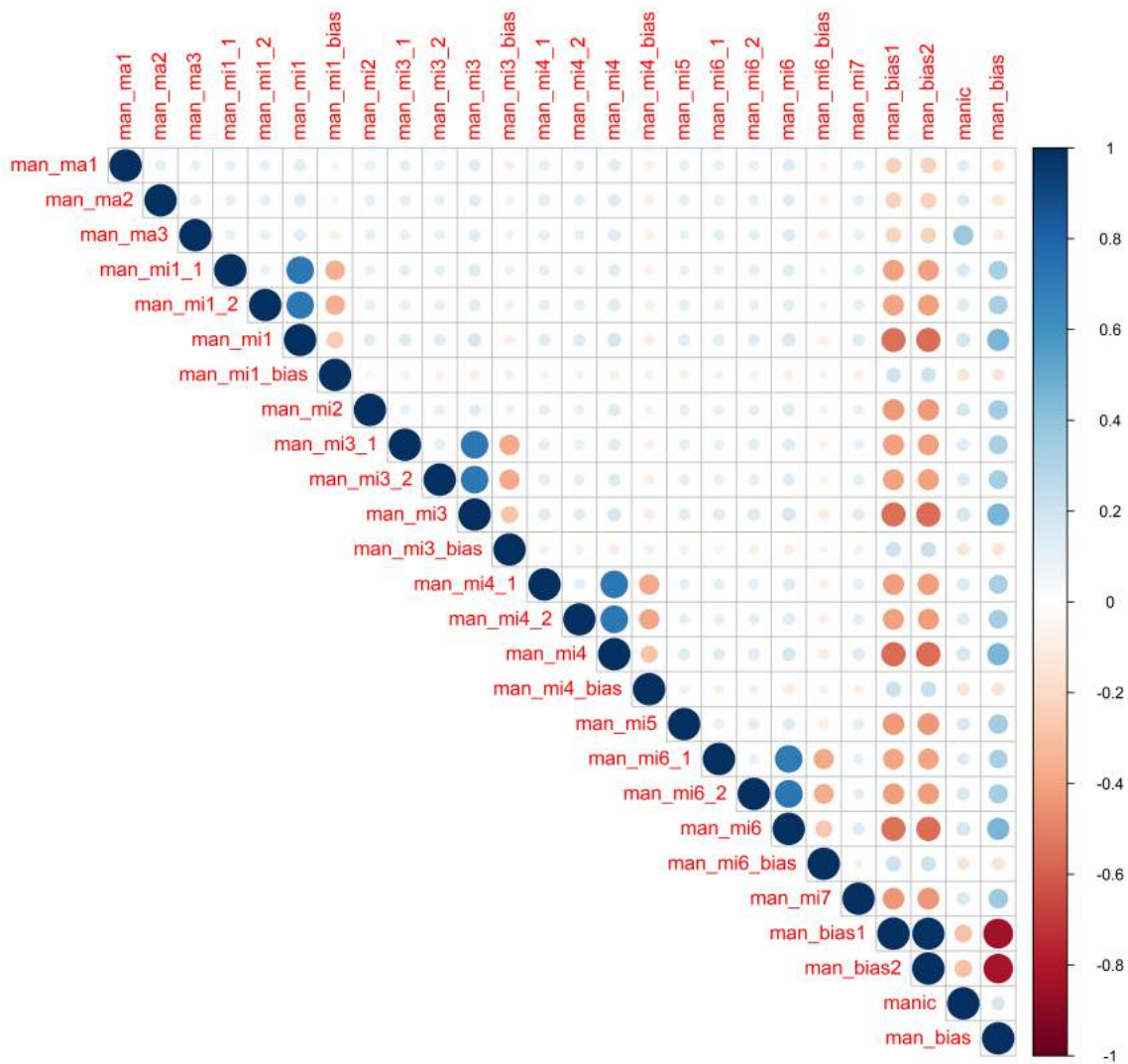


only

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

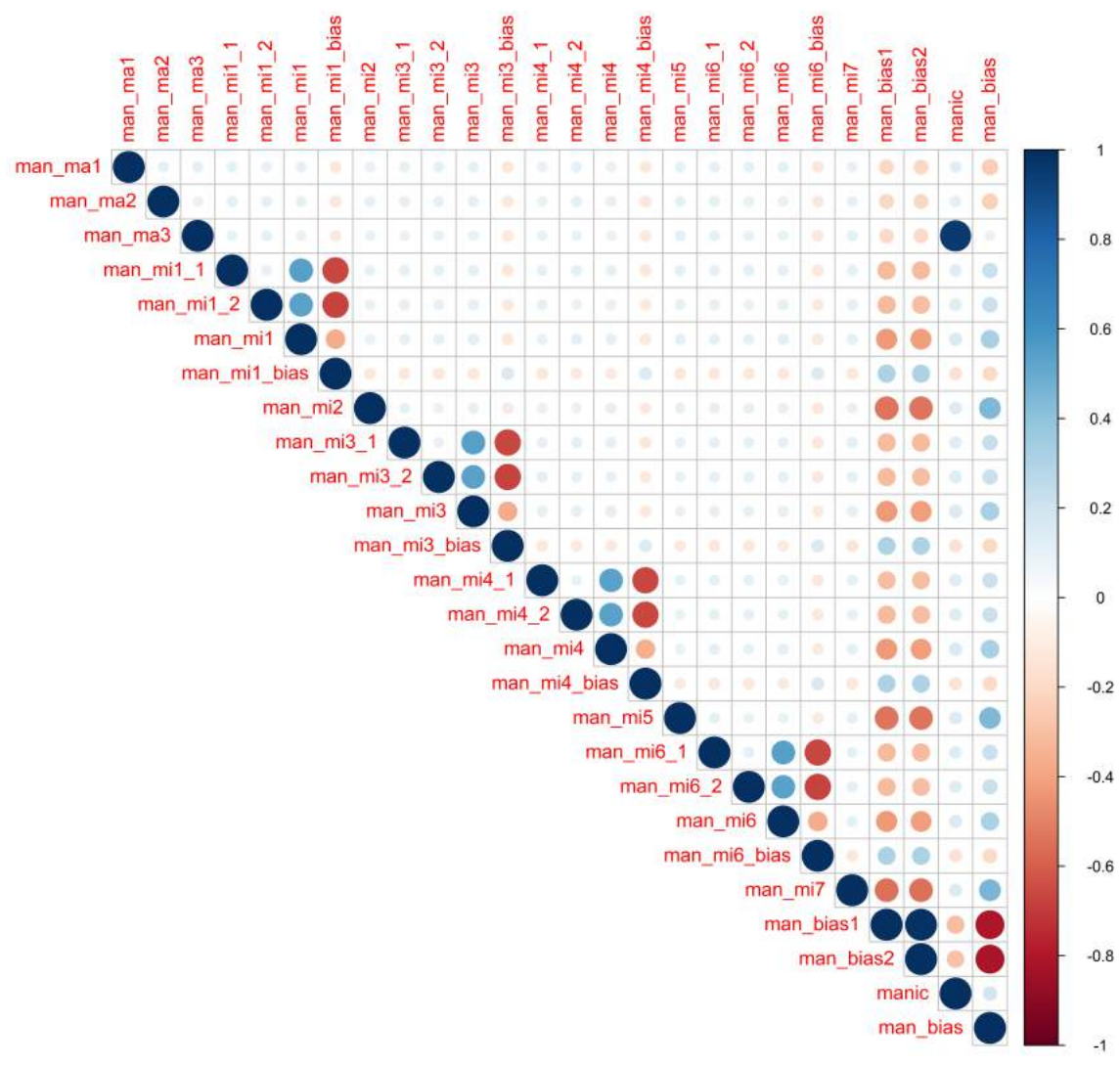


only

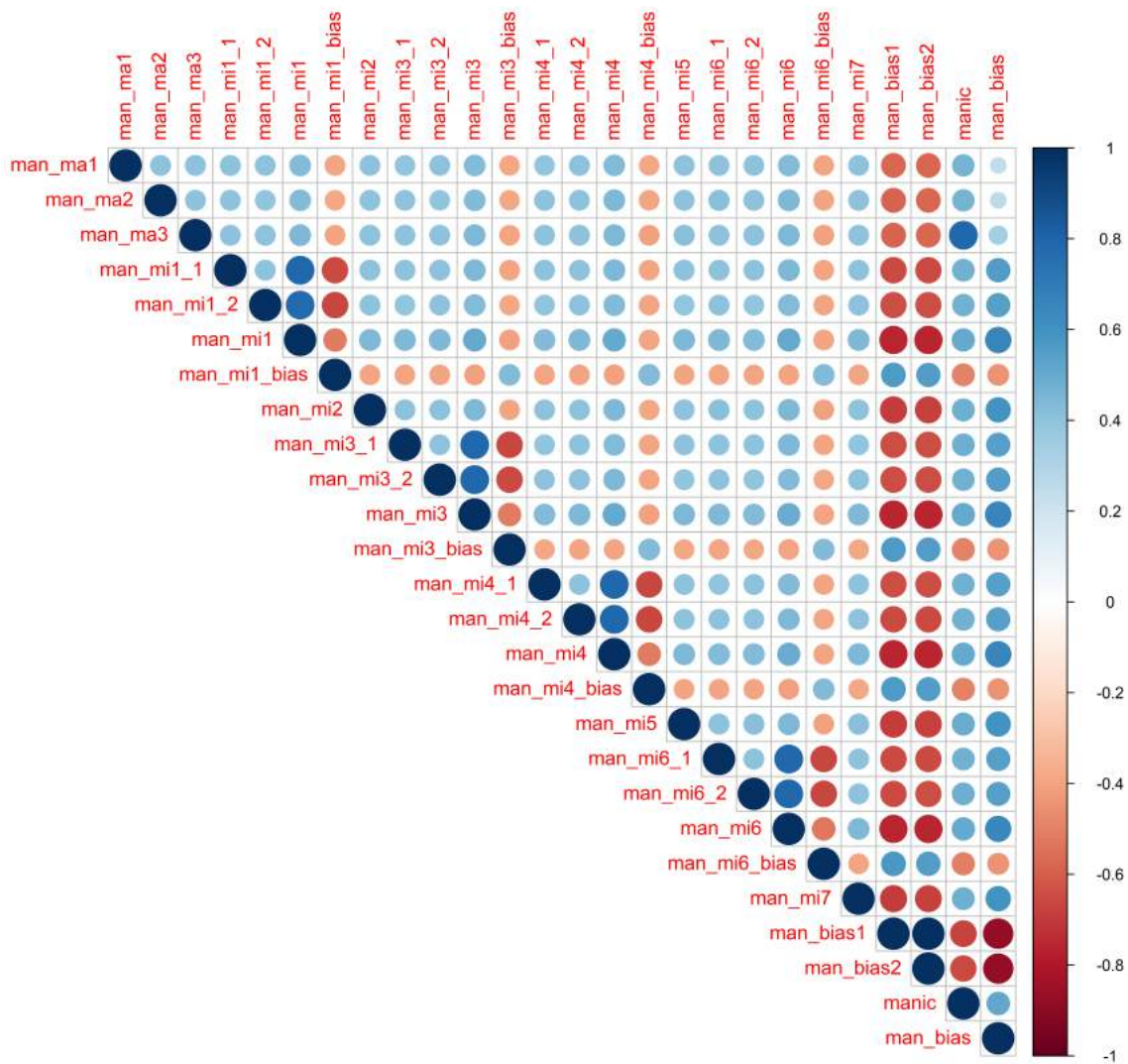


only

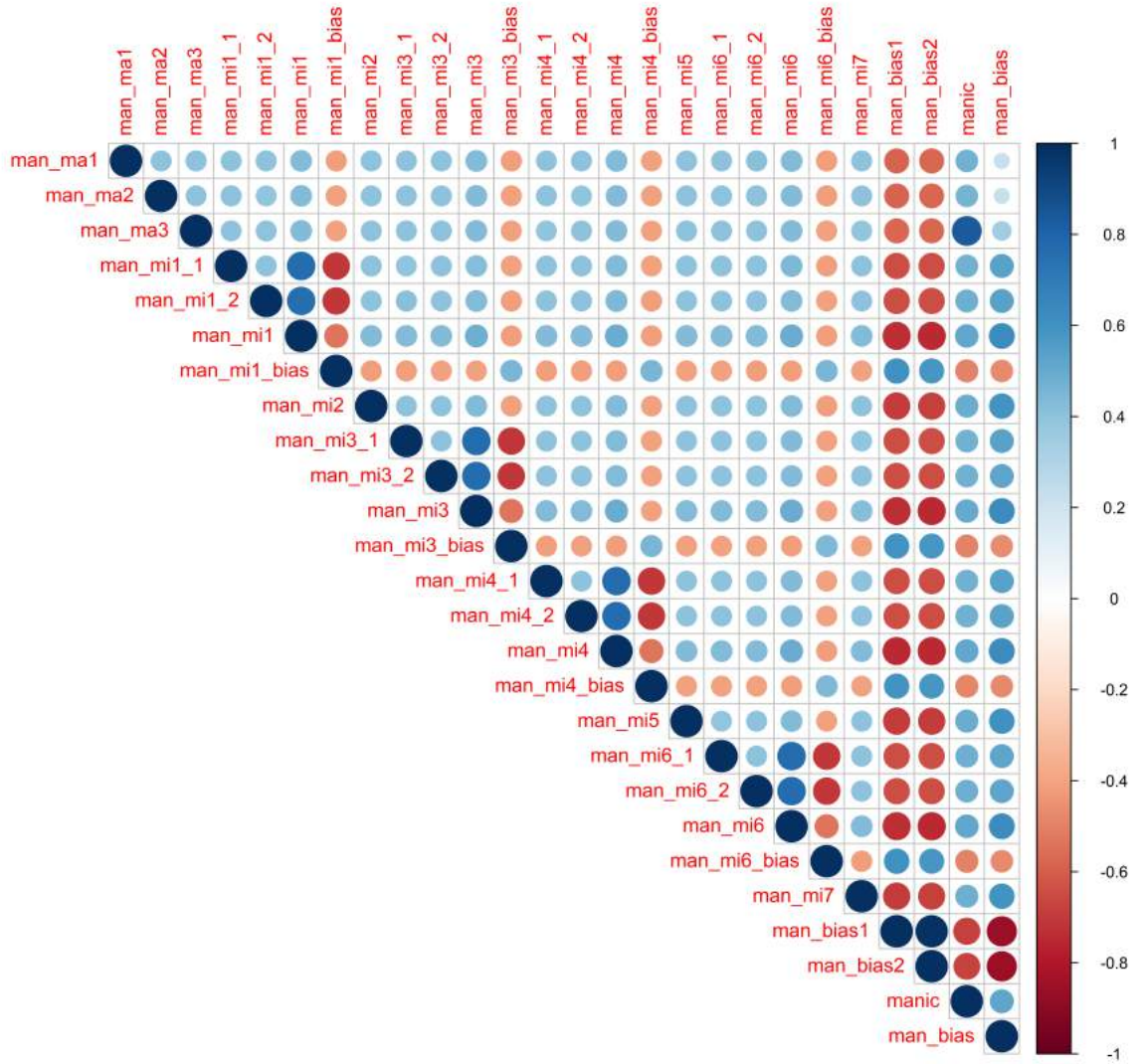
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



only

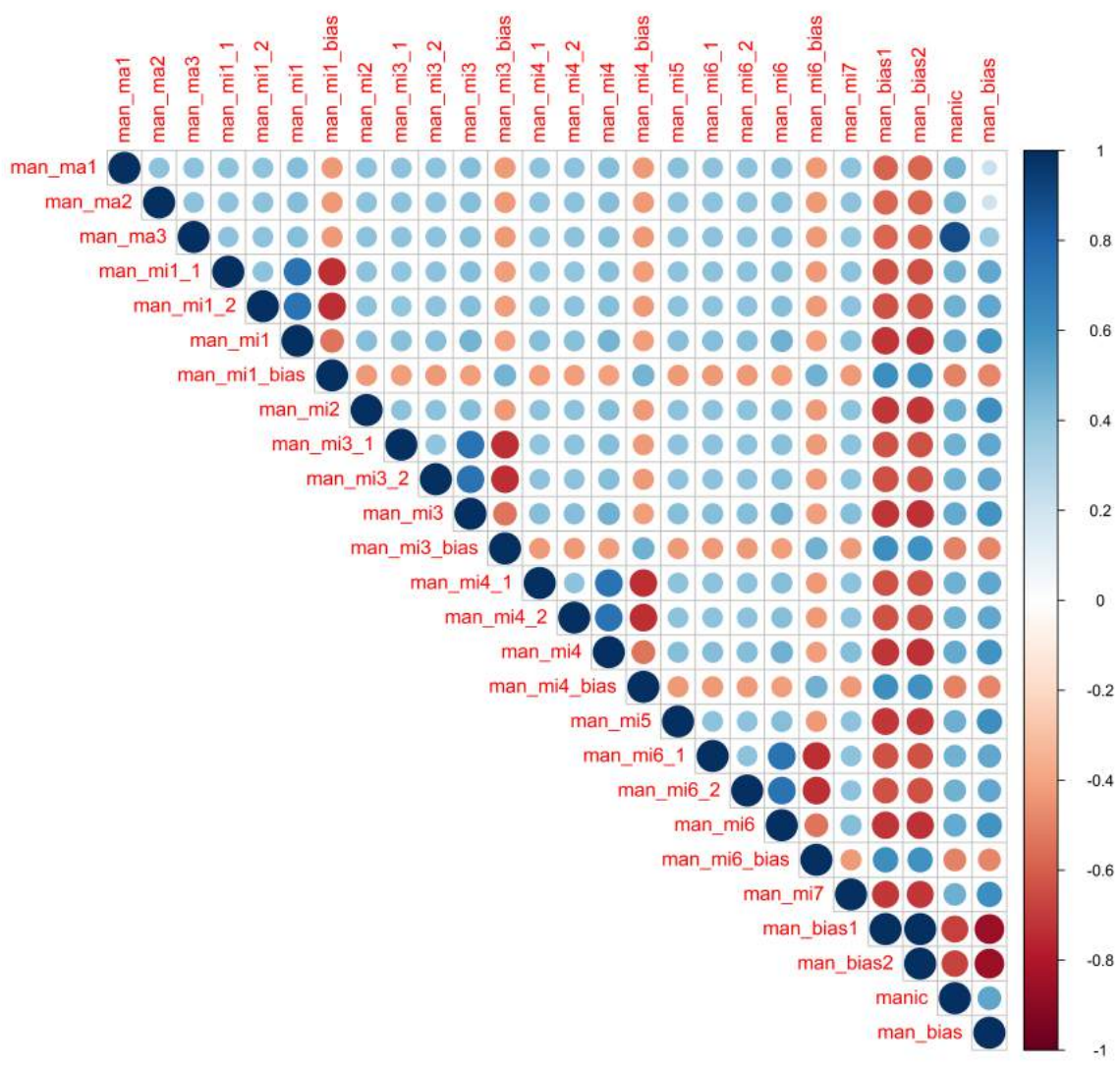


only



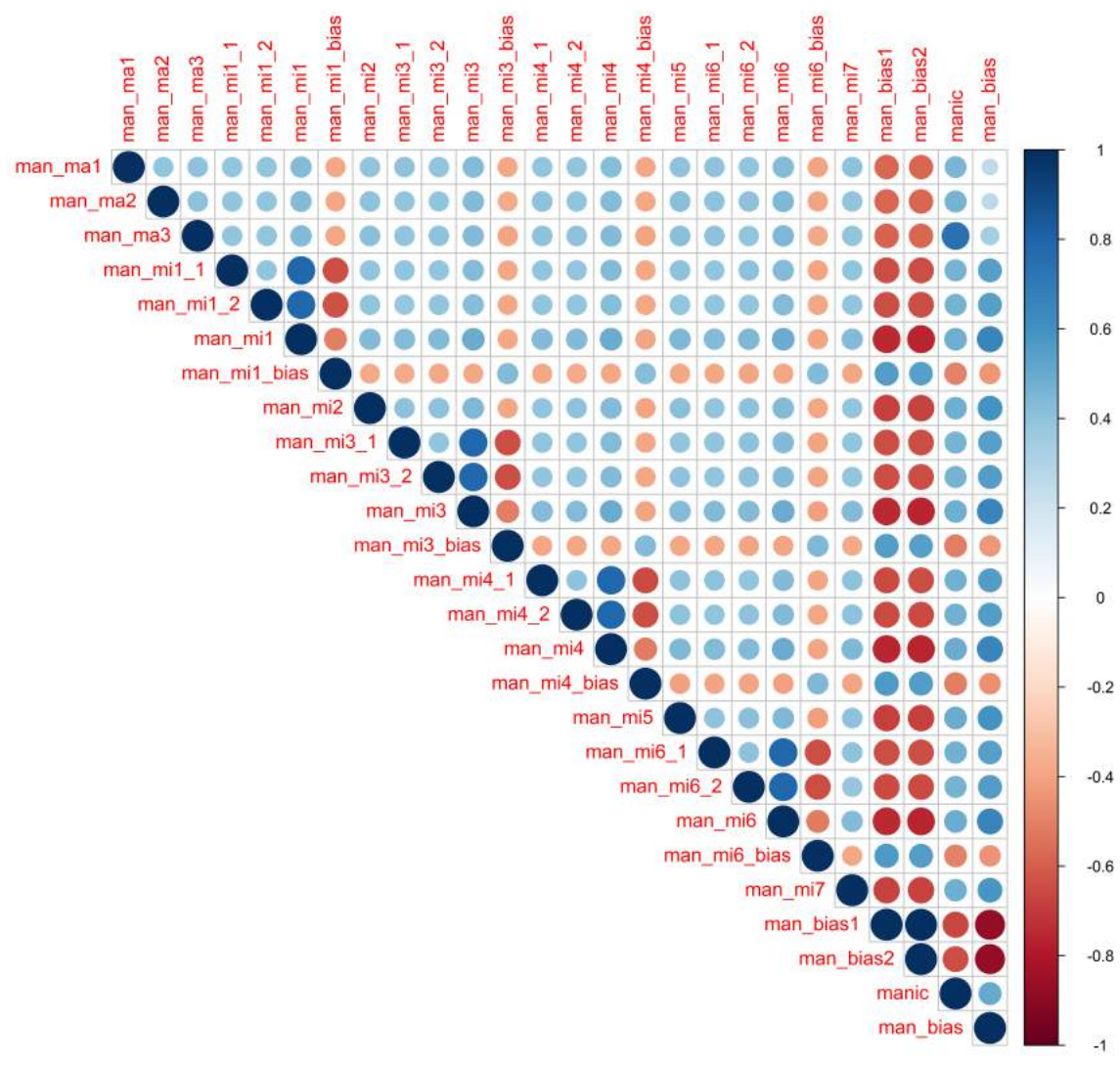
only

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

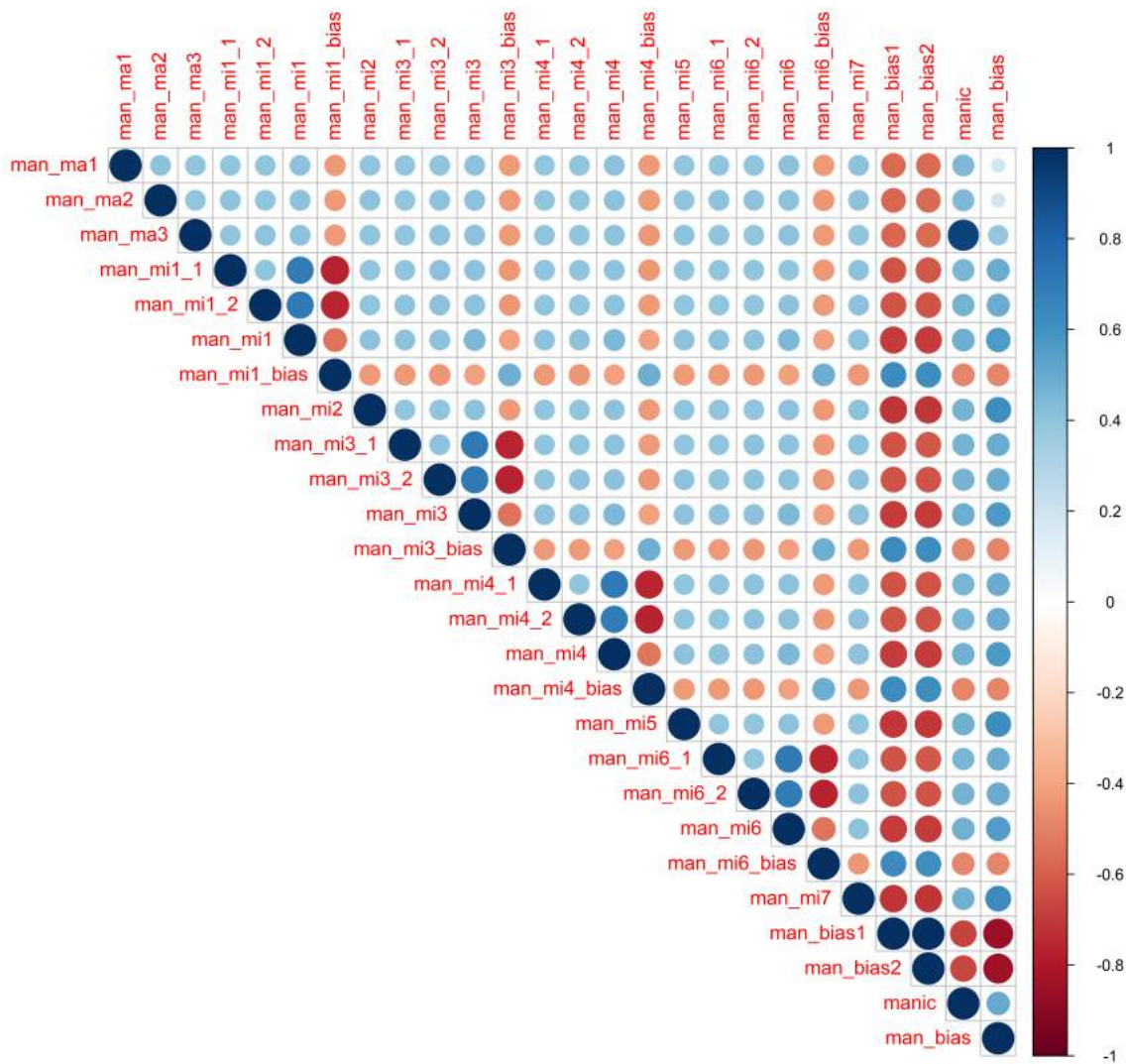


only

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



only



only

BMJ Open

A simulation study to demonstrate biases created by diagnostic criteria of mental illnesses: major depressive episodes, dysthymia, and manic episodes

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2020-037022.R1
Article Type:	Original research
Date Submitted by the Author:	24-Jun-2020
Complete List of Authors:	Chao, Yi-Sheng; Independent researcher Lin, Kuan-Fu; National Taiwan University Hospital Yun-Lin Branch, Psychiatry Wu, Chao-Jung; UQAM, Département d'informatique Wu, Hsing-Chien; Taipei Hospital, Internal Medicine Hsu, Hui-Ting; Changhua Christian Healthcare System, Pathology Tsao, Lien-Cheng; Changhua Christian Healthcare System, Surgery Cheng, Yen-Po; Changhua Christian Healthcare System, Surgery Lai, Yi-Chun; National Yang Ming University Hospital, Chest Medicine Chen, Wei-Chih; Taipei Veterans General Hospital, Chest Medicine; National Yang-Ming University, Institute of Emergency and Critical Care Medicine
Primary Subject Heading:	Mental health
Secondary Subject Heading:	Epidemiology, Research methods, Diagnostics
Keywords:	MENTAL HEALTH, Depression & mood disorders < PSYCHIATRY, EPIDEMIOLOGY, STATISTICS & RESEARCH METHODS

SCHOLARONE™
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

A simulation study to demonstrate biases created by diagnostic criteria of mental illnesses: major depressive episodes, dysthymia, and manic episodes

Yi-Sheng Chao¹, Kuan-Fu Lin,² Chao-Jung Wu³, Hsing-Chien Wu⁴, Hui-Ting Hsu⁵, Lien-Cheng Tsao⁵, Yen-Po Cheng⁵, Yi-Chun Lai⁶, Wei-Chih Chen^{7,8*}

¹Independent researcher, Montréal, H2X 0A8 Canada, ²National Taiwan University Hospital Yun-Lin Branch, Yunlin County, 640 Taiwan, ³Département d'informatique Université du Québec à Montréal, Montréal H3B 1B4 Canada, ⁴Taipei Hospital Ministry of Health and Welfare New Taipei city, 242 Taiwan, ⁵Changhua Christian Hospital, Changhua County 526, Taiwan, ⁶National Yang-Ming University Hospital, Yilan 260 Taiwan, ⁷Department of Chest Medicine, Taipei Veterans General Hospital, Taipei 112, Taiwan, ⁸Institute of Emergency and Critical Care Medicine, National Yang-Ming University, Taipei 112, Taiwan

*Corresponding author: Wei-Chih Chen
Department of Chest Medicine, Taipei Veterans General Hospital
No. 201, Section 2, Shih-Pai Road, Taipei 11217, Taiwan
Telephone: +886-2-28757456
Fax: +886-2-28757610
Email address: wiji.chen@gmail.com

Keywords: frailty; bias; forward-stepwise regression; the Health and Retirement Study; index mining

28 Abstract

29 Objectives

30 Composite diagnostic criteria alone are likely to create and introduce biases into diagnoses
31 that subsequently have poor relationships with input symptoms. This study aims to
32 understand the magnitudes of biases created by diagnostic criteria alone and introduced into
33 the diagnoses of mental illnesses with large disease burdens (major depressive episodes,
34 dysthymic disorder, and manic episodes) and the relationships between the diagnoses and
35 the input symptoms.

36 Settings

37 General psychiatric care.

38 Participants

39 Without real-world data available to the public, 100,000 subjects were simulated and the
40 input symptoms were assigned based on the assumed prevalence rates (0.05, 0.1, 0.3, 0.5,
41 and 0.7) and correlations between symptoms (0, 0.1, 0.4, 0.7, and 0.9). The input symptoms
42 were extracted from the diagnostic criteria. The diagnostic criteria were transformed into
43 mathematical equations to demonstrate the sources of biases and convert the input
44 symptoms into diagnoses.

45 Primary and secondary outcomes

46 Biases due to data censoring or categorization introduced into the intermediate variables,
47 and the three diagnoses were measured. The relationships between the input symptoms
48 and diagnoses were interpreted using forward stepwise linear regressions.

49 Results

50 The prevalence rates of the diagnoses were lower than those of the input symptoms and
51 proportional to the assumed prevalence rates and the correlations between the input
52 symptoms. Certain input or bias variables consistently explained the diagnoses better than
53 the others. Except for zero correlations and 0.7 prevalence rates of the input symptoms for
54 the diagnosis of dysthymic disorder, the input symptoms could not fully explain the
55 diagnoses.

56 Conclusions

57 There are biases created due to composite diagnostic criteria and introduced into the
58 diagnoses. The design of the diagnostic criteria determines the prevalence of the diagnoses,
59 the relationships between the input symptoms, the diagnoses, and the biases. The
60 importance of the input symptoms has been distorted largely by the diagnostic criteria.

61 Trial registration

62 Not applicable

63 Strength and limitation

- 64 1. The prevalence of three mental illnesses was determined by the prevalence of the
65 input symptoms and modified by the diagnostic criteria and correlations between the
66 input variables in simulated populations.

- 1
- 2
- 3 67 2. Biases due to data censoring or categorization were created by the diagnostic criteria
- 4 68 and introduced into the intermediate variables and the three diagnoses of mental
- 5 69 illnesses in simulated populations.
- 6 70
- 7 71 3. The diagnostic criteria modified the importance of the input symptoms; certain input
- 8 72 symptoms or bias variables were weighted more than expected in simulated
- 9 73 populations.
- 10 74 4. The design of diagnostic criteria influenced the diagnosis prevalence. With the same
- 11 75 input symptom prevalence, dysthymic disorder was the most prevalent among three
- 12 76 illnesses. Major depressive episodes were the least prevalent.
- 13 77 5. This study is based on simulated data and needs to be verified with real-world data.
- 14
- 15
- 16
- 17
- 18
- 19
- 20
- 21
- 22
- 23
- 24
- 25
- 26
- 27
- 28
- 29
- 30
- 31
- 32
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- 48
- 49
- 50
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60

For peer review only

78 Background

79 The diagnoses of several mental illnesses in patients are often made based on a
80 variety of criteria. These criteria often involve symptoms reported by the patients.[1, 2] For
81 example, the diagnosis of major depressive disorder defined in the Diagnostic and Statistical
82 Manual of Mental Disorders, 4th Edition, Text Revision (DSM-IV-TR) requires at least one
83 major depressive episode.[1, 2] For each major depressive episode, the major criteria are
84 “*depressive mood and/or loss of interest or pleasure in life activities for at least 2 weeks*”.[1,
85 2] In addition to the major criteria, the patients need to report at least five of the nine
86 symptoms that “*cause clinically significant impairment in social, work, or other important*
87 *areas of functioning almost every day,*” including insomnia or hypersomnia and fatigue or
88 loss of interest.[1, 2] In other words, patients need to match both the major and minor criteria
89 before being diagnosed with a major depressive episode.

90 Historically this symptom-based diagnostic approach developed by Feighner et al.
91 has been widely accepted.[3, 4] Since then, mental illnesses can be diagnosed through
92 different sets of criteria. This approach is important because clinicians become capable of
93 screening important symptoms before diagnosing and treating patients accordingly. In fact,
94 these criteria can also be seen as composite measures that use multiple measures to
95 capture disorders that may not be quantified with single variables.[5, 6] Recent studies on
96 composite measures have found them problematic because biases can be introduced while
97 aggregating information from input variables.[6] The biases emerge while the sums of input
98 variables are censored or while input variables are transformed inadequately.[6, 7] In other
99 words, biases can be created when there is information in the composite measures that is
100 not explained by and unrelated to the input variables.[6] For example, categorizing
101 continuous variables considers individuals in the same group the same and disregards the
102 heterogeneity between those in the same categories.[6] Such practices induce biases and
103 decrease measurement precision.[6, 7]

104 Currently there is no extensive review on the existence of these biases created by
105 composite measures or medical diagnoses, and only selected diagnoses have been studied
106 for such biases. These biases have been proven vital to another symptom-based composite
107 measure, the diagnosis of frailty, a condition that often occurs in the elderly and is
108 significantly associated with health outcomes, such as mortality, falls, and morbidity.[6]
109 Frailty is diagnosed based on several symptoms and characterized by weakness and
110 vulnerability to adverse health events.[6] While using one of the most widely used diagnostic
111 criteria, the Biological Syndrome Model scores, to diagnose frailty,[8] biases alone can
112 explain more than 71% of the variances of the frailty diagnosis.[6] The biases introduced by
113 data censoring and data categorization can better explain the frailty diagnosis than the input
114 symptoms.[6]

115 Mostly designed as symptom-based composite measures, it is possible that the
116 diagnostic criteria of mental illnesses also create and introduce biases into diagnoses so that
117 the diagnoses could not be fully explained by the input symptoms. In concern of the biases
118 created by the diagnostic criteria alone, this study aims first to understand the relationships
119 between mental symptoms and diagnoses and then to quantify the potential role of the
120 biases regarding the diagnoses by simulating populations with different prevalence rates and
121 between-variable correlations of mental symptoms.

122 Methods

123 Assumptions and simulation parameters

124 A file containing R codes to reproduce the simulations was attached in the
125 supplementary file. Simulated populations with mental symptoms of different prevalence
126 rates and between-variable correlations were created to interpret the diagnoses and
127 understand the potential magnitudes of biases that could be introduced via data processing
128 implied by the diagnostic criteria (reproducible using R codes and data in the supplementary
129 file). Three diagnoses of mental illnesses were chosen for the leading associated disease
130 burdens:[2] major depressive episodes for the diagnosis of major depressive disorder,
131 dysthymic disorder, and manic episodes for the diagnosis of bipolar disorder.[1]

132 There were assumptions made to simulate the populations (Table 1). First, for each
133 simulation, the prevalence rates of the input symptoms were assumed to be similar for the
134 three diagnoses in this study. Second, the input symptoms for the diagnoses of major
135 depressive episodes and dysthymic disorder correlated with the same correlation
136 coefficients.[9] The symptoms for the diagnosis of manic episodes correlated to one another.
137 Third, the input symptoms for the diagnosis of manic episodes were created independently
138 of those for the diagnosis of the other two mental illnesses. The assumptions of the
139 prevalence rates and between-variable correlations were made because there was no
140 acceptable-quality data on the symptoms of mental illnesses published. There were studies
141 on the prevalence of mental illnesses,[10, 11] but the information on the prevalence of
142 mental symptoms was very limited. There were variables about depression or anxiety
143 collected in national surveys, such as the items collected through the Center for
144 Epidemiologic Studies Depression Scale.[6, 12-18] However, these variables were not the
145 symptoms used in the DSM-IV-TR. Lastly, we assumed that the diagnoses were made
146 accurately based on the input symptoms reported precisely by patients. The diagnostic
147 criteria in the DSM-IV-TR were strictly followed. However, these assumptions did not hold in
148 the real world.[19] For simplicity and practicality reasons, we assumed perfect diagnostic
149 quality by physicians and accurate reporting of the input symptoms by patients in the
150 simulated populations.

151 Diagnostic criteria as mathematical functions

152 The input symptoms were extracted from the major and minor criteria of the diagnoses and
153 listed in Table 2 to Table 4. The input symptoms, major and minor criteria, and the
154 diagnoses were assigned variable names. All input symptoms, items or domains in the major
155 or minor criteria, and the diagnoses were binomial variables, presenting zero and one for the
156 absence and presence of the symptoms, criteria, and the diagnoses respectively. For
157 example, “*insomnia*” and “*hypersomnia*” were extracted from one of the minor criteria for the
158 diagnosis of major depressive episodes.[1] “*More talkative than usual*” and “*pressure to keep*
159 *talking*” were extracted from one of the minor criteria for the diagnosis of manic episodes.[1]

160 Mathematical functions were generated based on the diagnostic criteria to convert input
161 symptoms into diagnoses. For example, one of the minor criteria of dysthymic disorder was
162 “*poor appetite or overeating*.” This required two input symptoms and one bias variable to
163 generate the criterion.[6] In other words, “*poor appetite or overeating*” equaling the sum of
164 two input variables, “*poor appetite*” and “*overeating*,” and a bias variable to achieve
165 censoring of the sum of both variables.[6] The sum of two binomial variables could be zero,

1
2
3 166 one and two for the subjects. However, to derive a binomial variable (having at least one
4 167 symptom) based on a distribution of 0 to 2, the bias variable had values of -1 for subjects
5 168 with both symptoms to obtain values less than or equal to one in all subjects.[6] Therefore,
6 169 the bias variable had values of -1 for the subject with both symptoms and 0 for the other
7 170 subjects. In addition to adding variables together to derive an intermediate variable or a
8 171 diagnosis, multiplication, categorization, and other more complicated methods were used in
9 172 the diagnostic criteria to generate diagnosis variables and domain variables in the major or
10 173 minor criteria.

13
14 174 For example, the diagnosis of dysthymic disorder required the confirmation of both
15 175 the major criteria, "*depressed mood most of the day for more days than not, for at least 2*
16 176 *years*" and the minor criteria, "*the presence of two or more of the following symptoms,*" at the
17 177 same time.[1] The diagnosis based on whether subjects meeting both the major and minor
18 178 criteria of dysthymic disorder is the same as identifying those with a multiplicative product of
19 179 1 of two binomial variables (0 and 1 for absence and presence of the major or minor criteria).
20 180 In the equations, two binomial variables were multiplied to obtain the diagnosis of dysthymic
21 181 disorder among those with a multiplicative product of 1. Individuals could be assigned zero
22 182 or one for whether they met both criteria, while the sum of major and minor criteria were
23 183 zero, one or two for the individuals. Linearly, a bias variable with values of -1 or zero was
24 184 created and those meeting the major or minor criteria were assigned -1.[6] For
25 185 categorization of continuous variables, bias variables were required to remove the variations
26 186 between the subjects in the same categories.[6] Other equations to generate the
27 187 intermediate variables and the diagnoses were listed and explained in Table 2 to Table 4.

31 188 **Generation of bias variables**

32 189 Bias variables could be generated while binomial input symptoms were summed or
33 190 multiplied to obtain binomial intermediate or diagnosis variables (see the example in the
34 191 previous two paragraphs).[6] A visual presentation of how bias variables were generated
35 192 was published.[6] Therefore the number of bias variables depended on the complexity of
36 193 how the diagnoses were made. For example, six of the nine items or domains in the minor
37 194 criteria for the diagnosis of major depressive episodes were the censored sums of the input
38 195 symptoms and six bias variables were derived along with the intermediate variables that
39 196 represented the items in the minor criteria. The other bias variables were described in Table
40 197 2 to Table 4.

44 198 **Simulation parameters and simulated populations**

45 199 We simulated populations of 100,000 subjects. There were five prevalence rates to
46 200 simulate the input symptoms for the diagnosis of major depressive episodes, dysthymic
47 201 disorder, and manic episodes: 0.05, 0.1, 0.3, 0.5, and 0.7. The correlations between the
48 202 input symptoms were hypothesized to be 0, 0.1, 0.4, 0.7, and 0.9. There were 25
49 203 combinations of the assumed prevalence rates and between-variable correlations. The
50 204 presence of the input symptoms was randomly assigned to the subjects after specifying the
51 205 prevalence rates and between-variable correlations between the input symptoms.[20, 21]
52 206 The intermediate and diagnosis variables were derived according to the equations in Table 2
53 207 to Table 4. For each combination of prevalence rates and between-variable correlations, the
54 208 populations were simulated for 100 times to obtain the mean values and 95% confidence
55 209 intervals (CIs) of derived prevalence rates, as well as the adjusted R squared and p values
56 210 to approximate the diagnosis variables.

211 **Diagnosis approximation**

212 Due to the existence of the biases, the input symptoms were not likely to fully explain
213 the diagnoses.[6] Therefore, the diagnoses were approximated by the input, bias, and
214 intermediate variables individually and collectively.[6, 12, 14, 16] The approximation was
215 conducted using forward-stepwise linear regressions.[6, 12, 14, 16, 22] The interpretability of
216 the diagnoses by the input symptoms and bias variables was assessed via adjusted R
217 square: zero suggested that the input symptoms were unrelated to the diagnosis, and one
218 suggested that the input symptoms perfectly explained the diagnosis.[14, 15, 23-26]

219 All statistical analyses were conducted within the R environment (v3.4.1)[27] and
220 RStudio (v1.0.153).[28] Two-tailed P values less than 0.05 were considered statistical
221 significant.

222 **Patient and Public Involvement**

223 This is a simulation study that did not involve patients or human subjects.

224 **Results**

225 The derived prevalence rates of the input symptoms for the three mental illnesses
226 matched the assumed rates in the supplementary file. The derived correlations between the
227 input symptoms were close to assumed levels in the supplementary file. The simulations
228 were successful and accurate based on the assumed prevalence rates and correlations.

229 **Prevalence of intermediate variables**

230 The items in the major and minor criteria were the intermediate variables necessary
231 to create the diagnoses. The methods used to generate the intermediate variables were
232 important for the prevalence rates of the intermediate variables and the derived diagnoses in
233 Figure 1. For example, an intermediate variable, "*significant unintentional weight loss or*
234 *gain*," was created by summing and censoring two binomial variables with values of zero
235 and one (significant unintentional weight loss; significant unintentional weight gain). The
236 prevalence rates of the intermediate variables were larger than those of the two input
237 symptoms regardless of the assumed prevalence rates or between-variable correlations of
238 the input symptoms.

239 In contrast, the diagnosis of dysthymic disorder was a multiplication product of two
240 intermediate binomial variables, the major and minor criteria, and the prevalence rates of
241 dysthymic disorder were lower than those of the major or minor criteria under all
242 combinations of assumed correlations and prevalence rates in Figure 2.

243 **Prevalence of mental illnesses**

244 The derived prevalence rates of three diagnoses were plotted against the assumed
245 prevalence rates and correlations of the input symptoms in Figure 2 to Figure 4 and listed in
246 Table 5. None of the three diagnoses had prevalence rates exceeding those of the input
247 symptoms. In general, higher prevalence rates or between-variable correlations of the input
248 symptoms were associated with higher prevalence rates in the three diagnoses, except for
249 manic episodes that had higher prevalence rates (0.692) assuming zero correlations and 0.7
250 prevalence rates than the prevalence rate (0.679) assuming 0.1 correlations and 0.7
251 prevalence rates of the input symptoms. When compared across Figure 2 to Figure 4, given
252 the same assumed prevalence rates and between-variable correlations of the input

1
2
3 253 symptoms, the diagnostic criteria of dysthymic disorder consistently generated diagnoses of
4 254 the highest prevalence rates and the criteria of major depressive episodes created
5 255 diagnoses of the least prevalence rates (see Table 5 for details).

8 256 **Associations between the diagnoses and individual input symptoms or** 9 257 **bias variables**

10 258 The diagnoses were first interpreted with the input symptoms (including intermediate
11 259 variables) and the bias variables individually. The diagnosis of dysthymic disorder, for
12 260 example, was interpreted with the input symptoms, the bias variables, and both in Figure 5.
13 261 For each simulation, the diagnosis of dysthymic disorder was approximated with an
14 262 increasing number of the input symptoms, the bias variables, or both. After selecting the
15 263 variables that best approximated the diagnosis based on adjusted R-squared, the input
16 264 symptoms could explain a proportion of 0.955 of the diagnosis variance and the bias
17 265 variables could explain at most a proportion of 0.405 of the diagnosis variance in Figure 5.
18 266 With all variables used in the regression, the diagnosis could be perfectly explained by the
19 267 input symptoms and bias variables (adjusted R-squared = 1). The individual input symptoms
20 268 and the bias variables that individually best explained the diagnoses are listed in Table 6
21 269 and Table 7, respectively.

22 270 For the diagnosis of major depressive episodes, the first and second items in the major
23 271 criteria (variable names: mde_ma1 for or mde_ma2 in Table 2) individually best explained
24 272 the diagnosis depending on the assumed prevalence rates and correlations in Table 6. For
25 273 the diagnosis of dysthymic disorder, the major criteria (dys_ma in Table 3) consistently and
26 274 individually explained the diagnosis the best. For the diagnosis of manic episodes, the third
27 275 item of the major criteria (man_ma3 in Table 4) individually best explained the diagnosis in
28 276 all combinations of assumed prevalence rates and correlations. However, the proportions of
29 277 diagnosis variances best explained by individual input symptoms varied widely between
30 278 0.001 to 0.974, depending on the assumed prevalence rates and between-variable
31 279 correlations. Based on the adjusted R-squared for individual input symptoms, certain input
32 280 variables were more important than other symptoms due to a high correlation with the
33 281 diagnoses, such as the major criteria for the diagnosis of dysthymic disorder. The
34 282 prevalence rates and between-variable correlations were important to determine the
35 283 relationships between input symptoms and diagnoses.

36 284 Similarly, there were bias variables that consistently best explained the diagnoses in Table
37 285 7. For the diagnosis of major depressive episodes, the biases due to categorization of the
38 286 numbers of confirmed input symptoms (mde_bias1 and mde_bias2 in Table 2) were the
39 287 leading bias variable. The diagnosis of major depressive episodes not explained by the input
40 288 symptoms or information censoring (mde_bias in Table 2) was the leading bias variable in
41 289 two combinations of the assumed prevalence rates and correlations. For the diagnosis of
42 290 dysthymic disorder, the residual of the diagnosis not explained by the major and minor
43 291 criteria (dys_bias in Table 3) and the bias due to the categorization of the confirmed input
44 292 symptoms in the minor criteria (dys_mi_bias) were the leading bias variables. For the
45 293 diagnosis of manic episodes, the bias due to the categorization of the number of confirmed
46 294 input symptoms in the minor criteria up to three (man_bias1 in Table 4) was the leading bias
47 295 variables, except for two combinations of the assumed prevalence rates and correlations, in
48 296 which the bias due to categorization of the confirmed input symptoms in the minor criteria up
49 297 to four (man_bias2 in Table 4) best explained the diagnosis. However, the proportions of

1
2
3 298 diagnosis variances explained by individual bias variables varied widely from zero to 0.87.
4 299 Depending on the assumed prevalence rates and between-variable correlations of the input
5 300 symptoms, certain bias variables were more important than other bias variables and even
6 301 some input variables. The assumed prevalence rates and between-variable correlations
7 302 were important factors for the relationships between the bias variables and the diagnoses.
8
9

10 303 In general, the proportions of the diagnosis variance that could be explained by either
11 304 individual input symptoms or single bias variables were low when the prevalence rates and
12 305 between-variable correlations of the input symptoms were assumed to be low. With higher
13 306 assumed prevalence rates or correlations, the proportions of the diagnoses explained by the
14 307 single input symptoms or bias variables were higher. Across three diagnoses, the diagnosis
15 308 of dysthymic disorder could be better explained by single input variables (higher adjusted R-
16 309 squared), and the diagnosis of major depressive episodes was associated with the least
17 310 adjusted R-squared. The bias variables of the diagnosis of manic episodes could explain the
18 311 diagnosis individually better than the bias variables of the other two diagnoses.
19
20
21

22 312 **Approximating the diagnoses with input symptoms**

23 313 When the diagnoses were approximated by their own input symptoms (Table 8),
24 314 there were always some diagnosis variances that could not be explained by the input
25 315 symptoms. In other words, the input symptoms could not fully explain the diagnoses, except
26 316 for the diagnosis of dysthymic disorder that could be fully explained by the input symptoms
27 317 (adjusted R-squared = 1) assuming zero between-variable correlations and 0.7 prevalence
28 318 rates for the input symptoms. In Table 8, the proportions of diagnosis variances explained by
29 319 input symptoms increased with higher assumed prevalence rates or between-variable
30 320 correlations of the input symptoms in general. The input symptoms of dysthymic disorder
31 321 explained the diagnosis better than those of the other two diagnoses under all combinations
32 322 of assumed prevalence rates and between-variable correlations. However, the proportion of
33 323 diagnosis variance explained by own input symptoms varied widely from 0.003 to 1.0. The
34 324 assumed prevalence rates and between-variable correlations of the input symptoms and the
35 325 design of the diagnostic criteria were all important for the relationships between input
36 326 symptoms and diagnoses.
37
38
39
40

41 327 **Approximating the diagnoses with bias variables**

42 328 The diagnoses were approximated with the bias variables of their own. The bias
43 329 variables always explained some of the diagnosis variances, except for the diagnosis of
44 330 dysthymic disorder assuming zero between-variable correlations and 0.7 prevalence rates
45 331 for the input symptoms (adjusted R-squared = 0). With increasing assumed between-
46 332 variable correlations for the input symptoms, the adjusted R-squared increased. However,
47 333 given the same assumed between-variable correlations, the proportions of diagnosis
48 334 variances explained by the bias variables might increase or decrease with the assumed
49 335 prevalence rates. Compared to the adjusted R-squared in Table 8, the proportion of the
50 336 diagnosis variances explained by the bias variables was always smaller than that explained
51 337 by the input symptoms in Table 9. However, the proportions of the diagnosis variance
52 338 explained by bias variables also varied widely from zero to 0.89. The assumed prevalence
53 339 rates and between-variable correlations of input symptoms and the design of the diagnostic
54 340 criteria were important for the relationship between the bias variables and the diagnoses.
55 341 Only when the input symptoms for the diagnosis of dysthymic disorder were randomly and
56
57
58
59
60

342 independently prevalent to 70% of the simulated populations, the bias variables became
343 irrelevant to the diagnosis.

344 Discussion

345 This study is a first attempt to assess the biases created by mental illness diagnostic
346 criteria, as well as understand the relationships between input symptoms and the diagnoses
347 of three mental illnesses: major depressive episodes (at least one episode required for the
348 diagnosis of major depressive disorder), dysthymic disorder, and manic episodes. The
349 diagnostic criteria of these three mental illnesses have been reviewed and rewritten as
350 mathematical functions. Simulated populations of 100,000 for each of 100 simulations, with
351 input symptoms of the three diagnoses, were created. For simplicity and practicality, the
352 presence of the input symptoms was randomly assigned, and the input symptoms were
353 assumed to have uniform prevalence rates and between-variable correlations. There were
354 25 combinations of assumed prevalence rates and between-variable correlations simulated.

355 Mathematically, the diagnostic criteria are functions and composite measures to
356 transform information from the input symptoms to diagnoses. There are bias variables
357 created by the diagnostic criteria due to data processing.[6] There are three major
358 mechanisms of introducing biases: censoring, data categorization,[7] and multiplication of
359 input symptoms.[6] These mechanisms introduce information or biases that cannot be fully
360 explained by the input symptoms.[6] The introduced biases can sometimes explain more
361 than half of the variance in the diagnoses depending on the prevalence rates and between-
362 variable correlations of the input symptoms. The findings show that the design of the
363 diagnostic criteria is important for bias introduction and significant for the prevalence of the
364 diagnoses in populations, the relationships between the input symptoms and the diagnoses,
365 and the relationships between the bias variables and the diagnoses.

366 The role of the diagnostic criteria

367 With the same assumptions in the prevalence rates and between-variable
368 correlations of the input symptoms, the design of the diagnostic criteria of three mental
369 illnesses can be compared to each other. The design of diagnostic criteria transform input
370 symptoms into various diagnosis prevalence rates with implicit upper limits (i.e. no more
371 prevalent than the input symptoms), unacknowledged differential weights on the input
372 symptoms (i.e., certain input symptoms better explaining the diagnoses), and the
373 introduction of biases (i.e., due to censoring, data categorization, or multiplication).

374 We were the first to notice that the prevalence rates of the three diagnoses were
375 lower than those of the input symptoms if randomly distributed with uniform prevalence rates
376 and correlations. Given similar assumed input symptom prevalence and correlations,
377 dysthymic disorder is the most prevalent, and major depressive episodes are the least. The
378 diagnosis of dysthymic disorder can be better explained by its input symptoms individually or
379 collectively than the other two diagnoses. The diagnosis of major depressive episodes is
380 least explained by own input symptoms individually or collectively. As expected, the
381 diagnosis of the three mental illness is similar to composite measures or indices and is
382 subject to the biases introduced by data processing, given all combinations of the assumed
383 prevalence rates and between-variable correlations of the input symptoms.[6] There is only
384 one exception: dysthymic disorder with the input symptoms that are randomly and

1
2
3 385 independently present in 70% of the population. This is because the diagnosis of dysthymic
4 386 disorder is a multiplicative product of the major and minor criteria. Without correlations,
5 387 everyone in the population is certain to qualify for the minor criteria (probability of 100%
6 388 because having at least two out of the six items in the minor criteria: mathematically $[C(2,6)$
7 389 $+ C(3,6) + C(4,6) + C(5,6) + C(6,6)] \times (0.7)^6 = 37 \times 0.117 = 4.35 > 100\%$). If 70% of the
8 390 population were also randomly assigned with the major criteria and 100% were assigned
9 391 with the minor criteria, 70% would be diagnosed with dysthymic disorder, and the diagnosis
10 392 of dysthymic disorder can be fully explained by the major criteria. In fact, without correlations
11 393 between input symptoms, it only requires each of the six items in the minor criteria to be
12 394 randomly assigned to 54.8% $[(1/37)^{(1/6)}]$ of the population for everyone to qualify for the
13 395 minor criteria, and the diagnosis can be fully explained by the minor and major criteria.

396 **Distortion of the input symptoms**

397 The importance of the input symptoms has been distorted due to the diagnostic
398 criteria for the three mental illnesses. The same phenomenon has been proven in the
399 diagnosis of frailty based on three of the most commonly used scoring methods.[6] In other
400 words, based on the functions to generate the diagnoses, the input symptoms are
401 differentially weighted, and weights are not explicitly acknowledged. The most prominent is
402 the diagnosis of dysthymic disorder; more than 90% of the variance can be explained by its
403 major criteria assuming 0.7 or 0.9 between-variable correlations for the input symptoms in
404 Table 6. Another example is that the third item of the major criteria for the diagnosis of manic
405 episodes, "*irritable mood*," individually predicts the diagnosis better than any other input
406 symptoms or intermediate variables. The input symptom has been given more weight than
407 others and can explain more than 91.8% of the diagnosis variance, assuming 0.9
408 correlations between input symptoms. Based on the texts in the DSM-IV-TR, we do not think
409 this symptom should be emphasized to this degree. However, the diagnostic criteria impose
410 implicit and unequal weights to the input symptoms, and introduce biases into the
411 diagnoses.

412 **Future directions**

413 We think it important to rethink the role and importance of the diagnostic system.
414 Current approaches are embedded with implicit assumptions of the prevalence rates of the
415 diagnoses (no higher than input symptoms if similar symptom prevalence), unacknowledged
416 weights to input symptoms (certain input symptoms explaining the diagnoses much better),
417 and biases that were induced by data processing and could not be explained by the input
418 symptoms. It is unclear whether the diagnosis of dysthymic disorder was intentionally
419 designed to be more prevalent than those of major depressive episodes or manic episodes
420 assuming input symptoms with the same prevalence rates.

421 In the real world, there are other important issues related to the diagnostic criteria.
422 For example, diagnoses are not closely linked to treatment,[19, 29] diagnoses are not well
423 made particularly by non-psychiatrists,[30] and there are two diagnostic systems (the DSM
424 and the International Classification of Disease) that require efforts to harmonize.[31] Amid
425 these issues, we think the diagnostic criteria for mental illnesses should be reviewed and
426 improved to be easier to understand and use without introducing biases, and closely linked
427 to clinical decisions. Certain measures and biomarkers have been proven useful to identify
428 mental illnesses.[32, 33] We are developing methods to detect symptom-based conditions
429 better and propose methods to search for neglected mental symptoms.

430 **Limitations**

431 The strength of this study is the use of simple assumptions in simulated populations
432 that enables the comparison of the diagnostic criteria of three mental illnesses. However, the
433 assumptions in the prevalence rates and between-variable correlations for the input
434 symptoms might not be realistic. Some of the assumptions are unlikely to hold in the real
435 world. However, simulations are the only option for us due to the lack of real-world data on
436 the prevalence of the input symptoms. In addition, the translation from symptoms to
437 diagnoses was assumed to be perfect based on the diagnostic criteria. The simulations in
438 this study only reflect the problems in the design of the diagnostic criteria and are not
439 designed to review the impact of how they are used in the real world.

440 **Conclusion**

441 To the best of our knowledge, there is no study on the relationships between the
442 input symptoms and diagnoses. The input symptoms were extracted from the diagnostic
443 criteria and the diagnostic criteria were transformed into mathematical functions. Without
444 mental illness data available to the public, 100,000 subjects were simulated with different
445 assumptions on the prevalence rates (0.05, 0.1, 0.3, 0.5, and 0.7) and correlations (0, 0.1,
446 0.4, 0.7, and 0.9) of the input symptoms. We found that biases were introduced into the
447 diagnoses of three mental illnesses: major depressive episodes, dysthymic disorder, and
448 manic episodes. The prevalence rates of the diagnoses were proportional to the assumed
449 prevalence rates and between-variable correlations of the input symptoms. Certain input
450 symptoms were more important than the others in explaining the diagnoses. However, the
451 input symptoms could not fully explain the diagnoses, except when the input symptoms
452 independent of each other with 0.7 symptom prevalence rates were used for the diagnosis of
453 dysthymic disorder. In conclusion, the criteria used to diagnose these three mental illnesses
454 may fail to represent the concepts they are based on, in a similar manner to three of the
455 most commonly used scoring methods to diagnose frailty.

456 **Declarations**

457 **Acknowledgments**

458 Not applicable.

459 **Ethical Statement**

460 Not applicable.

461 **Funding Statement**

462 The authors received no specific funding for this work.

463 **Consent to participate**

464 Not applicable

465 **Consent for publication**

466 Not applicable

1
2
3 467 **Data Availability**

4 468 No real-world data used - all analysis are based on simulations reproducible with the files in
5 469 the supplemental materials.
6
7

8 470 **Competing Interests**

9 471 YSC is currently employed by the Canadian Agency for Drugs and Technologies in Health.
10 472 The other authors have declared that no competing interests exist.
11
12

13 473 **Authors' Contributions**

14 474 YSC conceptualized and designed this study, managed and analyzed data and
15 475 drafted the manuscript. KFL assisted in the interpretation of the diagnostic criteria. CJW
16 476 assisted in data management and computation. HCW, HTH, LCT, YPC, YCL, and WCC
17 477 participated in the design of this study. All authors reviewed and approved the manuscript.
18
19

20 478
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

References

1. American Psychiatric Association. Diagnostic and Statistical Manual of Mental Disorders, Fourth Edition, Text Revision (DSM-IV-TR®). Washington, DC: American Psychiatric Association Publishing; 2010.
2. Center for Substance Abuse Treatment. Managing Depressive Symptoms in Substance Abuse Clients During Early Recovery. Rockville, MD: Substance Abuse and Mental Health Services Administration (US); 2008.
3. Feighner JP, Robins E, Guze SB, Woodruff RA, Jr., Winokur G, Munoz R. Diagnostic criteria for use in psychiatric research. *Arch Gen Psychiatry*. 1972;26(1):57-63. Epub 1972/01/01.
4. Kendler KS, Munoz RA, Murphy G. The development of the Feighner criteria: a historical perspective. *The American journal of psychiatry*. 2010;167(2):134-42. Epub 2009/12/17. doi: 10.1176/appi.ajp.2009.09081155.
5. Chao Y-S, Wu C-J. PP46 When Composite Measures Or Indices Fail: Data Processing Lessons. *International Journal of Technology Assessment in Health Care*. 2019;34(S1):83-. Epub 01/03. doi: 10.1017/S0266462318002088.
6. Chao Y-S, Wu H-C, Wu C-J, Chen W-C. Index or illusion: The case of frailty indices in the Health and Retirement Study. *PLOS ONE*. 2018;13(7):e0197859. doi: 10.1371/journal.pone.0197859.
7. Barnwell-Menard JL, Li Q, Cohen AA. Effects of categorization method, regression type, and variable distribution on the inflation of Type-I error rate when categorizing a confounding variable. *Stat Med*. 2015;34(6):936-49. Epub 2014/12/17. doi: 10.1002/sim.6387.
8. Cigolle CT, Ofstedal MB, Tian Z, Blaum CS. Comparing models of frailty: the Health and Retirement Study. *J Am Geriatr Soc*. 2009;57(5):830-9. Epub 2009/05/21. doi: 10.1111/j.1532-5415.2009.02225.x.
9. Brown TA, Chorpita BF, Korotitsch W, Barlow DH. Psychometric properties of the Depression Anxiety Stress Scales (DASS) in clinical samples. *Behaviour research and therapy*. 1997;35(1):79-89.
10. Lim GY, Tam WW, Lu Y, Ho CS, Zhang MW, Ho RC. Prevalence of Depression in the Community from 30 Countries between 1994 and 2014. *Scientific reports*. 2018;8(1):2861-. doi: 10.1038/s41598-018-21243-x. PubMed PMID: 29434331.
11. Smith DJ, Nicholl BI, Cullen B, Martin D, Ul-Haq Z, Evans J, et al. Prevalence and Characteristics of Probable Major Depression and Bipolar Disorder within UK Biobank: Cross-Sectional Study of 172,751 Participants. *PLOS ONE*. 2013;8(11):e75362. doi: 10.1371/journal.pone.0075362.
12. Chao Y-S, Wu C-J. PD26 Principal Component Approximation: Canadian Health Measures Survey. *International Journal of Technology Assessment in Health Care*. 2019;34(S1):138-9. Epub 01/03. doi: 10.1017/S026646231800301X.
13. Chao Y-S, Wu C-J, Wu H-C, Chen W-C. Trend analysis for national surveys: Application to all variables from the Canadian Health Measures Survey cycle 1 to 4. *PLOS ONE*. 2018;13(8):e0200127. doi: 10.1371/journal.pone.0200127.
14. Chao Y-S, Wu H-C, Wu C-J, Chen W-C. Principal Component Approximation and Interpretation in Health Survey and Biobank Data. *Frontiers in Digital Humanities*. 2018;5(11). doi: 10.3389/fdigh.2018.00011.
15. Chao YS, Wu HC, Wu CJ, Chen WC. Stages of Biological Development across Age: An Analysis of Canadian Health Measure Survey 2007-2011. *Front Public Health*. 2018;5(2296-2565 (Print)). doi: 10.3389/fpubh.2017.00355. eCollection 2017.
16. Chao Y-S, Wu C-J. PD25 Principal Component Approximation: Medical Expenditure Panel Survey. *International Journal of Technology Assessment in Health Care*. 2019;34(S1):138-. Epub 01/03. doi: 10.1017/S0266462318003008.
17. Chao Y-S, Wu C-J, Chen T-S. Risk adjustment and observation time: comparison between cross-sectional and 2-year panel data from the Medical Expenditure Panel Survey (MEPS). *Health Information Science and Systems*. 2014;2:5. doi: 10.1186/2047-2501-2-5. PubMed PMID: PMC4340859.

- 1
2
3 530 18. Chao YS, Wu HT, Scutari M, Chen TS, Wu CJ, Durand M, et al. A network perspective on
4 531 patient experiences and health status: the Medical Expenditure Panel Survey 2004 to 2011. BMC
5 532 health services research. 2017;17(1472-6963 (Electronic)). doi: 10.1186/s12913-017-2496-5.
6 533
7 534 19. Bonnin JE. Treating without diagnosis: psychoanalysis in medical settings in Argentina. 2015.
8 535 20. Leisch F, Weingessel A, Hornik K. On the generation of correlated artificial binary data. 1998.
9 536 21. Leisch F, Weingessel A, Hornik K. bindata: Generation of Artificial Binary Data, 2012. URL
10 537 [http://CRAN.R-project.org/package=](http://CRAN.R-project.org/package=bindata) bindata R package version 09-19.
11 538 22. Lumley T, Lumley MT. Package 'leaps'. Regression Subset Selection Thomas Lumley Based on
12 539 Fortran Code by Alan Miller Available online: [http://CRAN.R-project.org/package=](http://CRAN.R-project.org/package=leaps) leaps (accessed on
13 540 18 March 2018). 2013.
14 541 23. Hastie T, Tibshirani R, Friedman J. The Elements of Statistical Learning: Data Mining,
15 542 Inference, and Prediction, Second Edition: Springer New York; 2009.
16 543 24. James G, Witten D, Hastie T, Tibshirani R. An Introduction to Statistical Learning: with
17 544 Applications in R: Springer New York; 2013.
18 545 25. Chao Y-S, Wu C-J. Principal component-based weighted indices and a framework to evaluate
19 546 indices: Results from the Medical Expenditure Panel Survey 1996 to 2011. PLoS ONE.
20 547 2017;12(9):e0183997. doi: 10.1371/journal.pone.0183997. PubMed PMID: PMC5590867.
21 548 26. Chao YS, Wu HT, Wu CJ. Feasibility of Classifying Life Stages and Searching for the
22 549 Determinants: Results from the Medical Expenditure Panel Survey 1996-2011. Front Public Health.
23 550 2017;5:247(2296-2565 (Print)). doi: 10.3389/fpubh.2017.00247. eCollection 2017.
24 551 27. R Development Core Team. R: A language and environment for statistical computing.
25 552 Vienna, Austria: R Foundation for Statistical Computing; 2016.
26 553 28. RStudio Team. RStudio: Integrated Development for R. Boston, MA: RStudio, Inc.; 2016.
27 554 29. Demyttenaere K, Bonnewyn A, Bruffaerts R, De Girolamo G, Gasquet I, Kovess V, et al.
28 555 Clinical factors influencing the prescription of antidepressants and benzodiazepines: Results from
29 556 the European study of the epidemiology of mental disorders (ESEMeD). Journal of affective
30 557 disorders. 2008;110(1-2):84-93.
31 558 30. Margolis RL. Nonpsychiatrist house staff frequently misdiagnose psychiatric disorders in
32 559 general hospital inpatients. Psychosomatics. 1994;35(5):485-91.
33 560 31. First MB. Harmonisation of ICD-11 and DSM-V: opportunities and challenges. The British
34 561 Journal of Psychiatry. 2009;195(5):382-90.
35 562 32. Husain SF, Tang T-B, Yu R, Tam WW, Tran B, Quek TT, et al. Cortical haemodynamic response
36 563 measured by functional near infrared spectroscopy during a verbal fluency task in patients with
37 564 major depression and borderline personality disorder. EBioMedicine. 2020;51:102586.
38 565 33. Ho CSH, Zhang MWB, Ho R. Optical topography in psychiatry: a chip off the old block or a
39 566 new look beyond the mind-brain frontiers? Frontiers in psychiatry. 2016;7:74.
40
41
42
43
44 566
45
46 567
47
48 568
49
50
51
52
53
54
55
56
57
58
59
60

569 **Table 1. The assumptions and parameters in the simulations**

Assumptions		
1	Equal prevalence rates for the input symptoms of the same diagnosis; presence of input symptoms assigned randomly	
2	Same correlations between the input symptoms of the diagnoses of major depressive episodes and dysthymic disorder; same correlations between the input symptoms of manic episodes	
3	The input symptoms of manic episodes created independent of those of major depressive episodes and dysthymic disorder	
4	Diagnoses made accurately based on the diagnostic criteria and symptoms reported accurately by patients	
Parameters of input symptoms of the same diagnosis for each simulation		
1	Population sizes	10,000
2	Prevalence rates (uniform for all input symptoms in a simulation)	0.05, 0.1, 0.3, 0.5, and 0.7
3	Correlations (uniform between all input symptoms of the same diagnosis in a simulation)	0, 0.1, 0.4, 0.7, and 0.9
4	Number of simulations for each combination of the assumed prevalence rates and between-variable correlations of the input symptoms	100

570

571

572

573

574

575

576 **Table 2. The input symptoms, intermediate variables, and bias variables for the diagnosis of major depressive episodes.**

Classification of symptoms	Criterion variable	Domains in the major or minor criteria	Domain variables	Symptoms	Symptom variables	Equations to derive diagnosis or domain variables	Approximation by linear regression	Mechanisms related to introducing biases
Major depressive episode (variable = mde)						$mde = mde_ma1 \times mde_ma2 \times (mde_mi3 + mde_mi4 + mde_mi5 + mde_mi6 + mde_mi7 + mde_mi8 + mde_mi9 + mde_bias1) + (1 - mde_ma1 \times mde_ma2) \times (mde_ma1 \times mde_ma2) \times (mde_mi3 + mde_mi4 + mde_mi5 + mde_mi6 + mde_mi7 + mde_mi8 + mde_mi9 + mde_bias2)$	$mde = intercept + coef1 \times mde_ma1 + coef2 \times mde_ma2 + coef3 \times mde_mi3 + coef4 \times mde_mi4 + coef5 \times mde_mi5 + coef6 \times mde_mi6 + coef7 \times mde_mi7 + coef8 \times mde_mi8 + coef9 \times mde_mi9 + coef10 \times mde_bias$	1) Multiplication to create the situations when one or two symptoms in the major criteria confirmed and the bias (mde_bias) calculated by extracting the information of the diagnosis not explained by the input symptoms and two bias variables generated by censoring (mde_bias1 and mde_bias2) 2) Categorizing of the sum of the input symptoms in the minor criteria at the threshold of three or four (mde_bias1 and mde_bias2)
Major criteria, essential for diagnosis		Depressed mood or a loss of interest or pleasure in daily activities for more than two weeks.						
		Depressed mood for more than two weeks.	mde_ma1					
		Loss of interest or pleasure in daily activities for more than two weeks.	mde_ma2					
Minor criteria (at least 5 of the symptoms including the two in major criteria)	mde_mi							
		Significant unintentional weight loss or gain	mde_mi3			$mde_mi3 = mde_mi3_1 + mde_mi3_2 + mde_mi3_bias$		Censoring of the sum of multiple input variables
				Significant unintentional weight gain	mde_mi3_1			
				Significant unintentional weight loss	mde_mi3_2			
				Information of the domain not explained by the input variables	mde_mi3_bias			
		Insomnia or sleeping too much	mde_mi4			$mde_mi4 = mde_mi4_1 + mde_mi4_2 + mde_mi4_bias$		Censoring of the sum of multiple input variables
				Insomnia	mde_mi4_1			
				Sleeping too much	mde_mi4_2			
				Information of the domain not explained by the input variables	mde_mi4_bias			
		Agitation or psychomotor retardation noticed by others	mde_mi5			$mde_mi5 = mde_mi5_1 + mde_mi5_2 + mde_mi5_bias$		Censoring of the sum of multiple input variables
				Agitation	mde_mi5_1			

			Psychomotor retardation noticed by others	mde_mi5_2	
			Information of the domain not explained by the input variables	mde_mi5_bias	
	Fatigue or loss of energy	mde_mi6		$mde_mi6 = mde_mi6_1 + mde_mi6_2 + mde_mi6_bias$	Censoring of the sum of multiple input variables
			Fatigue	mde_mi6_1	
			Loss of energy	mde_mi6_2	
			Information of the domain not explained by the input variables	mde_mi6_bias	
	Feelings of worthlessness or excessive guilt	mde_mi7		$mde_mi7 = mde_mi7_1 + mde_mi7_2 + mde_mi7_bias$	Censoring of the sum of multiple input variables
			Feelings of worthlessness	mde_mi7_1	
			Feelings of excessive guilt	mde_mi7_2	
			Information of the domain not explained by the input variables	mde_mi7_bias	
	Diminished ability to think or concentrate, or indecisiveness+	mde_mi8		$mde_mi8 = mde_mi8_1 + mde_mi8_2 + mde_mi8_bias$	Censoring of the sum of multiple input variables
			Diminished ability to think or concentrate	mde_mi8_1	
			Indecisiveness	mde_mi8_2	
			Information of the domain not explained by the input variables	mde_mi8_bias	
	Recurrent thoughts of death	mde_mi9			
	Information due to categorization (choosing three domains in minor criteria)	mde_bias1			Bias introduced to categorize the sum of the number of confirmed symptoms in the minor criteria
	Information due to categorization (choosing four domains in minor criteria)	mde_bias2			Bias introduced to categorize the sum of the number of confirmed symptoms in the minor criteria
	Information of diagnosis not explained by the domains	mde_bias			Information of the diagnosis not explained by the input variables and two bias variables generated due to data categorization

577

578

580 Table 3. The input symptoms, intermediate variables, and bias variables for the diagnosis of dysthymic disorder.

Classification of symptoms	Criterion variable	Major or minor criteria (domains)	Intermediate variables	Symptoms	Symptom variables	Equations to generate diagnosis or domain variables	Approximation	Mechanisms related to introducing biases
Dysthymia (variable = dys)						$dys = dys_ma \times dys_mi$	$dys = intercept + coef1 \times dys_ma + coef2 \times dys_mi + coef3 \times dys_bias$	Multiplication to create the situations where both the major and minor criteria met (union of two binomial variables, $mde_ma \times mde_mi$) and the bias variable (dys_bias) equivalent to the residual of the diagnosis not explained by the input symptoms and the bias variables due to censoring and categorization
Major criteria, essential for diagnosis		Depressed mood most of the day for more days than not, for at least 2 years	dys_ma					
Minor criteria (at least 2 items)			dys_mi			$dys_mi = dys_mi1 + dys_mi2 + dys_mi3 + dys_mi4 + dys_mi5 + dys_mi6 + dys_mi_bias$		Categorizing of the sum of multiple input variables
		Poor appetite or overeating	dys_mi1			$dys_mi1 = dys_mi1_1 + dys_mi1_2 + dys_mi1_bias$		Censoring of the sum of multiple input variables
				Poor appetite Overeating Information of the domain not explained by the input variables	dys_mi1_1 dys_mi1_2 dys_mi1_bias			
		Insomnia or sleeping too much*	dys_mi2/mde_mi4			$dys_mi2 = mde_mi4 = mde_mi4_1 + mde_mi4_2 + mde_mi4_bias$		Censoring of the sum of multiple input variables
				Insomnia Sleeping too much Information of the domain not explained by the input variables	mde_mi4_1 mde_mi4_2 mde_mi4_bias			
		Low energy or fatigue*	dys_mi3/mde_mi6			$dys_mi3 = mde_mi6 = mde_mi6_1 + mde_mi6_2 + mde_mi6_bias$		Censoring of the sum of multiple input variables
				Fatigue Loss of energy (low energy) Information of the domain not explained by the input variables	mde_mi6_1 mde_mi6_2 mde_mi6_bias			
		Low self-esteem Poor concentration or difficulty making decisions*	dys_mi4 dys_mi5/mde_mi8			$dys_mi5 = mde_mi8 = mde_mi8_1 + mde_mi8_2 + mde_mi8_bias$		Censoring of the sum of multiple input variables
				Diminished ability to think or concentrate (Poor concentration) difficulty making decisions (indecisiveness) Information of the domain not explained by the input variables	mde_mi8_1 mde_mi8_2 mde_mi8_bias			
		Feelings of hopelessness	dys_mi6					

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

	Information of minor criteria not explained by input variables	dys_mi_bias	Bias introduced by categorizing the number of input symptoms confirmed in the minor criteria
Information of diagnosis not explained by major or minor criteria	dys_bias		Information of the diagnosis not explained by the input symptoms and the bias variables generated due to data categorization (dys_mi_bias)

581 *The same input symptoms used for the diagnosis of major depressive episodes.

582

For peer review only

584 Table 4. The input symptoms, intermediate variables, and bias variables for the diagnosis of manic episodes.

Classification of symptoms	Criterion variable	Major or minor criteria (domains)	Domain variables	Symptoms	Symptom variables	Equations	Approximation	Mechanisms related to introducing biases	
Manic episode (variable = manic)						$\text{manic} = (1 - \text{man_ma1} \times \text{man_ma2}) \times (\text{man_ma1} + \text{man_ma2}) \times \text{man_ma3} \times (\text{man_mi1} + \text{man_mi2} + \text{man_mi3} + \text{man_mi4} + \text{man_mi5} + \text{man_mi6} + \text{man_mi7} + \text{man_bias1}) + [1 - (1 - \text{man_ma1} \times \text{man_ma2})(\text{man_ma1} + \text{man_ma2})] \times \text{man_ma3} \times (\text{man_mi1} + \text{man_mi2} + \text{man_mi3} + \text{man_mi4} + \text{man_mi5} + \text{man_mi6} + \text{man_mi7} + \text{man_bias2})$	$\text{manic} = \text{intercept} + \text{coef1} \times \text{man_ma1} + \text{coef2} \times \text{man_ma2} + \text{coef3} \times \text{man_ma3} + \text{coef4} \times \text{man_mi1} + \text{coef5} \times \text{man_mi2} + \text{coef6} \times \text{man_mi3} + \text{coef7} \times \text{man_mi4} + \text{coef8} \times \text{man_mi5} + \text{coef9} \times \text{man_mi6} + \text{coef10} \times \text{man_mi7} + \text{coef11} \times \text{man_bias}$	<ol style="list-style-type: none"> 1) Multiplication to create the situations where one of the symptom in the major criteria met (union of three binomial variables, such as $\text{man_ma1} + \text{man_ma2}$ and $\text{man_ma1} \times \text{man_ma2}$), \n 2) multiplication for the condition of presenting irritable mood (... x man_ma3), and 3) and the bias variable (man_bias) equivalent to the residual of the diagnosis not explained by the input symptoms and the bias variables due to censoring; 4) the bias variables introduced by categorizing the number of input symptoms confirmed in the minor criteria (man_bias1 and man_bias2) 	
	Major criteria, essential for the diagnosis of a manic episode (more than one bipolar episode required to diagnose bipolar disorder)		A distinct period of abnormally and persistently elevated, expansive, or irritable mood, lasting at least 1 week (or any duration if hospitalization is necessary)		Elevated mood, lasting at least 1 week	man_ma1			
				Expansive mood, lasting at least 1 week	man_ma2				
				Irritable mood, lasting at least 1 week	man_ma3				
Minor criteria (3 or more of the following symptoms have persisted; 4 if the mood is only irritable)		Increased self-esteem or grandiosity	man_mi1	Increased self-esteem	man_mi1_1	$\text{man_mi1} = \text{man_mi1_1} + \text{man_mi1_2} + \text{man_mi1_bias}$		Censoring of the sum of multiple input variables	
			Grandiosity	man_mi1_2					

			Information of the domain not explained by the input variables	man_mi1_bias	
	Decreased need for sleep (e.g., feels rested after only 3 hours of sleep)	man_mi2			
	More talkative than usual or pressure to keep talking	man_mi3	More talkative than usual Pressure to keep talking	man_mi3_1 man_mi3_2	man_mi3 = man_mi3_1 + man_mi3_2 + man_mi3_bias
			Information of the domain not explained by the input variables	man_mi3_bias	Censoring of the sum of multiple input variables
	Flight of ideas or subjective experience that thoughts are racing	man_mi4	Flight of ideas Subjective experience that thoughts are racing	man_mi4_1 man_mi4_2	man_mi4 = man_mi4_1 + man_mi4_2 + man_mi4_bias
			Information of the domain not explained by the input variables	man_mi4_bias	Censoring of the sum of multiple input variables
	Distractibility (i.e., attention too easily drawn to unimportant or irrelevant external stimuli)	man_mi5			
	Increase in goal-directed activity (either socially, at work or school, or sexually) or psychomotor agitation	man_mi6	Increase in goal-directed activity Psychomotor agitation	man_mi6_1 man_mi6_2	man_mi6 = man_mi6_1 + man_mi6_2 + man_mi6_bias
			Information of the domain not explained by the input variables	man_mi6_bias	Censoring of the sum of multiple input variables
	Excessive involvement in pleasurable activities that have a high potential for painful consequences (e.g., engaging in unrestrained buying sprees, sexual indiscretions, or foolish business investments)	man_mi7			
	Information of diagnosis due to categorization (choosing at least three symptoms)	man_bias1			Bias introduced by categorizing the number of input symptoms confirmed in the minor criteria
	Information of diagnosis due to categorization (choosing at least four symptoms)	man_bias2			Bias introduced by categorizing the number of input symptoms confirmed in the minor criteria
	Information of diagnosis not	man_bias			Information of the diagnosis not explained by the input symptoms and

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

explained by symptoms

the bias variables generated due to data categorization, man_bias1 and man_bias2

For peer review only

586 **Table 5. The derived prevalence rates of the diagnoses of major depressive episodes, dysthymic**
 587 **disorder, and manic episodes based on the assumed prevalence rates and between-variable correlations**
 588 **of the input symptoms**

Assumed correlations between input symptoms	Assumed prevalence of input symptoms	Major depressive episodes	Dysthymic disorder	Manic episodes
0	0.05	0 (95% CI = 0 to 0)	0.004 (95% CI = 0.004 to 0.004)	0 (95% CI = 0 to 0)
0	0.1	0.001 (95% CI = 0.001 to 0.001)	0.025 (95% CI = 0.025 to 0.025)	0.002 (95% CI = 0.002 to 0.002)
0	0.3	0.067 (95% CI = 0.067 to 0.067)	0.249 (95% CI = 0.249 to 0.249)	0.136 (95% CI = 0.135 to 0.136)
0	0.5	0.245 (95% CI = 0.244 to 0.245)	0.493 (95% CI = 0.493 to 0.493)	0.436 (95% CI = 0.436 to 0.436)
0	0.7	0.49 (95% CI = 0.49 to 0.49)	0.7 (95% CI = 0.7 to 0.7)	0.692 (95% CI = 0.692 to 0.693)
0.1	0.05	0.004 (95% CI = 0.004 to 0.004)	0.018 (95% CI = 0.018 to 0.018)	0.007 (95% CI = 0.007 to 0.007)
0.1	0.1	0.011 (95% CI = 0.011 to 0.011)	0.049 (95% CI = 0.049 to 0.049)	0.022 (95% CI = 0.021 to 0.022)
0.1	0.3	0.094 (95% CI = 0.094 to 0.094)	0.25 (95% CI = 0.25 to 0.25)	0.172 (95% CI = 0.171 to 0.172)
0.1	0.5	0.267 (95% CI = 0.267 to 0.268)	0.482 (95% CI = 0.482 to 0.482)	0.425 (95% CI = 0.425 to 0.425)
0.1	0.7	0.51 (95% CI = 0.509 to 0.51)	0.697 (95% CI = 0.697 to 0.697)	0.679 (95% CI = 0.679 to 0.679)
0.4	0.05	0.019 (95% CI = 0.019 to 0.019)	0.037 (95% CI = 0.037 to 0.037)	0.029 (95% CI = 0.029 to 0.029)
0.4	0.1	0.042 (95% CI = 0.042 to 0.042)	0.078 (95% CI = 0.078 to 0.078)	0.062 (95% CI = 0.062 to 0.062)
0.4	0.3	0.166 (95% CI = 0.166 to 0.167)	0.267 (95% CI = 0.267 to 0.267)	0.231 (95% CI = 0.231 to 0.231)
0.4	0.5	0.344 (95% CI = 0.344 to 0.344)	0.476 (95% CI = 0.476 to 0.476)	0.44 (95% CI = 0.44 to 0.441)
0.4	0.7	0.57 (95% CI = 0.57 to 0.57)	0.689 (95% CI = 0.688 to 0.689)	0.666 (95% CI = 0.666 to 0.666)
0.7	0.05	0.035 (95% CI = 0.035 to 0.035)	0.046 (95% CI = 0.046 to 0.046)	0.042 (95% CI = 0.042 to 0.042)
0.7	0.1	0.071 (95% CI = 0.071 to 0.071)	0.092 (95% CI = 0.092 to 0.092)	0.085 (95% CI = 0.085 to 0.085)
0.7	0.3	0.233 (95% CI = 0.233 to 0.234)	0.285 (95% CI = 0.285 to 0.285)	0.27 (95% CI = 0.27 to 0.27)
0.7	0.5	0.422 (95% CI = 0.421 to 0.422)	0.486 (95% CI = 0.485 to 0.486)	0.469 (95% CI = 0.468 to 0.469)
0.7	0.7	0.635 (95% CI = 0.635 to 0.635)	0.69 (95% CI = 0.69 to 0.691)	0.678 (95% CI = 0.677 to 0.678)
0.9	0.05	0.042 (95% CI = 0.042 to 0.042)	0.048 (95% CI = 0.048 to 0.048)	0.046 (95% CI = 0.046 to 0.046)
0.9	0.1	0.085 (95% CI = 0.085 to 0.085)	0.096 (95% CI = 0.096 to 0.097)	0.093 (95% CI = 0.093 to 0.093)
0.9	0.3	0.268 (95% CI = 0.268 to 0.268)	0.293 (95% CI = 0.293 to 0.293)	0.286 (95% CI = 0.286 to 0.287)
0.9	0.5	0.463 (95% CI = 0.463 to 0.463)	0.493 (95% CI = 0.492 to 0.493)	0.485 (95% CI = 0.485 to 0.486)
0.9	0.7	0.669 (95% CI = 0.669 to 0.669)	0.695 (95% CI = 0.694 to 0.695)	0.688 (95% CI = 0.688 to 0.688)

589

590

592 **Table 6. The individual input symptoms that best explained the diagnoses: major depressive episodes,**
 593 **dysthymic disorder, and manic episodes**

Assumed correlations between input symptoms	Assumed prevalence of input symptoms	Major depressive episodes	Dysthymic disorder	Manic episodes
0	0.05	mde_ma1	dys_ma	man_ma3
0	0.05	0.001 (95% CI = 0.001 to 0.001)	0.076 (95% CI = 0.075 to 0.077)	0.002 (95% CI = 0.002 to 0.002)
0	0.1	mde_ma1	dys_ma	man_ma3
0	0.1	0.01 (95% CI = 0.01 to 0.01)	0.228 (95% CI = 0.227 to 0.229)	0.021 (95% CI = 0.02 to 0.021)
0	0.3	mde_ma1	dys_ma	man_ma3
0	0.3	0.167 (95% CI = 0.167 to 0.167)	0.774 (95% CI = 0.773 to 0.774)	0.366 (95% CI = 0.366 to 0.367)
0	0.5	mde_ma2	dys_ma	man_ma3
0	0.5	0.324 (95% CI = 0.324 to 0.325)	0.971 (95% CI = 0.971 to 0.971)	0.773 (95% CI = 0.772 to 0.773)
0	0.7	mde_ma2	dys_ma	man_ma3
0	0.7	0.412 (95% CI = 0.412 to 0.412)	0.999 (95% CI = 0.999 to 0.999)	0.964 (95% CI = 0.964 to 0.964)
0.1	0.05	mde_ma2	dys_ma	man_ma3
0.1	0.05	0.07 (95% CI = 0.07 to 0.071)	0.353 (95% CI = 0.352 to 0.355)	0.136 (95% CI = 0.135 to 0.137)
0.1	0.1	mde_ma1	dys_ma	man_ma3
0.1	0.1	0.101 (95% CI = 0.1 to 0.101)	0.462 (95% CI = 0.461 to 0.463)	0.199 (95% CI = 0.198 to 0.199)
0.1	0.3	mde_ma2	dys_ma	man_ma3
0.1	0.3	0.242 (95% CI = 0.242 to 0.243)	0.777 (95% CI = 0.777 to 0.778)	0.483 (95% CI = 0.483 to 0.484)
0.1	0.5	mde_ma2	dys_ma	man_ma3
0.1	0.5	0.365 (95% CI = 0.365 to 0.366)	0.932 (95% CI = 0.931 to 0.932)	0.74 (95% CI = 0.74 to 0.741)
0.1	0.7	mde_ma2	dys_ma	man_ma3
0.1	0.7	0.445 (95% CI = 0.445 to 0.446)	0.986 (95% CI = 0.986 to 0.986)	0.906 (95% CI = 0.906 to 0.907)
0.4	0.05	mde_ma1	dys_ma	man_ma3
0.4	0.05	0.375 (95% CI = 0.373 to 0.376)	0.731 (95% CI = 0.729 to 0.732)	0.561 (95% CI = 0.559 to 0.562)
0.4	0.1	mde_ma1	dys_ma	man_ma3
0.4	0.1	0.395 (95% CI = 0.394 to 0.396)	0.763 (95% CI = 0.762 to 0.764)	0.595 (95% CI = 0.594 to 0.596)
0.4	0.3	mde_ma1	dys_ma	man_ma3
0.4	0.3	0.465 (95% CI = 0.465 to 0.466)	0.851 (95% CI = 0.85 to 0.851)	0.701 (95% CI = 0.701 to 0.702)
0.4	0.5	mde_ma2	dys_ma	man_ma3
0.4	0.5	0.525 (95% CI = 0.524 to 0.525)	0.908 (95% CI = 0.908 to 0.908)	0.787 (95% CI = 0.786 to 0.787)
0.4	0.7	mde_ma2	dys_ma	man_ma3
0.4	0.7	0.568 (95% CI = 0.568 to 0.569)	0.946 (95% CI = 0.946 to 0.947)	0.855 (95% CI = 0.854 to 0.855)
0.7	0.05	mde_ma2	dys_ma	man_ma3
0.7	0.05	0.688 (95% CI = 0.687 to 0.69)	0.909 (95% CI = 0.908 to 0.909)	0.831 (95% CI = 0.83 to 0.832)
0.7	0.1	mde_ma1	dys_ma	man_ma3
0.7	0.1	0.688 (95% CI = 0.687 to 0.689)	0.912 (95% CI = 0.911 to 0.913)	0.836 (95% CI = 0.835 to 0.836)
0.7	0.3	mde_ma2	dys_ma	man_ma3
0.7	0.3	0.71 (95% CI = 0.709 to 0.711)	0.93 (95% CI = 0.93 to 0.93)	0.862 (95% CI = 0.861 to 0.862)
0.7	0.5	mde_ma2	dys_ma	man_ma3
0.7	0.5	0.729 (95% CI = 0.728 to 0.729)	0.944 (95% CI = 0.943 to 0.944)	0.882 (95% CI = 0.882 to 0.883)
0.7	0.7	mde_ma1	dys_ma	man_ma3
0.7	0.7	0.745 (95% CI = 0.744 to 0.745)	0.954 (95% CI = 0.954 to 0.955)	0.9 (95% CI = 0.9 to 0.9)
0.9	0.05	mde_ma1	dys_ma	man_ma3
0.9	0.05	0.828 (95% CI = 0.827 to 0.829)	0.958 (95% CI = 0.957 to 0.958)	0.918 (95% CI = 0.917 to 0.919)
0.9	0.1	mde_ma2	dys_ma	man_ma3
0.9	0.1	0.838 (95% CI = 0.838 to 0.839)	0.961 (95% CI = 0.961 to 0.961)	0.925 (95% CI = 0.924 to 0.925)
0.9	0.3	mde_ma2	dys_ma	man_ma3
0.9	0.3	0.856 (95% CI = 0.856 to 0.857)	0.969 (95% CI = 0.968 to 0.969)	0.937 (95% CI = 0.936 to 0.937)
0.9	0.5	mde_ma2	dys_ma	man_ma3
0.9	0.5	0.862 (95% CI = 0.862 to 0.863)	0.972 (95% CI = 0.972 to 0.972)	0.942 (95% CI = 0.942 to 0.943)
0.9	0.7	mde_ma2	dys_ma	man_ma3
0.9	0.7	0.865 (95% CI = 0.865 to 0.866)	0.974 (95% CI = 0.974 to 0.974)	0.946 (95% CI = 0.946 to 0.946)

594

595

596

597

598

599

600

601

602

603

604

605

606 Table 7. The individual bias variables that best explained the diagnoses: major depressive
607 episodes, dysthymic disorder, and manic episodes

Assumed correlations between input symptoms	Assumed prevalence of input symptoms	Major depressive episodes	Dysthymic disorder	Manic episodes
0	0.05	mde_bias2	dys_bias	man_bias2
0	0.05	0 (95% CI = 0 to 0)	0.028 (95% CI = 0.028 to 0.028)	0.001 (95% CI = 0.001 to 0.001)
0	0.1	mde_bias2	dys_bias	man_bias2
0	0.1	0.004 (95% CI = 0.004 to 0.004)	0.053 (95% CI = 0.053 to 0.054)	0.011 (95% CI = 0.011 to 0.011)
0	0.3	mde_bias2	dys_bias	man_bias1
0	0.3	0.015 (95% CI = 0.015 to 0.015)	0.045 (95% CI = 0.045 to 0.045)	0.089 (95% CI = 0.089 to 0.09)
0	0.5	mde_bias	dys_bias	man_bias1
0	0.5	0.013 (95% CI = 0.013 to 0.014)	0.007 (95% CI = 0.007 to 0.007)	0.035 (95% CI = 0.034 to 0.035)
0	0.7	mde_bias	dys_bias	man_bias1
0	0.7	0.01 (95% CI = 0.01 to 0.01)	0 (95% CI = 0 to 0)	0.002 (95% CI = 0.002 to 0.002)
0.1	0.05	mde_bias2	dys_bias	man_bias1
0.1	0.05	0.037 (95% CI = 0.037 to 0.037)	0.113 (95% CI = 0.113 to 0.114)	0.083 (95% CI = 0.083 to 0.084)
0.1	0.1	mde_bias2	dys_bias	man_bias1
0.1	0.1	0.047 (95% CI = 0.047 to 0.048)	0.122 (95% CI = 0.121 to 0.122)	0.116 (95% CI = 0.115 to 0.116)
0.1	0.3	mde_bias2	dys_mi_bias	man_bias1
0.1	0.3	0.077 (95% CI = 0.077 to 0.077)	0.105 (95% CI = 0.105 to 0.106)	0.198 (95% CI = 0.197 to 0.198)
0.1	0.5	mde_bias2	dys_mi_bias	man_bias1
0.1	0.5	0.079 (95% CI = 0.079 to 0.08)	0.073 (95% CI = 0.073 to 0.073)	0.166 (95% CI = 0.166 to 0.167)
0.1	0.7	mde_bias2	dys_mi_bias	man_bias1
0.1	0.7	0.065 (95% CI = 0.065 to 0.065)	0.047 (95% CI = 0.046 to 0.047)	0.094 (95% CI = 0.093 to 0.094)
0.4	0.05	mde_bias1	dys_mi_bias	man_bias1
0.4	0.05	0.294 (95% CI = 0.293 to 0.295)	0.415 (95% CI = 0.413 to 0.416)	0.432 (95% CI = 0.431 to 0.433)
0.4	0.1	mde_bias1	dys_mi_bias	man_bias1
0.4	0.1	0.304 (95% CI = 0.303 to 0.304)	0.419 (95% CI = 0.418 to 0.42)	0.445 (95% CI = 0.444 to 0.445)
0.4	0.3	mde_bias1	dys_mi_bias	man_bias1
0.4	0.3	0.335 (95% CI = 0.334 to 0.335)	0.411 (95% CI = 0.411 to 0.412)	0.473 (95% CI = 0.472 to 0.473)
0.4	0.5	mde_bias1	dys_mi_bias	man_bias1
0.4	0.5	0.354 (95% CI = 0.354 to 0.355)	0.395 (95% CI = 0.395 to 0.396)	0.475 (95% CI = 0.474 to 0.475)
0.4	0.7	mde_bias1	dys_mi_bias	man_bias1
0.4	0.7	0.356 (95% CI = 0.355 to 0.356)	0.367 (95% CI = 0.366 to 0.367)	0.451 (95% CI = 0.45 to 0.451)
0.7	0.05	mde_bias1	dys_mi_bias	man_bias1
0.7	0.05	0.616 (95% CI = 0.615 to 0.617)	0.705 (95% CI = 0.704 to 0.706)	0.723 (95% CI = 0.722 to 0.724)
0.7	0.1	mde_bias1	dys_mi_bias	man_bias1
0.7	0.1	0.611 (95% CI = 0.611 to 0.612)	0.699 (95% CI = 0.698 to 0.699)	0.72 (95% CI = 0.72 to 0.721)
0.7	0.3	mde_bias1	dys_mi_bias	man_bias1
0.7	0.3	0.623 (95% CI = 0.623 to 0.624)	0.699 (95% CI = 0.699 to 0.7)	0.728 (95% CI = 0.728 to 0.729)
0.7	0.5	mde_bias1	dys_mi_bias	man_bias1
0.7	0.5	0.632 (95% CI = 0.632 to 0.633)	0.696 (95% CI = 0.696 to 0.697)	0.731 (95% CI = 0.731 to 0.732)
0.7	0.7	mde_bias1	dys_mi_bias	man_bias1
0.7	0.7	0.639 (95% CI = 0.638 to 0.639)	0.693 (95% CI = 0.692 to 0.693)	0.732 (95% CI = 0.731 to 0.732)
0.9	0.05	mde_bias1	dys_mi_bias	man_bias1
0.9	0.05	0.777 (95% CI = 0.776 to 0.778)	0.835 (95% CI = 0.834 to 0.835)	0.847 (95% CI = 0.847 to 0.848)
0.9	0.1	mde_bias1	dys_mi_bias	man_bias1
0.9	0.1	0.788 (95% CI = 0.788 to 0.789)	0.842 (95% CI = 0.841 to 0.843)	0.855 (95% CI = 0.854 to 0.855)
0.9	0.3	mde_bias1	dys_mi_bias	man_bias1
0.9	0.3	0.807 (95% CI = 0.806 to 0.807)	0.854 (95% CI = 0.853 to 0.854)	0.867 (95% CI = 0.867 to 0.868)
0.9	0.5	mde_bias1	dys_mi_bias	man_bias1
0.9	0.5	0.811 (95% CI = 0.811 to 0.811)	0.855 (95% CI = 0.855 to 0.856)	0.87 (95% CI = 0.87 to 0.871)
0.9	0.7	mde_bias1	dys_mi_bias	man_bias1
0.9	0.7	0.812 (95% CI = 0.811 to 0.812)	0.853 (95% CI = 0.853 to 0.853)	0.869 (95% CI = 0.869 to 0.87)

608

609

610

612 Table 8. Approximating the diagnoses using input symptoms and derived adjusted R-
 613 squared

Assumed correlations between input symptoms	Assumed prevalence of input symptoms	Major depressive episodes	Dysthymic disorder	Manic episodes
0	0.05	0.003 (95% CI = 0.002 to 0.003)	0.122 (95% CI = 0.121 to 0.123)	0.004 (95% CI = 0.004 to 0.005)
0	0.1	0.024 (95% CI = 0.023 to 0.024)	0.305 (95% CI = 0.304 to 0.306)	0.039 (95% CI = 0.038 to 0.039)
0	0.3	0.348 (95% CI = 0.348 to 0.349)	0.842 (95% CI = 0.841 to 0.842)	0.483 (95% CI = 0.482 to 0.483)
0	0.5	0.649 (95% CI = 0.649 to 0.649)	0.986 (95% CI = 0.986 to 0.986)	0.817 (95% CI = 0.817 to 0.817)
0	0.7	0.823 (95% CI = 0.823 to 0.823)	1 (95% CI = 1 to 1)	0.967 (95% CI = 0.967 to 0.967)
0.1	0.05	0.143 (95% CI = 0.141 to 0.144)	0.435 (95% CI = 0.433 to 0.436)	0.212 (95% CI = 0.211 to 0.213)
0.1	0.1	0.198 (95% CI = 0.197 to 0.199)	0.539 (95% CI = 0.538 to 0.54)	0.29 (95% CI = 0.289 to 0.291)
0.1	0.3	0.45 (95% CI = 0.45 to 0.451)	0.826 (95% CI = 0.826 to 0.827)	0.588 (95% CI = 0.588 to 0.589)
0.1	0.5	0.663 (95% CI = 0.663 to 0.664)	0.952 (95% CI = 0.952 to 0.952)	0.799 (95% CI = 0.799 to 0.799)
0.1	0.7	0.809 (95% CI = 0.809 to 0.809)	0.991 (95% CI = 0.991 to 0.991)	0.922 (95% CI = 0.922 to 0.922)
0.4	0.05	0.587 (95% CI = 0.585 to 0.588)	0.782 (95% CI = 0.781 to 0.783)	0.675 (95% CI = 0.674 to 0.676)
0.4	0.1	0.607 (95% CI = 0.606 to 0.608)	0.807 (95% CI = 0.807 to 0.808)	0.698 (95% CI = 0.697 to 0.698)
0.4	0.3	0.688 (95% CI = 0.688 to 0.689)	0.878 (95% CI = 0.877 to 0.878)	0.775 (95% CI = 0.774 to 0.775)
0.4	0.5	0.761 (95% CI = 0.761 to 0.762)	0.925 (95% CI = 0.924 to 0.925)	0.838 (95% CI = 0.838 to 0.838)
0.4	0.7	0.821 (95% CI = 0.821 to 0.822)	0.956 (95% CI = 0.956 to 0.956)	0.887 (95% CI = 0.887 to 0.888)
0.7	0.05	0.813 (95% CI = 0.812 to 0.814)	0.925 (95% CI = 0.925 to 0.926)	0.877 (95% CI = 0.877 to 0.878)
0.7	0.1	0.826 (95% CI = 0.826 to 0.827)	0.928 (95% CI = 0.927 to 0.928)	0.881 (95% CI = 0.881 to 0.882)
0.7	0.3	0.86 (95% CI = 0.86 to 0.86)	0.942 (95% CI = 0.942 to 0.942)	0.9 (95% CI = 0.9 to 0.9)
0.7	0.5	0.88 (95% CI = 0.88 to 0.88)	0.953 (95% CI = 0.953 to 0.953)	0.913 (95% CI = 0.913 to 0.913)
0.7	0.7	0.895 (95% CI = 0.895 to 0.895)	0.962 (95% CI = 0.962 to 0.962)	0.925 (95% CI = 0.925 to 0.925)
0.9	0.05	0.903 (95% CI = 0.903 to 0.904)	0.965 (95% CI = 0.965 to 0.966)	0.941 (95% CI = 0.94 to 0.941)
0.9	0.1	0.91 (95% CI = 0.91 to 0.911)	0.968 (95% CI = 0.968 to 0.968)	0.945 (95% CI = 0.945 to 0.945)
0.9	0.3	0.923 (95% CI = 0.923 to 0.923)	0.974 (95% CI = 0.974 to 0.974)	0.954 (95% CI = 0.953 to 0.954)
0.9	0.5	0.928 (95% CI = 0.928 to 0.928)	0.976 (95% CI = 0.976 to 0.977)	0.958 (95% CI = 0.957 to 0.958)
0.9	0.7	0.932 (95% CI = 0.932 to 0.932)	0.978 (95% CI = 0.978 to 0.978)	0.96 (95% CI = 0.96 to 0.96)

614

615

617 Table 9. Approximating the diagnoses using bias variables and derived R-squared

Assumed correlations between input symptoms	Assumed prevalence of input symptoms	Major depressive episodes	Dysthymic disorder	Manic episodes
0	0.05	0.003 (95% CI = 0.002 to 0.003)	0.029 (95% CI = 0.029 to 0.03)	0.004 (95% CI = 0.004 to 0.004)
0	0.1	0.013 (95% CI = 0.012 to 0.013)	0.056 (95% CI = 0.056 to 0.056)	0.017 (95% CI = 0.017 to 0.017)
0	0.3	0.083 (95% CI = 0.083 to 0.083)	0.047 (95% CI = 0.047 to 0.047)	0.098 (95% CI = 0.098 to 0.099)
0	0.5	0.111 (95% CI = 0.111 to 0.112)	0.007 (95% CI = 0.007 to 0.007)	0.039 (95% CI = 0.038 to 0.039)
0	0.7	0.095 (95% CI = 0.095 to 0.095)	0 (95% CI = 0 to 0)	0.012 (95% CI = 0.012 to 0.013)
0.1	0.05	0.083 (95% CI = 0.082 to 0.084)	0.145 (95% CI = 0.144 to 0.146)	0.126 (95% CI = 0.125 to 0.127)
0.1	0.1	0.096 (95% CI = 0.095 to 0.097)	0.156 (95% CI = 0.155 to 0.156)	0.154 (95% CI = 0.153 to 0.154)
0.1	0.3	0.145 (95% CI = 0.144 to 0.145)	0.139 (95% CI = 0.138 to 0.139)	0.216 (95% CI = 0.216 to 0.216)
0.1	0.5	0.172 (95% CI = 0.172 to 0.173)	0.097 (95% CI = 0.097 to 0.097)	0.182 (95% CI = 0.181 to 0.182)
0.1	0.7	0.175 (95% CI = 0.175 to 0.175)	0.065 (95% CI = 0.064 to 0.065)	0.115 (95% CI = 0.115 to 0.116)
0.4	0.05	0.421 (95% CI = 0.419 to 0.423)	0.455 (95% CI = 0.453 to 0.456)	0.505 (95% CI = 0.504 to 0.506)
0.4	0.1	0.422 (95% CI = 0.421 to 0.423)	0.454 (95% CI = 0.453 to 0.455)	0.507 (95% CI = 0.506 to 0.508)
0.4	0.3	0.435 (95% CI = 0.434 to 0.435)	0.442 (95% CI = 0.442 to 0.443)	0.512 (95% CI = 0.512 to 0.513)
0.4	0.5	0.452 (95% CI = 0.452 to 0.453)	0.427 (95% CI = 0.427 to 0.427)	0.506 (95% CI = 0.505 to 0.506)
0.4	0.7	0.46 (95% CI = 0.459 to 0.46)	0.403 (95% CI = 0.402 to 0.403)	0.481 (95% CI = 0.481 to 0.482)
0.7	0.05	0.728 (95% CI = 0.727 to 0.729)	0.729 (95% CI = 0.728 to 0.731)	0.764 (95% CI = 0.763 to 0.765)
0.7	0.1	0.722 (95% CI = 0.721 to 0.723)	0.723 (95% CI = 0.722 to 0.724)	0.76 (95% CI = 0.759 to 0.761)
0.7	0.3	0.726 (95% CI = 0.726 to 0.727)	0.722 (95% CI = 0.722 to 0.723)	0.761 (95% CI = 0.761 to 0.762)
0.7	0.5	0.732 (95% CI = 0.731 to 0.732)	0.72 (95% CI = 0.719 to 0.72)	0.76 (95% CI = 0.76 to 0.761)
0.7	0.7	0.737 (95% CI = 0.736 to 0.737)	0.717 (95% CI = 0.716 to 0.717)	0.758 (95% CI = 0.758 to 0.759)
0.9	0.05	0.852 (95% CI = 0.851 to 0.853)	0.85 (95% CI = 0.849 to 0.851)	0.871 (95% CI = 0.871 to 0.872)
0.9	0.1	0.86 (95% CI = 0.859 to 0.861)	0.857 (95% CI = 0.856 to 0.857)	0.876 (95% CI = 0.876 to 0.877)
0.9	0.3	0.872 (95% CI = 0.871 to 0.872)	0.867 (95% CI = 0.867 to 0.868)	0.886 (95% CI = 0.886 to 0.886)
0.9	0.5	0.874 (95% CI = 0.874 to 0.875)	0.869 (95% CI = 0.868 to 0.869)	0.888 (95% CI = 0.887 to 0.888)
0.9	0.7	0.874 (95% CI = 0.874 to 0.875)	0.867 (95% CI = 0.866 to 0.867)	0.886 (95% CI = 0.886 to 0.886)

1
2
3 619 Figure 1. The prevalence rates of an intermediate variable for the diagnosis of major
4 620 depressive episodes.
5

6 621
7

8 622 Note: The intermediate variable is “*significant unintentional weight loss or gain*” and the input
9 623 symptoms are “*significant unintentional weight loss*” and “*significant unintentional weight*
10 624 *gain*.” The black line represents the situation where the prevalence rates of the input
11 625 symptoms are the same as that of the intermediate variable. Lines above the black lines
12 626 have prevalence rates larger than those of the input symptoms.
13
14

15 627
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

1
2
3 629 Figure 2. The prevalence rates of dysthymic disorder.
4
5 630
6
7 631 Note: Dysthymic disorder is diagnosed when both the major (depressed mood most of the
8 632 day for more days than not, for at least 2 years) and minor criteria (at least two of the six
9 633 items) are confirmed. The black line represents the situation where the prevalence rates of
10 634 the input symptoms are the same as those of the intermediate variable. Lines below the
11 635 black lines have prevalence rates lower than those of the input symptoms.
12
13
14 636
15
16 637
17
18 638
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

1
2
3 639 **Figure 3. The prevalence rates of major depressive episodes.**
4
5 640

6
7 641 Note: Major depressive episodes are diagnosed when both major and minor criteria are
8 642 confirmed. The black line represents the situation where the prevalence rates of the input
9 643 symptoms are the same as that of the intermediate variable. Lines below the black lines
10 644 have prevalence rates lower than those of the input symptoms.
11

12 645

13
14 646

15
16 647
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

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

648 **Figure 4. The prevalence rates of manic episodes**

649

650 Note: Manic episodes are diagnosed when the symptoms present as described in the
651 diagnostic manual. The black line represents the situation where the prevalence rates of the
652 input symptoms are the same as those of the input symptoms. Lines below the black lines
653 have prevalence rates lower than those of the input symptoms.

654

655

656

For peer review only

1
2
3 657 **Figure 5. The approximation of the diagnosis of dysthymic disorder by the input symptoms, the bias**
4 658 **variables, and both, measured by R-squared**

5
6 659

7
8 660 Note: the diagnosis of dysthymic disorder is approximated by the input symptoms, the bias
9 661 variables, and both using forward-stepwise regression. The selection of the variables was
10 662 determined by adjusted R-squared. See Table 4 for the details in the input symptoms and
11 663 the bias variables. The assumed correlations between the input symptoms are 0.4 and the
12 664 assumed prevalence rates of the input symptoms are 0.7.

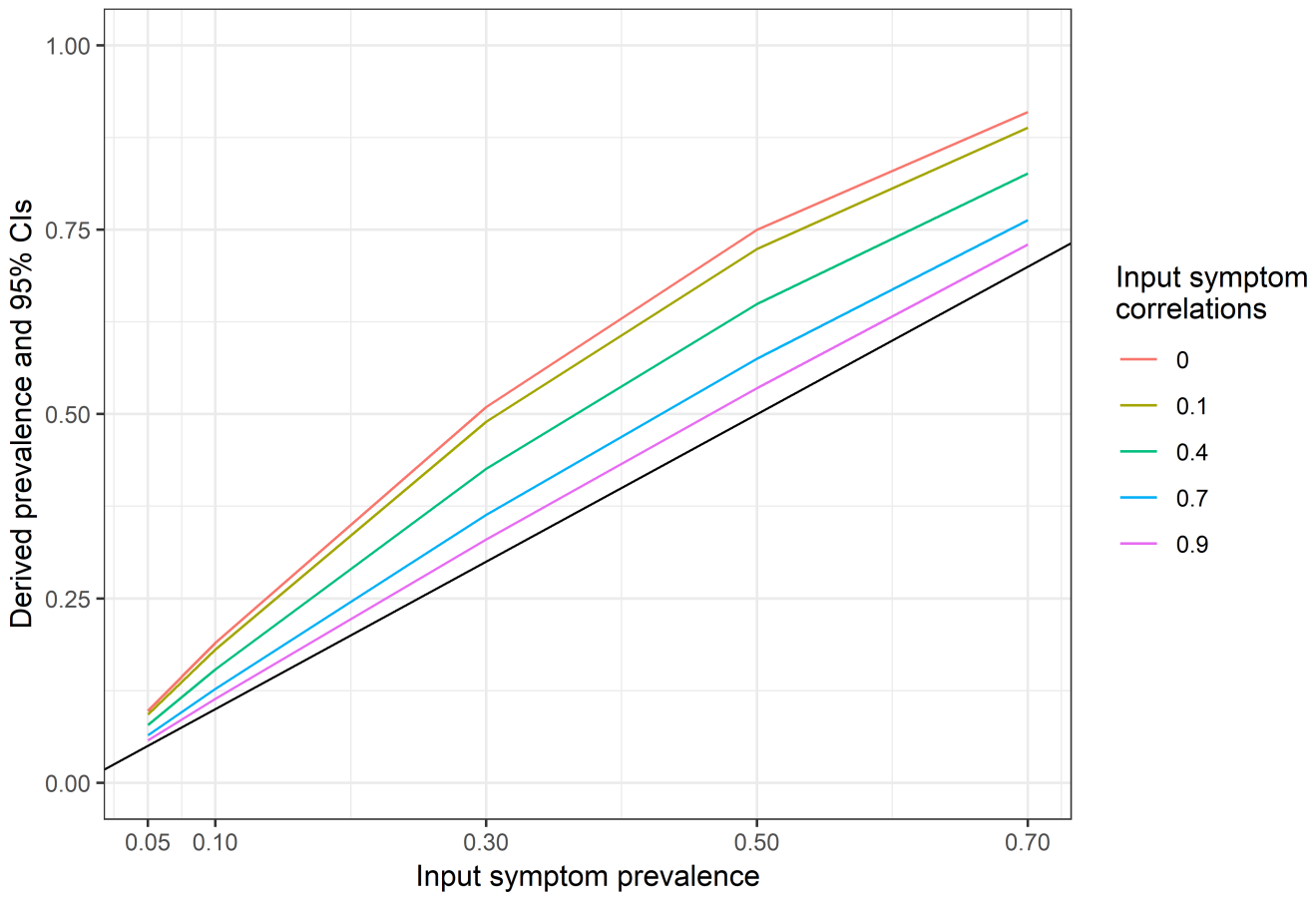
13 665

14 666

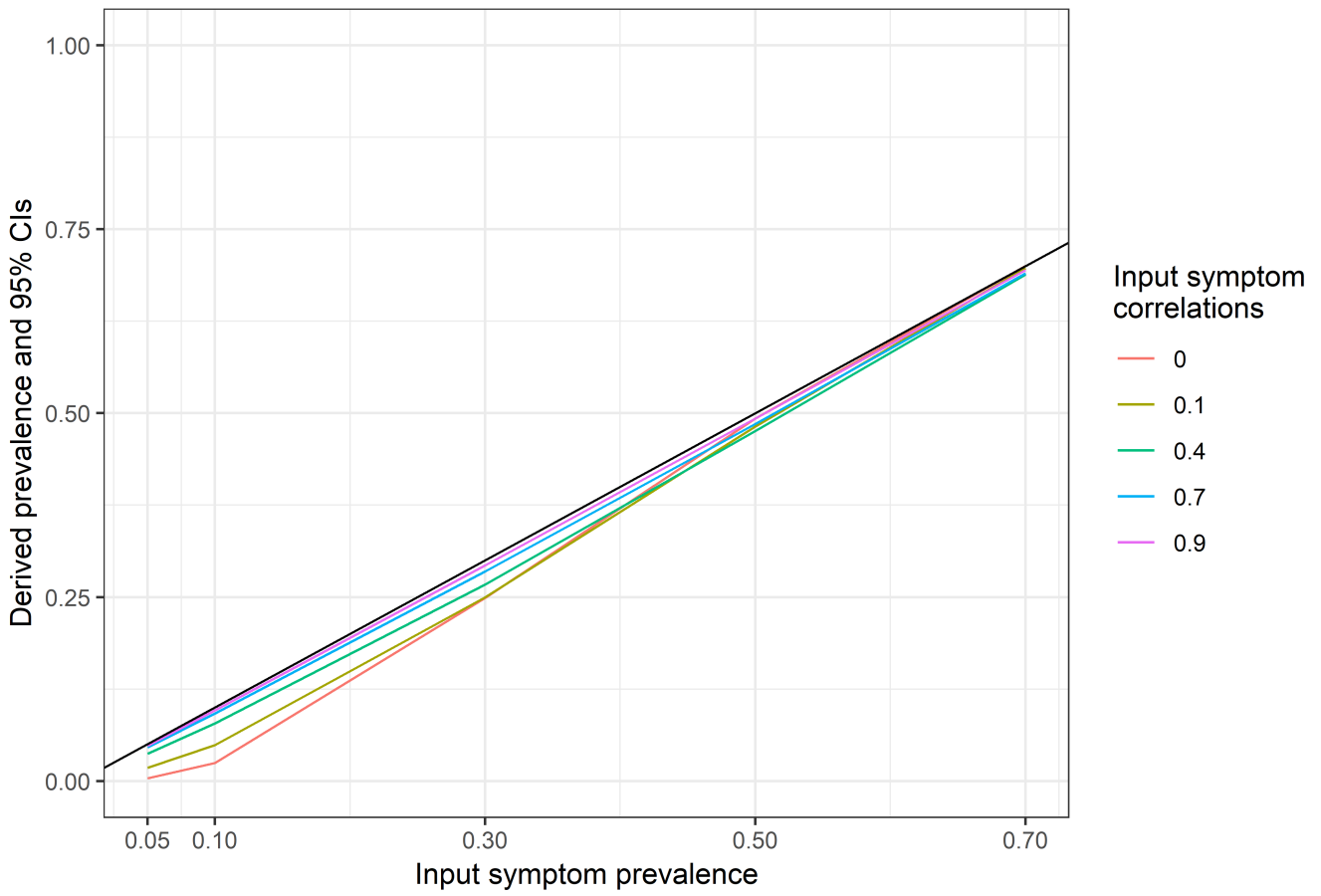
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

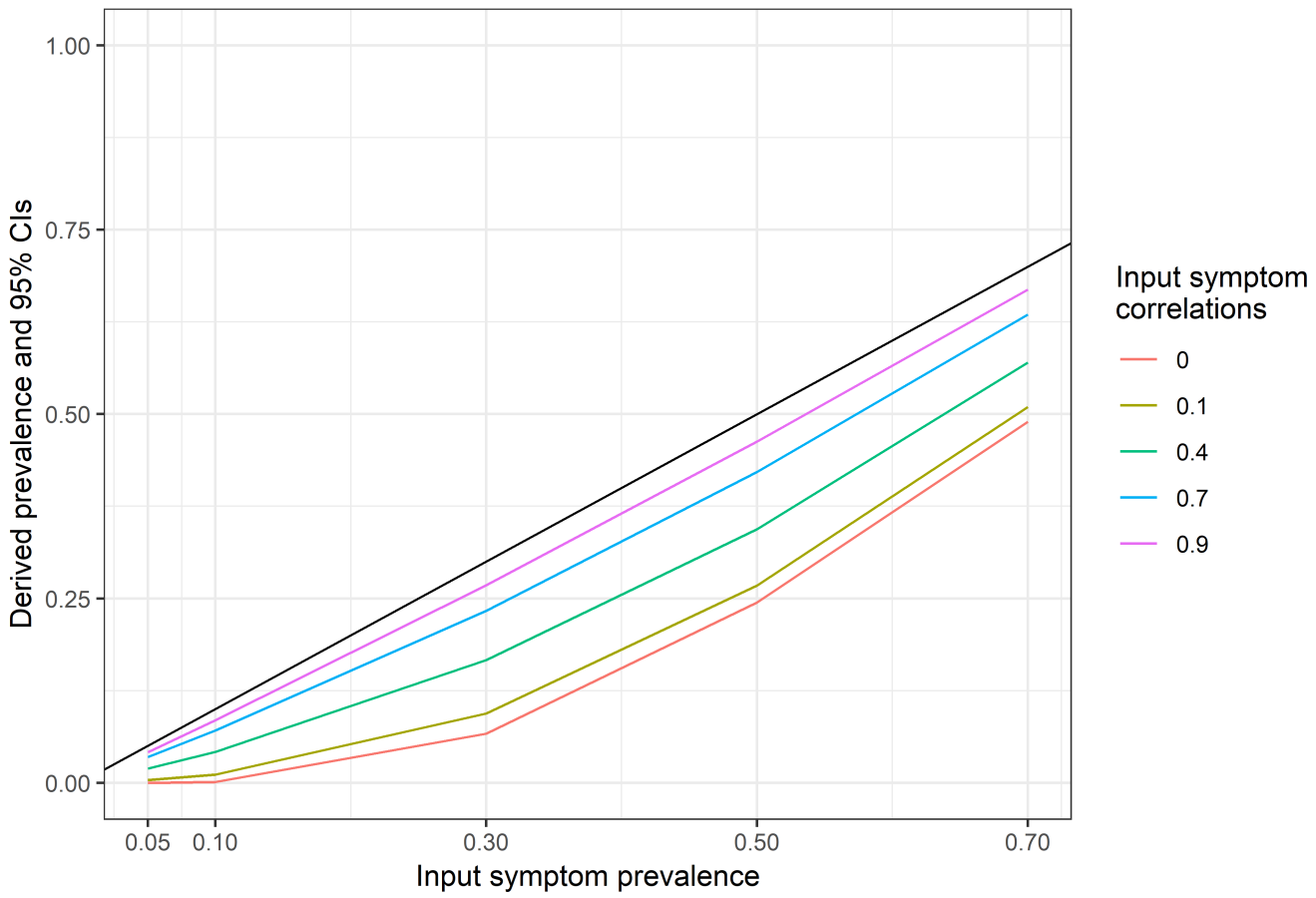
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



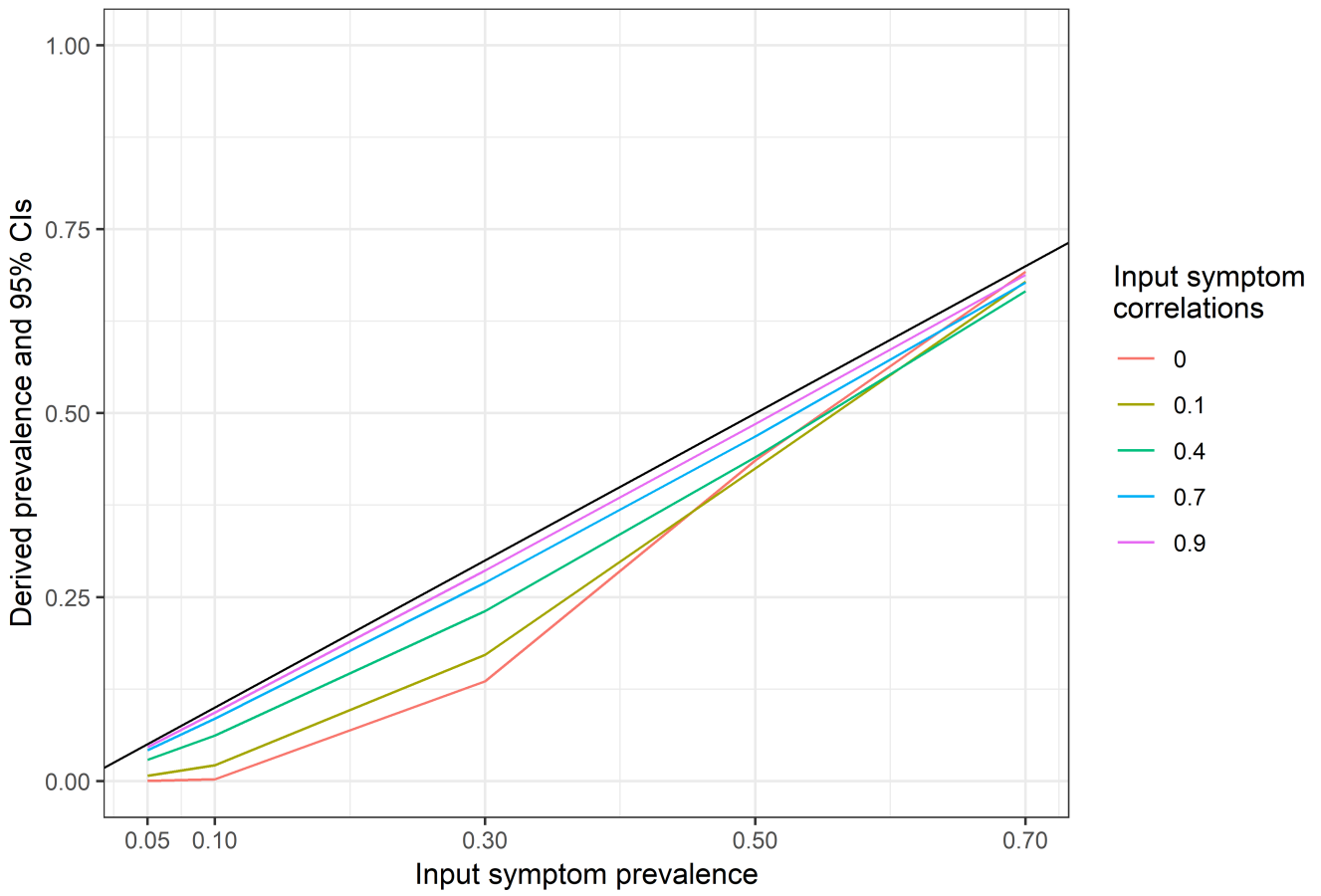
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



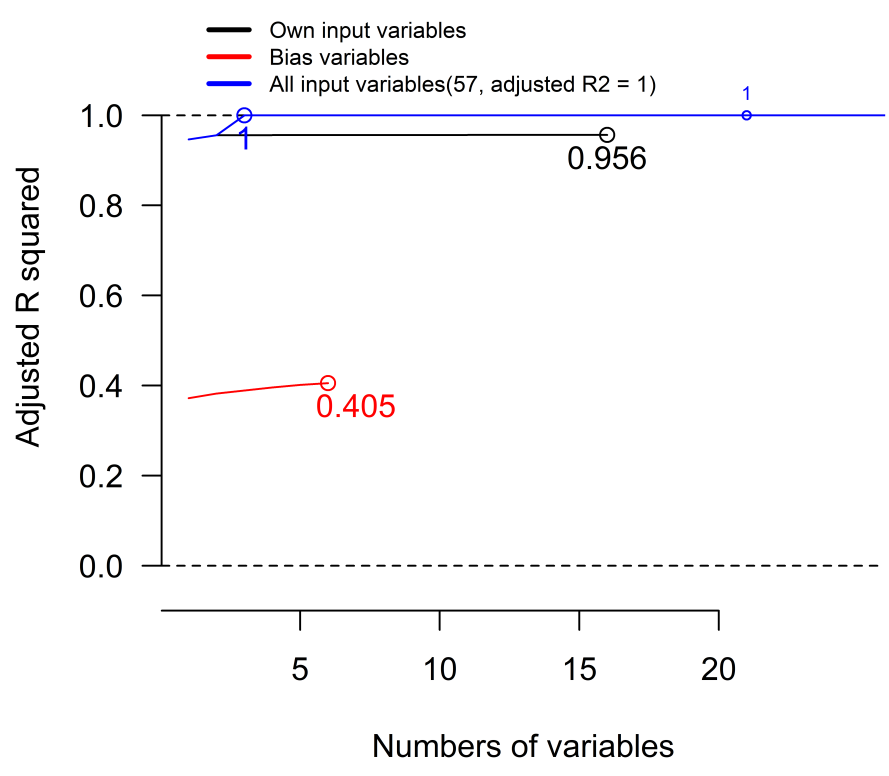
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



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



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



1
2
3
4 1 A simulation study to demonstrate the
5
6
7 2 biases in the diagnoses of mental
8
9
10 3 illnesses: major depressive episodes,
11
12 4 dysthymia, and manic episodes
13
14

15 5 **Yi-Sheng Chao^{1*}, Kuan-Fu Lin,² Chao-Jung Wu³, Hsing-Chien Wu⁴, Hui-Ting**
16 6 **Hsu⁵, Lien-Cheng Tsao⁵, Yen-Po Cheng⁵, Yi-Chun Lai⁶, Wei-Chih Chen^{7,8}**

17
18 7 *¹Independent researcher, Montréal, H2X 0A8 Canada, ²National Taiwan*
19 8 *University Hospital Yun-Lin Branch, Yunlin County, 640 Taiwan, ³Département*
20 9 *d'informatique Université du Québec à Montréal, Montréal H3B 1B4 Canada,*
21 10 *⁴Taipei Hospital Ministry of Health and Welfare New Taipei city, 242 Taiwan,*
22 11 *⁵Changhua Christian Hospital, Changhua County 526, Taiwan, ⁶National Yang-*
23 12 *Ming University Hospital, Yilan 260 Taiwan, ⁷Department of Chest Medicine,*
24 13 *Taipei Veterans General Hospital, Taipei 112, Taiwan, ⁸Institute of Emergency*
25 14 *and Critical Care Medicine, National Yang-Ming University, Taipei 112, Taiwan*
26 15 **chaoyisheng@post.harvard.edu*

27
28 16 **Keywords:** Frailty; bias; forward-stepwise regression; the Health and
29 17 Retirement Study; index mining
30 18
31
32 19


```

1
2
3
4 title: "2019_09_06 simulated mental illnesses"
5 author: "Yi-Sheng Chao"
6 date: "November 22, 2018"
7 output: pdf_document
8 editor_options:
9   chunk_output_type: inline
10
11
12
13 ##Adding correlations to the random variables
14
15 ```{r}
16 library(bindata)
17
18 library(openxlsx)
19 resu = read.xlsx("A simulation study to demonstrate the biases in three
20 diagnoses of mental illnesses.xlsx", sheet = "Prob 1")
21 names(resu)
22 unique(resu$variable)
23 memory.limit(size = 10^13)
24 ssize = 10^5
25 times = 10^2
26
27 prevalence = c(0.05, 0.1, 0.3, 0.5, 0.7)
28 rho = c(0, 0.1, 0.4, 0.7, 0.9)#correlation coefficients of the input
29 symptoms
30
31 collect = c("mean", "max",
32 "min", "derivedprevalence", "coef", "coefse", "p", "intercept",
33 "interceptp", "r2", "subcoef", "subcoefse", "subp", "subintercept",
34 "subinterceptp", "subr2", "appbyownr2", "appbybiasr2", "appbyallr2",
35 "appbyownvar", "appbybiasvar", "appbyallvar", "appbyownn", "appbybiasn",
36 "appbyalln")
37
38
39 set.seed(1)
40
41
42 ##Create a simulated data set to extract variables
43 for(preval in 1:length(prevalence)){
44   for(rh in 1:length(rho)){
45
46
47     library(openxlsx)
48     resu = read.xlsx("A simulation study to demonstrate the biases in three
49 diagnoses of mental illnesses.xlsx", sheet = "Prob 1")
50
51     # foreach(c = 1:times) %dopar% {
52     for(c in 1:times){
53
54       library(bindata)
55       bindata = as.data.frame(rmvbin(ssize, rep(prevalence[preval], 40),
56 bincorr=(1 - rho[rh])*diag(40) + rho[rh]))
57       bindata2 = as.data.frame(rmvbin(ssize, rep(prevalence[preval], 20),
58 bincorr=(1 - rho[rh])*diag(20) + rho[rh]))
59
60       ##demographic characteristics

```

```
1
2
3   sim = data.frame(1:ssize)
4   names(sim) = "id"
5   sim$female = rbinom(n = ssize, size = 1, prob = 0.51)
6   sim$age = sample(30:60, ssize, replace = TRUE)
7   sim$edu = rnorm(ssize, mean = 12, sd = 5)
8   sim$edu[which(sim$edu <= 0)] = 0
9   sim$id = NULL
10
11
12   sim$mde_ma1 = bindata[,1]
13   sim$mde_ma2 = bindata[,2]
14
15   sim$mde_mi3_1 = bindata[,3]
16   sim$mde_mi3_2 = bindata[,4]
17   sim$mde_mi3 = 1*((sim$mde_mi3_1 + sim$mde_mi3_2) > 0)
18   sim$mde_mi3_bias = sim$mde_mi3 - sim$mde_mi3_1 - sim$mde_mi3_2
19
20   sim$mde_mi4_1 = bindata[,5]
21   sim$mde_mi4_2 = bindata[,6]
22   sim$mde_mi4 = 1*((sim$mde_mi4_1 + sim$mde_mi4_2) > 0)
23   sim$mde_mi4_bias = sim$mde_mi4 - sim$mde_mi4_1 - sim$mde_mi4_2
24
25   sim$mde_mi5_1 = bindata[,7]
26   sim$mde_mi5_2 = bindata[,8]
27   sim$mde_mi5 = 1*((sim$mde_mi5_1 + sim$mde_mi5_2) > 0)
28   sim$mde_mi5_bias = sim$mde_mi5 - sim$mde_mi5_1 - sim$mde_mi5_2
29
30   sim$mde_mi6_1 = bindata[,9]
31   sim$mde_mi6_2 = bindata[,10]
32   sim$mde_mi6 = 1*((sim$mde_mi6_1 + sim$mde_mi6_2) > 0)
33   sim$mde_mi6_bias = sim$mde_mi6 - sim$mde_mi6_1 - sim$mde_mi6_2
34
35   sim$mde_mi7_1 = bindata[,11]
36   sim$mde_mi7_2 = bindata[,12]
37   sim$mde_mi7 = 1*((sim$mde_mi7_1 + sim$mde_mi7_2) > 0)
38   sim$mde_mi7_bias = sim$mde_mi7 - sim$mde_mi7_1 - sim$mde_mi7_2
39
40   sim$mde_mi8_1 = bindata[,13]
41   sim$mde_mi8_2 = bindata[,14]
42   sim$mde_mi8 = 1*((sim$mde_mi8_1 + sim$mde_mi8_2) > 0)
43   sim$mde_mi8_bias = sim$mde_mi8 - sim$mde_mi8_1 - sim$mde_mi8_2
44
45   sim$mde_mi9 = bindata[,15]
46
47   sim$mde_bias1 = 1 * ((sim$mde_mi3 + sim$mde_mi4 + sim$mde_mi5 +
48   sim$mde_mi6 + sim$mde_mi7 + sim$mde_mi8 + sim$mde_mi9)>2) - (sim$mde_mi3
49   + sim$mde_mi4 + sim$mde_mi5 + sim$mde_mi6 + sim$mde_mi7 + sim$mde_mi8 +
50   sim$mde_mi9)
51   sim$mde_bias2 = 1 * ((sim$mde_mi3 + sim$mde_mi4 + sim$mde_mi5 +
52   sim$mde_mi6 + sim$mde_mi7 + sim$mde_mi8 + sim$mde_mi9)>3) - (sim$mde_mi3
53   + sim$mde_mi4 + sim$mde_mi5 + sim$mde_mi6 + sim$mde_mi7 + sim$mde_mi8 +
54   sim$mde_mi9)
55
56   sim$mde = sim$mde_ma1 * sim$mde_ma2 * (sim$mde_mi3 + sim$mde_mi4 +
57   sim$mde_mi5 + sim$mde_mi6 + sim$mde_mi7 + sim$mde_mi8 + sim$mde_mi9 +
58   sim$mde_bias1) + (1- sim$mde_ma1 * sim$mde_ma2) * (sim$mde_ma1 *
59   sim$mde_ma2) * (sim$mde_mi3 + sim$mde_mi4 + sim$mde_mi5 + sim$mde_mi6 +
60   sim$mde_mi7 + sim$mde_mi8 + sim$mde_mi9 + sim$mde_bias2)
```

```

1
2
3
4   sim$mde_bias = sim$mde - (sim$mde_ma1 + sim$mde_ma2) - (sim$mde_mi3 +
5   sim$mde_mi4 + sim$mde_mi5 + sim$mde_mi6 + sim$mde_mi7 + sim$mde_mi8 +
6   sim$mde_mi9 + sim$mde_bias1) - (sim$mde_bias2)
7
8   ##Definition Below: even the bias and own input variables could not fully
9   explain the diagnosis
10  # sim$mde_bias = sim$mde - (sim$mde_mi3 + sim$mde_mi4 + sim$mde_mi5 +
11  sim$mde_mi6 + sim$mde_mi7 + sim$mde_mi8 + sim$mde_mi9 + sim$mde_bias1) -
12  (sim$mde_mi3 + sim$mde_mi4 + sim$mde_mi5 + sim$mde_mi6 + sim$mde_mi7 +
13  sim$mde_mi8 + sim$mde_mi9 + sim$mde_bias2)
14
15  # sim$mde_bias = sim$mde - (sim$mde_ma1 + sim$mde_ma2 + sim$mde_mi3 +
16  sim$mde_mi4 + sim$mde_mi5 + sim$mde_mi6 + sim$mde_mi7 + sim$mde_mi8 +
17  sim$mde_mi9)
18
19
20  # sim$mde_bias = resid(lm(sim$mde ~ sim$mde_ma1 + sim$mde_ma2 +
21  sim$mde_mi3 + sim$mde_mi4 + sim$mde_mi5 + sim$mde_mi6 + sim$mde_mi7 +
22  sim$mde_mi8 + sim$mde_mi9, data=sim))
23
24  ##DYS
25
26  sim$dys_ma = bindata[,16]
27
28  sim$dys_mi1_1 = bindata[,17]
29  sim$dys_mi1_2 = bindata[,18]
30  sim$dys_mi1 = 1*((sim$dys_mi1_1 + sim$dys_mi1_2) > 0)
31  sim$dys_mi1_bias = sim$dys_mi1 - sim$dys_mi1_1 - sim$dys_mi1_2
32
33  sim$dys_mi4 = bindata[,19]
34
35  sim$dys_mi6 = bindata[,20]
36
37  sim$dys_mi = 1*((sim$dys_mi1 + sim$mde_mi4 + sim$mde_mi6 + sim$dys_mi4 +
38  sim$mde_mi8 + sim$dys_mi6)>1)
39
40  sim$dys_mi_bias = sim$dys_mi - (sim$dys_mi1 + sim$mde_mi4 + sim$mde_mi6 +
41  sim$dys_mi4 + sim$mde_mi8 + sim$dys_mi6)
42
43  sim$dys = sim$dys_ma * sim$dys_mi
44
45
46  sim$dys_bias = sim$dys - (sim$dys_ma + sim$dys_mi)
47
48  # sim$dys_bias = resid(lm(sim$dys ~ sim$dys_ma + sim$dys_mi, data=sim))
49
50
51  ##Manic
52  sim$man_ma1 = bindata2[,1]
53  sim$man_ma2 = bindata2[,2]
54  sim$man_ma3 = bindata2[,3]
55
56  sim$man_mi1_1 = bindata2[,4]
57  sim$man_mi1_2 = bindata2[,5]
58  sim$man_mi1 = 1*((sim$man_mi1_1 + sim$man_mi1_2) > 0)
59  sim$man_mi1_bias = sim$man_mi1 - (sim$man_mi1_1 + sim$man_mi1_2)
60  sim$man_mi2 = bindata2[,6]

```

```

1
2
3   sim$man_mi3_1 = bindata2[,7]
4   sim$man_mi3_2 = bindata2[,8]
5   sim$man_mi3 = 1*((sim$man_mi3_1 + sim$man_mi3_2) > 0)
6   sim$man_mi3_bias = sim$man_mi3 - (sim$man_mi3_1 + sim$man_mi3_2)
7   sim$man_mi4_1 = bindata2[,9]
8   sim$man_mi4_2 = bindata2[,10]
9   sim$man_mi4 = 1*((sim$man_mi4_1 + sim$man_mi4_2) > 0)
10  sim$man_mi4_bias = sim$man_mi4 - (sim$man_mi4_1 + sim$man_mi4_2)
11  sim$man_mi5 = bindata2[,11]
12  sim$man_mi6_1 = bindata2[,12]
13  sim$man_mi6_2 = bindata2[,13]
14  sim$man_mi6 = 1*((sim$man_mi6_1 + sim$man_mi6_2) > 0)
15  sim$man_mi6_bias = sim$man_mi6 - (sim$man_mi6_1 + sim$man_mi6_2)
16  sim$man_mi7 = bindata2[,14]
17  sim$man_bias1 = 1*((sim$man_mi1 + sim$man_mi2 + sim$man_mi3 + sim$man_mi4
18  + sim$man_mi5 + sim$man_mi6 + sim$man_mi7) > 2) - (sim$man_mi1 +
19  sim$man_mi2 + sim$man_mi3 + sim$man_mi4 + sim$man_mi5 + sim$man_mi6 +
20  sim$man_mi7)
21  sim$man_bias2 = 1*((sim$man_mi1 + sim$man_mi2 + sim$man_mi3 + sim$man_mi4
22  + sim$man_mi5 + sim$man_mi6 + sim$man_mi7) > 3) - (sim$man_mi1 +
23  sim$man_mi2 + sim$man_mi3 + sim$man_mi4 + sim$man_mi5 + sim$man_mi6 +
24  sim$man_mi7)
25
26
27  sim$manic = (1- sim$man_ma1 * sim$man_ma2) * (sim$man_ma1 + sim$man_ma2)
28  * sim$man_ma3 * (sim$man_mi1 + sim$man_mi2 + sim$man_mi3 + sim$man_mi4 +
29  sim$man_mi5 + sim$man_mi6 + sim$man_mi7 + sim$man_bias1) + (1 - (1 -
30  sim$man_ma1 * sim$man_ma2) * (sim$man_ma1 + sim$man_ma2)) * sim$man_ma3 *
31  (sim$man_mi1 + sim$man_mi2 + sim$man_mi3 + sim$man_mi4 + sim$man_mi5 +
32  sim$man_mi6 + sim$man_mi7 + sim$man_bias2)
33
34
35  sim$man_bias = sim$manic - (sim$man_ma1 + sim$man_ma2 + sim$man_ma3) -
36  (sim$man_mi1 + sim$man_mi2 + sim$man_mi3 + sim$man_mi4 + sim$man_mi5 +
37  sim$man_mi6 + sim$man_mi7 + sim$man_bias1) - (sim$man_bias2)
38
39  ##end of generate data
40
41
42  resu[, paste(collect, "_", c, sep = "")] = NA
43  for(r in 1:nrow(resu)){
44    #variable characteristics
45    if(is.na(resu$variable[r]) == FALSE){
46      resu[r, paste0("derivedprevalence_", c, collapse = "")] =
47  nrow(sim[which(sim[, resu$variable[r]] == 1),])/ssize
48      resu[r, paste0("mean_", c, collapse = "")] =
49  mean(sim[,resu$variable[r]])
50      resu[r, paste0("max_", c, collapse = "")] =
51  max(sim[,resu$variable[r]])
52      resu[r, paste0("min_", c, collapse = "")] =
53  min(sim[,resu$variable[r]])
54    }
55    ##regression for the diagnosis
56    if(is.na(resu$variable[r]) == FALSE & resu$variable[r] !=
57  resu$outcome[r]){
58      eval(parse(text = paste0("templm = summary(lm(", resu$outcome[r],
59  " ~ ", resu$variable[r], ", data = sim))", collapse = "")))
60

```

```

1
2
3     resu[r, paste0("coef_", c, collapse = "")] =
4 templm$coefficients[resu$variable[r], "Estimate"]
5     resu[r, paste0("coefse_", c, collapse = "")] =
6 templm$coefficients[resu$variable[r], "Std. Error"]
7     resu[r, paste0("p_", c, collapse = "")] =
8 templm$coefficients[resu$variable[r], "Pr(>|t|)"]
9     resu[r, paste0("intercept_", c, collapse = "")] =
10 templm$coefficients["(Intercept)", "Estimate"]
11     resu[r, paste0("interceptp_", c, collapse = "")] =
12 templm$coefficients["(Intercept)", "Pr(>|t|)"]
13     resu[r, paste0("r2_", c, collapse = "")] = templm$r.squared
14 }
15     ##regression for the suboutcome/domain variables
16     if(is.na(resu$variable[r]) == FALSE & is.na(resu$suboutcome[r]) ==
17 FALSE & resu$variable[r] != resu$outcome[r] & resu$variable[r] !=
18 resu$suboutcome[r]){
19         eval(parse(text = paste0("templm = summary(lm(",
20 resu$suboutcome[r], "~ ", resu$variable[r], ", data = sim))", collapse =
21 "")))
22         resu[r, paste0("subcoef_", c, collapse = "")] =
23 templm$coefficients[resu$variable[r], "Estimate"]
24         resu[r, paste0("subcoefse_", c, collapse = "")] =
25 templm$coefficients[resu$variable[r], "Std. Error"]
26         resu[r, paste0("subp_", c, collapse = "")] =
27 templm$coefficients[resu$variable[r], "Pr(>|t|)"]
28         resu[r, paste0("subintercept_", c, collapse = "")] =
29 templm$coefficients["(Intercept)", "Estimate"]
30         resu[r, paste0("subinterceptp_", c, collapse = "")] =
31 templm$coefficients["(Intercept)", "Pr(>|t|)"]
32         resu[r, paste0("subr2_", c, collapse = "")] = templm$r.squared
33     }
34
35
36     if(r %in% as.character(1:100*50)){print(c("r:", r))}
37 }#r = rows of the variable list
38
39
40     ##Approximation by own, bias or all variables
41
42
43     #plotting area_start: only the last simulation data set used for
44 plotting
45     library(leaps)
46     #MDE
47     #own variables only
48     mdeown = NA
49     library(car)
50     sim.new = sim[,c("mde",
51 names(summary((lm(as.formula(paste0("mde ~ ", paste0(names(sim)
52 [grepl("mde_", names(sim)) == TRUE & grepl("bias", names(sim)) == FALSE],
53 collapse = " + ")), collapse = "")), data = sim))$aliases)
54 [summary((lm(as.formula(paste0("mde ~ ", paste0(names(sim)[grepl("mde_",
55 names(sim)) == TRUE & grepl("bias", names(sim)) == FALSE], collapse = " +
56 ")), collapse = "")), data = sim))$aliases == FALSE &
57 names(summary((lm(as.formula(paste0("mde ~ ", paste0(names(sim)
58 [grepl("mde_", names(sim)) == TRUE & grepl("bias", names(sim)) == FALSE],
59 collapse = " + ")), collapse = "")), data = sim))$aliases) !=
60 "(Intercept)"])]

```

```

1
2
3
4
5     for(repe in 1:40){
6         tempvif = vif(lm(mde~., data = sim.new))
7         if(any(tempvif > 10)){
8             sim.new = sim.new[,which(names(sim.new) != names(tempvif)
9 [which(tempvif == tempvif[order(-tempvif)][1]])]]
10        }
11    }
12
13
14    try(
15        (mdeown = regsubsets(mde~., data = sim.new, really.big=T,
16 method = "forward", nvmax = ncol(sim.new))), silent = F
17    )
18
19    mdeownsummary = NA
20    if(any(is.na(mdeown)) == FALSE){
21        mdeownsummary = summary(mdeown)
22    }
23
24    mdeownsummary$adjr2
25
26
27
28    ##own and bias variables
29    mdebias = NA
30
31    mdebias = regsubsets(as.formula(paste0("mde ~ ",
32 paste0(names(sim)[grepl("mde_", names(sim)) == TRUE & grepl("bias",
33 names(sim)) == TRUE], collapse = " + ")), collapse = "")), data = sim,
34 nvmax = 100, really.big=T, method = "forward")
35    mdebiassummary = summary(mdebias)
36    mdebiassummary$adjr2
37
38
39    ##all variables
40    ###in case of collinearity
41    mdeall = NA
42
43
44    ##Deal with collinearity
45    library(car)
46    sim.new = sim[,c("mde", names(summary((lm(mde~., data = sim)))
47 $aliased)[summary((lm(mde~., data = sim)))$aliased == FALSE &
48 names(summary((lm(mde~., data = sim)))$aliased) != "(Intercept)"))]
49
50
51    for(repe in 1:40){
52        tempvif = vif(lm(mde~., data = sim.new))
53        if(any(tempvif > 10)){
54            sim.new = sim.new[,which(names(sim.new) != names(tempvif)
55 [which(tempvif == tempvif[order(-tempvif)][1]])]]
56        }
57    }
58
59    ##Somehow there are problems in executing regsubsets even after
60 removing collinear variables

```

```

1
2
3
4     try(
5         (mdeall = regsubsets(mde~., data = sim.new, really.big=T,
6 method = "forward", nvmax = ncol(sim.new))), silent = F
7     )
8
9     mdeallsummary = NA
10    if(any(is.na(mdeall)) == FALSE){
11        mdeallsummary = summary(mdeall)
12    }
13
14    # mdeallsummary$adjr2
15
16
17    #DYS
18    #dys
19    #own variables only
20    dysown = NA
21    library(car)
22    sim.new = sim[,c("dys",
23 names(summary((lm(as.formula(paste0("dys ~ ", paste0(c(names(sim)
24 [(grepl("dys_", names(sim)) == TRUE | grepl("mde_mi4", names(sim)) ==
25 TRUE | grepl("mde_mi6", names(sim)) == TRUE | grepl("mde_mi8",
26 names(sim)) == TRUE) & grepl("bias", names(sim)) == FALSE]), collapse = "
27 + "), collapse = "")), data = sim))$aliased)
28 [summary((lm(as.formula(paste0("dys ~ ", paste0(c(names(sim)
29 [(grepl("dys_", names(sim)) == TRUE | grepl("mde_mi4", names(sim)) ==
30 TRUE | grepl("mde_mi6", names(sim)) == TRUE | grepl("mde_mi8",
31 names(sim)) == TRUE) & grepl("bias", names(sim)) == FALSE]), collapse = "
32 + "), collapse = "")), data = sim))$aliased == FALSE &
33 names(summary((lm(as.formula(paste0("dys ~ ", paste0(c(names(sim)
34 [(grepl("dys_", names(sim)) == TRUE | grepl("mde_mi4", names(sim)) ==
35 TRUE | grepl("mde_mi6", names(sim)) == TRUE | grepl("mde_mi8",
36 names(sim)) == TRUE) & grepl("bias", names(sim)) == FALSE]), collapse = "
37 + "), collapse = "")), data = sim))$aliased) != "(Intercept)"]]]
38
39    for(repe in 1:40){
40        tempvif = vif(lm(dys~., data = sim.new))
41        if(any(tempvif > 10)){
42            sim.new = sim.new[,which(names(sim.new) != names(tempvif)
43 [which(tempvif == tempvif[order(-tempvif)][1]])]]
44        }
45    }
46    ##Somehow there are problems in executing regsubsets even after
47 removing collinear variables
48    try(
49        (dysown = regsubsets(dys~., data = sim.new, really.big=T,
50 method = "forward", nvmax = 100)), silent = T
51    )
52
53    if(any(is.na(dysown)) == FALSE){
54        dysownsummary = summary(dysown)
55    }
56
57    ##own and bias variables
58    dysbias = NA
59
60

```

```

1
2
3     dysbias = regsubsets(as.formula(paste0("dys ~ ",
4 paste0(names(sim)[(grepl("dys_", names(sim)) == TRUE | grepl("mde_mi4",
5 names(sim)) == TRUE | grepl("mde_mi6", names(sim)) == TRUE |
6 grepl("mde_mi8", names(sim)) == TRUE) & grepl("bias", names(sim)) ==
7 TRUE], collapse = " + "), collapse = "")), data = sim, nvmax = 100,
8 really.big=T, method = "forward")
9     dysbiassummary = summary(dysbias)
10    dysbiassummary$adjr2
11
12
13    ##all variables
14    ###in case of collinearity
15    dysall = NA
16    library(car)
17    sim.new = sim[,c("dys", names(summary((lm(dys~., data = sim)))
18 $aliased)[summary((lm(dys~., data = sim)))$aliased == FALSE &
19 names(summary((lm(dys~., data = sim)))$aliased) != "(Intercept)"])]
20
21    for(repe in 1:40){
22        tempvif = vif(lm(dys~., data = sim.new))
23        if(any(tempvif > 10)){
24            sim.new = sim.new[,which(names(sim.new) != names(tempvif)
25 [which(tempvif == tempvif[order(-tempvif)][1]])]]
26        }
27    }
28    ##Somehow there are problems in executing regsubsets even after
29 removing collinear variables
30    try(
31        (dysall = regsubsets(dys~., data = sim.new, really.big=T,
32 method = "forward", nvmax = 100)), silent = T
33    )
34
35    if(any(is.na(dysall)) == FALSE){
36        dysallsummary = summary(dysall)
37    }
38    # dysallsummary$adjr2
39
40
41
42
43    #manic
44    #own variables only
45    manown = NA
46    library(car)
47    sim.new = sim[,c("manic",
48 names(summary((lm(as.formula(paste0("manic ~ ", paste0(names(sim)
49 [grepl("man_", names(sim)) == TRUE & grepl("bias", names(sim)) == FALSE],
50 collapse = " + "), collapse = "")), data = sim)))$aliased)
51 [summary((lm(as.formula(paste0("manic ~ ", paste0(names(sim)
52 [grepl("man_", names(sim)) == TRUE & grepl("bias", names(sim)) == FALSE],
53 collapse = " + "), collapse = "")), data = sim)))$aliased == FALSE &
54 names(summary((lm(as.formula(paste0("manic ~ ", paste0(names(sim)
55 [grepl("man_", names(sim)) == TRUE & grepl("bias", names(sim)) == FALSE],
56 collapse = " + "), collapse = "")), data = sim)))$aliased) !=
57 "(Intercept)"]]
58    for(repe in 1:40){
59        tempvif = vif(lm(manic~., data = sim.new))
60        if(any(tempvif > 10)){

```



```

1
2
3           sim.new = sim.new[,which(names(sim.new) != names(tempvif)
4 [which(tempvif == tempvif[order(-tempvif)][1]])]]
5       }
6   }
7   try(
8       (manown = regsubsets(manic~., data = sim.new, really.big=T,
9 method = "forward", nvmax = 100)), silent = T
10      )
11      manownsummary = NA
12      if(any(is.na(manown)) == FALSE){
13          manownsummary = summary(manown)
14      }
15
16      ##own and bias variables
17      manbias = NA
18
19      manbias = regsubsets(as.formula(paste0("manic ~ ",
20 paste0(names(sim)[grepl("man_", names(sim)) == TRUE & grepl("bias",
21 names(sim)) == TRUE], collapse = " + "), collapse = "")), data = sim,
22 nvmax = 100, really.big=T, method = "forward")
23      manbiassummary = summary(manbias)
24      manbiassummary$adjr2
25
26
27      ##all variables
28      ###in case of collinearity
29      manall = NA
30      library(car)
31      sim.new = sim[,c("manic", names(summary((lm(manic~., data =
32 sim))))$aliased)[summary((lm(manic~., data = sim))))$aliased == FALSE &
33 names(summary((lm(manic~., data = sim))))$aliased) != "(Intercept)"]]
34      for(repe in 1:40){
35          tempvif = vif(lm(manic~., data = sim.new))
36          if(any(tempvif > 10)){
37              sim.new = sim.new[,which(names(sim.new) != names(tempvif)
38 [which(tempvif == tempvif[order(-tempvif)][1]])]]
39          }
40      }
41      ##Somehow there are problems in executing regsubsets even after
42 removing collinear variables
43      try(
44          (manall = regsubsets(manic~., data = sim.new, really.big=T,
45 method = "forward", nvmax = 100)), silent = T
46          )
47          manallsummary = NA
48          if(any(is.na(manall)) == FALSE){
49              manallsummary = summary(manall)
50          }
51
52      ##extract information from the outmat
53      #MDE
54
55      resu[which(resu$variable == "mde"), paste0("appbyownr2_", c,
56 collapse = "")] = mdeownsummary$adjr2[which.max(mdeownsummary$adjr2)]
57      resu[which(resu$variable == "mde"), paste0("appbyownn_", c,
58 collapse = "")] = which.max(mdeownsummary$adjr2)
59
60

```

```

1
2
3     resu[which(resu$variable == "mde"), paste0("appbyownvar_", c,
4 collapse = "")] = paste0(dimnames(mdeownsummary$outmat)[[2]]
5 [which(mdeownsummary$outmat[which.max(mdeownsummary$adjr2),] == "*")],
6 collapse = ",")
7
8     resu[which(resu$variable == "mde"), paste0("appbybiasr2_", c,
9 collapse = "")] = mdebiassummary$adjr2[which.max(mdebiassummary$adjr2)]
10    resu[which(resu$variable == "mde"), paste0("appbybiasn_", c,
11 collapse = "")] = which.max(mdebiassummary$adjr2)
12    resu[which(resu$variable == "mde"), paste0("appbybiasvar_", c,
13 collapse = "")] = paste0(dimnames(mdebiassummary$outmat)[[2]]
14 [which(mdebiassummary$outmat[which.max(mdebiassummary$adjr2),] == "*")],
15 collapse = ",")
16
17    if(any(is.na(mdeall)) == FALSE){
18        resu[which(resu$variable == "mde"), paste0("appbyallr2_",
19 c, collapse = "")] = mdeallsummary$adjr2[which.max(mdeallsummary$adjr2)]
20        resu[which(resu$variable == "mde"), paste0("appbyalln_", c,
21 collapse = "")] = which.max(mdeallsummary$adjr2)
22        resu[which(resu$variable == "mde"), paste0("appbyallvar_",
23 c, collapse = "")] = paste0(dimnames(mdeallsummary$outmat)[[2]]
24 [which(mdeallsummary$outmat[which.max(mdeallsummary$adjr2),] == "*")],
25 collapse = ",")
26    }
27
28
29
30    #DYS
31    resu[which(resu$variable == "dys"), paste0("appbyownr2_", c,
32 collapse = "")] = dysownsummary$adjr2[which.max(dysownsummary$adjr2)]
33    resu[which(resu$variable == "dys"), paste0("appbyownn_", c,
34 collapse = "")] = which.max(dysownsummary$adjr2)
35    resu[which(resu$variable == "dys"), paste0("appbyownvar_", c,
36 collapse = "")] = paste0(dimnames(dysownsummary$outmat)[[2]]
37 [which(dysownsummary$outmat[which.max(dysownsummary$adjr2),] == "*")],
38 collapse = ",")
39
40    resu[which(resu$variable == "dys"), paste0("appbybiasr2_", c,
41 collapse = "")] = dysbiassummary$adjr2[which.max(dysbiassummary$adjr2)]
42    resu[which(resu$variable == "dys"), paste0("appbybiasn_", c,
43 collapse = "")] = which.max(dysbiassummary$adjr2)
44    resu[which(resu$variable == "dys"), paste0("appbybiasvar_", c,
45 collapse = "")] = paste0(dimnames(dysbiassummary$outmat)[[2]]
46 [which(dysbiassummary$outmat[which.max(dysbiassummary$adjr2),] == "*")],
47 collapse = ",")
48
49    if(any(is.na(dysall)) == FALSE){
50        resu[which(resu$variable == "dys"), paste0("appbyallr2_", c,
51 collapse = "")] = dysallsummary$adjr2[which.max(dysallsummary$adjr2)]
52        resu[which(resu$variable == "dys"), paste0("appbyalln_", c,
53 collapse = "")] = which.max(dysallsummary$adjr2)
54        resu[which(resu$variable == "dys"), paste0("appbyallvar_", c,
55 collapse = "")] = paste0(dimnames(dysallsummary$outmat)[[2]]
56 [which(dysallsummary$outmat[which.max(dysallsummary$adjr2),] == "*")],
57 collapse = ",")
58    }
59
60

```

```

1
2
3     #MANIC
4     resu[which(resu$variable == "manic"), paste0("appbyownr2_", c,
5 collapse = "")] = manownsummary$adjr2[which.max(manownsummary$adjr2)]
6     resu[which(resu$variable == "manic"), paste0("appbyownn_", c,
7 collapse = "")] = which.max(manownsummary$adjr2)
8     resu[which(resu$variable == "manic"), paste0("appbyownvar_", c,
9 collapse = "")] = paste0(dimnames(manownsummary$outmat)[[2]]
10 [which(manownsummary$outmat[which.max(manownsummary$adjr2),] == "*")],
11 collapse = ",")
12
13     resu[which(resu$variable == "manic"), paste0("appbybiasr2_", c,
14 collapse = "")] = manbiassummary$adjr2[which.max(manbiassummary$adjr2)]
15     resu[which(resu$variable == "manic"), paste0("appbybiasn_", c,
16 collapse = "")] = which.max(manbiassummary$adjr2)
17     resu[which(resu$variable == "manic"), paste0("appbybiasvar_",
18 c, collapse = "")] = paste0(dimnames(manbiassummary$outmat)[[2]]
19 [which(manbiassummary$outmat[which.max(manbiassummary$adjr2),] == "*")],
20 collapse = ",")
21
22
23     if(any(is.na(manall)) == FALSE){
24         resu[which(resu$variable == "manic"), paste0("appbyallr2_",
25 c, collapse = "")] = manallsummary$adjr2[which.max(manallsummary$adjr2)]
26         resu[which(resu$variable == "manic"), paste0("appbyalln_", c,
27 collapse = "")] = which.max(manallsummary$adjr2)
28         resu[which(resu$variable == "manic"), paste0("appbyallvar_",
29 c, collapse = "")] = paste0(dimnames(manallsummary$outmat)[[2]]
30 [which(manallsummary$outmat[which.max(manallsummary$adjr2),] == "*")],
31 collapse = ",")
32     }
33
34
35
36
37     print(c("c:", c))
38     print(c("cor:", rho[rh]))
39     print(c("Prevalence: ", prevalence[preval]))
40
41 }#c
42
43
44
45     ##adding summary statistics to the result data frame
46     resu[, paste(collect, "_mean", sep = "")] = NA
47     resu[, paste(collect, "_sd", sep = "")] = NA
48     resu[, paste(collect, "_se", sep = "")] = NA
49     resu[, paste(collect, "_95up", sep = "")] = NA
50     resu[, paste(collect, "_95lo", sep = "")] = NA
51     resu[, paste(collect, "_rangeup", sep = "")] = NA
52     resu[, paste(collect, "_rangelo", sep = "")] = NA
53
54     for(co in 1:length(collect)){
55         for(r in 1:nrow(resu)){
56             if((collect[co] %in% c("appbyownvar", "appbybiasvar", "appbyallvar"))
57 == FALSE){
58                 resu[r,paste0(collect[co], "_mean", collapse = "")] =
59 mean(unlist(resu[r, paste(collect[co], "_", 1:times, sep = "")])[which(!
60 is.na(unlist(resu[r, paste(collect[co], "_", 1:times, sep = ""))]))))

```

```

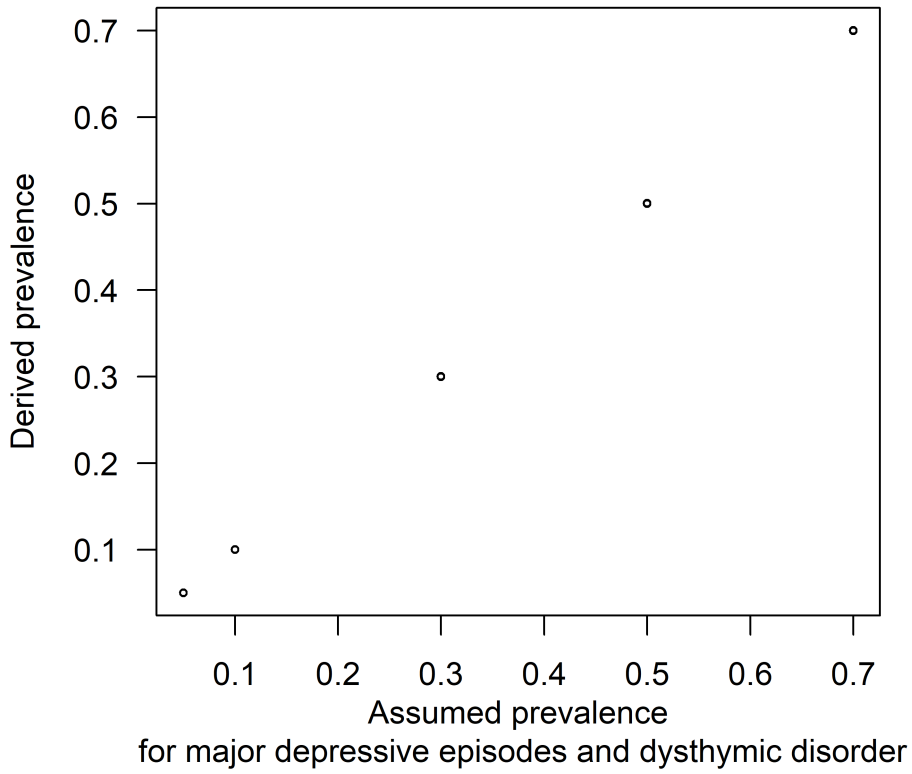
1
2
3     resu[r,paste0(collect[co], "_sd", collapse = "")] = sd(unlist(resu[r,
4 paste(collect[co], "_", 1:times, sep = "")])[which(!is.na(unlist(resu[r,
5 paste(collect[co], "_", 1:times, sep = ""))]))))
6     resu[r,paste0(collect[co], "_se", collapse = "")] = sd(unlist(resu[r,
7 paste(collect[co], "_", 1:times, sep = "")])[which(!is.na(unlist(resu[r,
8 paste(collect[co], "_", 1:times, sep = ""))])))))/(times^0.5)
9
10    #95% CIs
11    resu[r,paste0(collect[co], "_95up", collapse = "")] =
12 resu[r,paste0(collect[co], "_mean", collapse = "")] +
13 1.96*resu[r,paste0(collect[co], "_se", collapse = "")]
14    resu[r,paste0(collect[co], "_95lo", collapse = "")] =
15 resu[r,paste0(collect[co], "_mean", collapse = "")] -
16 1.96*resu[r,paste0(collect[co], "_se", collapse = "")]
17
18    #range
19    resu[r,paste0(collect[co], "_rangelo", collapse = "")] =
20 min(resu[r,paste(collect[co], "_", 1:times, sep = "")])
21    resu[r,paste0(collect[co], "_rangeup", collapse = "")] =
22 max(resu[r,paste(collect[co], "_", 1:times, sep = "")])
23 }#r
24
25    ##Add information about the aliased variables
26
27
28
29    ##save in another data set
30    eval(parse(text = paste0("resu_cor", rho[rh], "_preval",
31 prevalence[preval], " = resu", collapse = "")))
32
33
34 }
35
36 #export results
37 write.csv(cbind(resu[,c("definition", "variable", "mean_mean",
38 "mean_95up", "mean_95lo", "max_rangeup",
39 "min_rangelo", "derivedprevalence_mean", "derivedprevalence_95up",
40 "derivedprevalence_95lo", "coef_mean", "coef_95up", "coef_95lo",
41 "p_mean", "p_95up", "p_95lo", "r2_mean", "r2_95up",
42 "r2_95lo", "subcoef_mean", "subcoef_95up", "subcoef_95lo",
43 "subp_mean", "subp_95up", "subp_95lo", "subr2_mean", "subr2_95up",
44 "subr2_95lo", "appbyownr2_mean", "appbyownr2_95up", "appbyownr2_95lo",
45 "appbyownn_mean", "appbyownn_95up", "appbyownn_95lo", "appbybiasr2_mean",
46 "appbybiasr2_95up", "appbybiasr2_95lo", "appbybiasn_mean",
47 "appbybiasn_95up", "appbybiasn_95lo", "appbyallr2_mean",
48 "appbyallr2_95up", "appbyallr2_95lo", "appbyalln_mean", "appbyalln_95up",
49 "appbyalln_95lo"
50
51 )], resu), file = paste0("simulation results_cor", rho[rh], "_preval",
52 prevalence[preval], ".csv"))
53
54
55
56 }#co
57
58
59
60

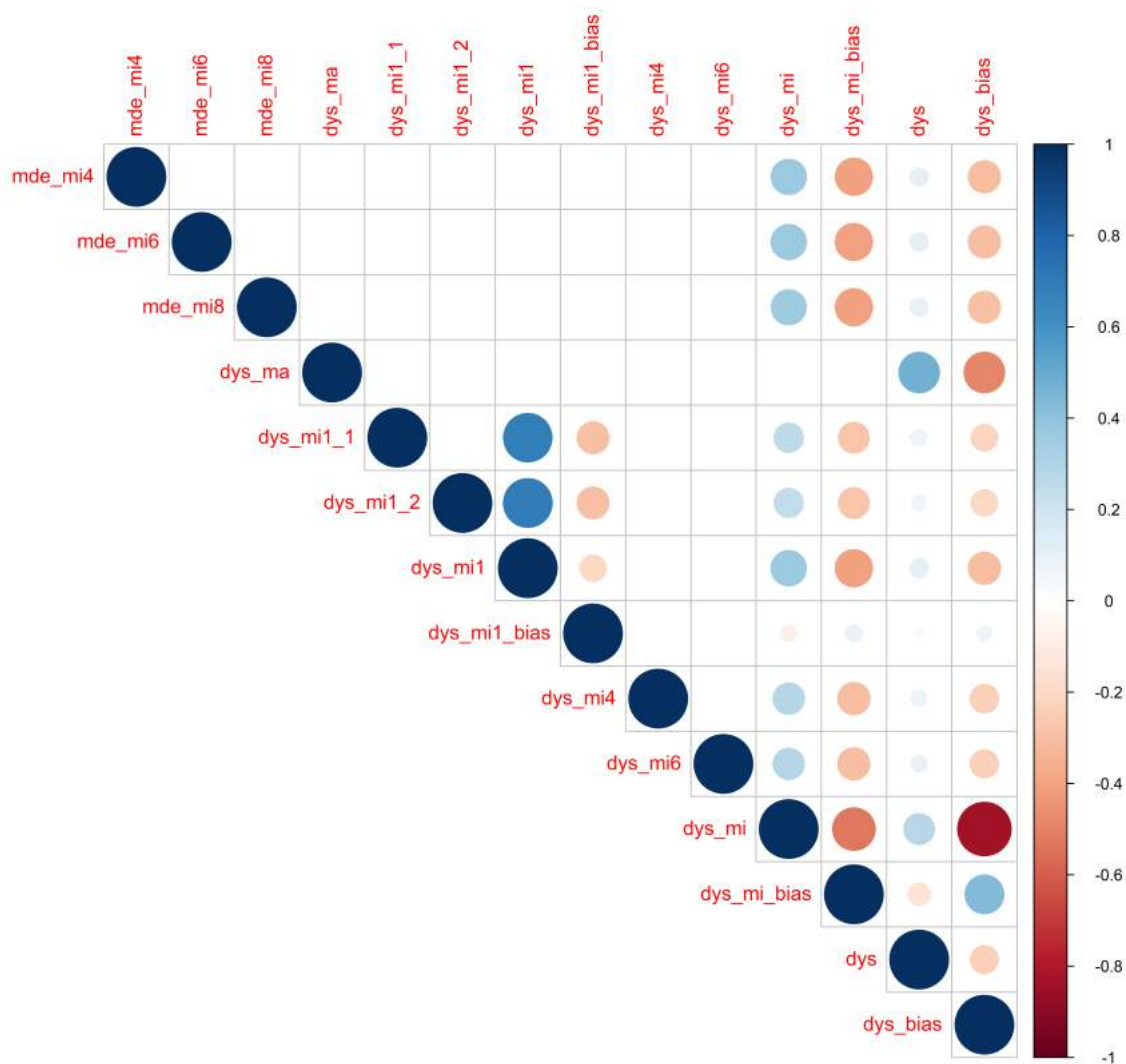
```

```
1
2
3     print(c("cor:", rho[rh]))
4     print(c("Prevalence: ", prevalence[preval]))
5
6         }#rho
7     #store data
8
9     print(c("Prevalence: ", prevalence[preval]))
10        }#prevalence
11
12
13     ...
```

For peer review only

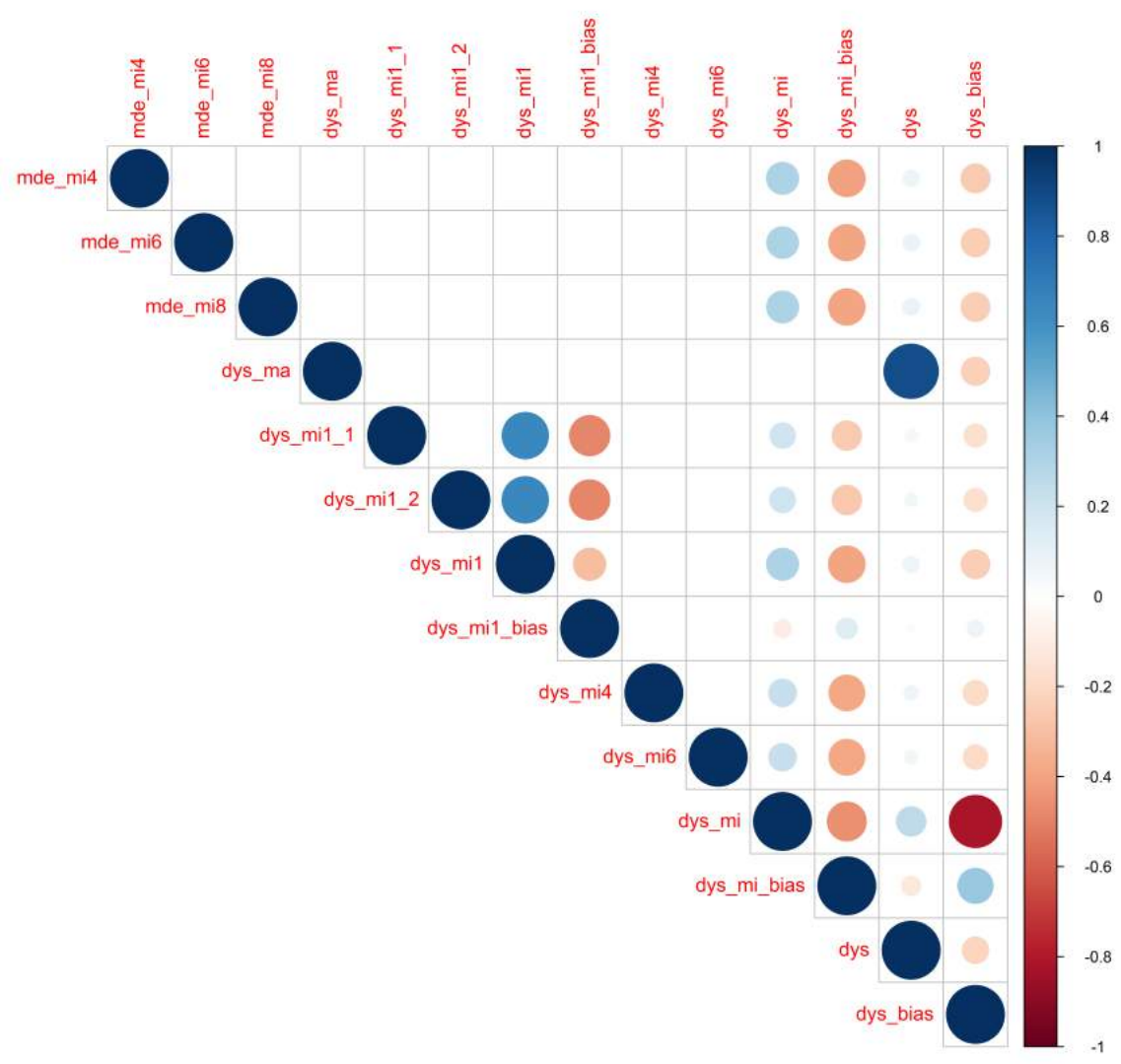
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



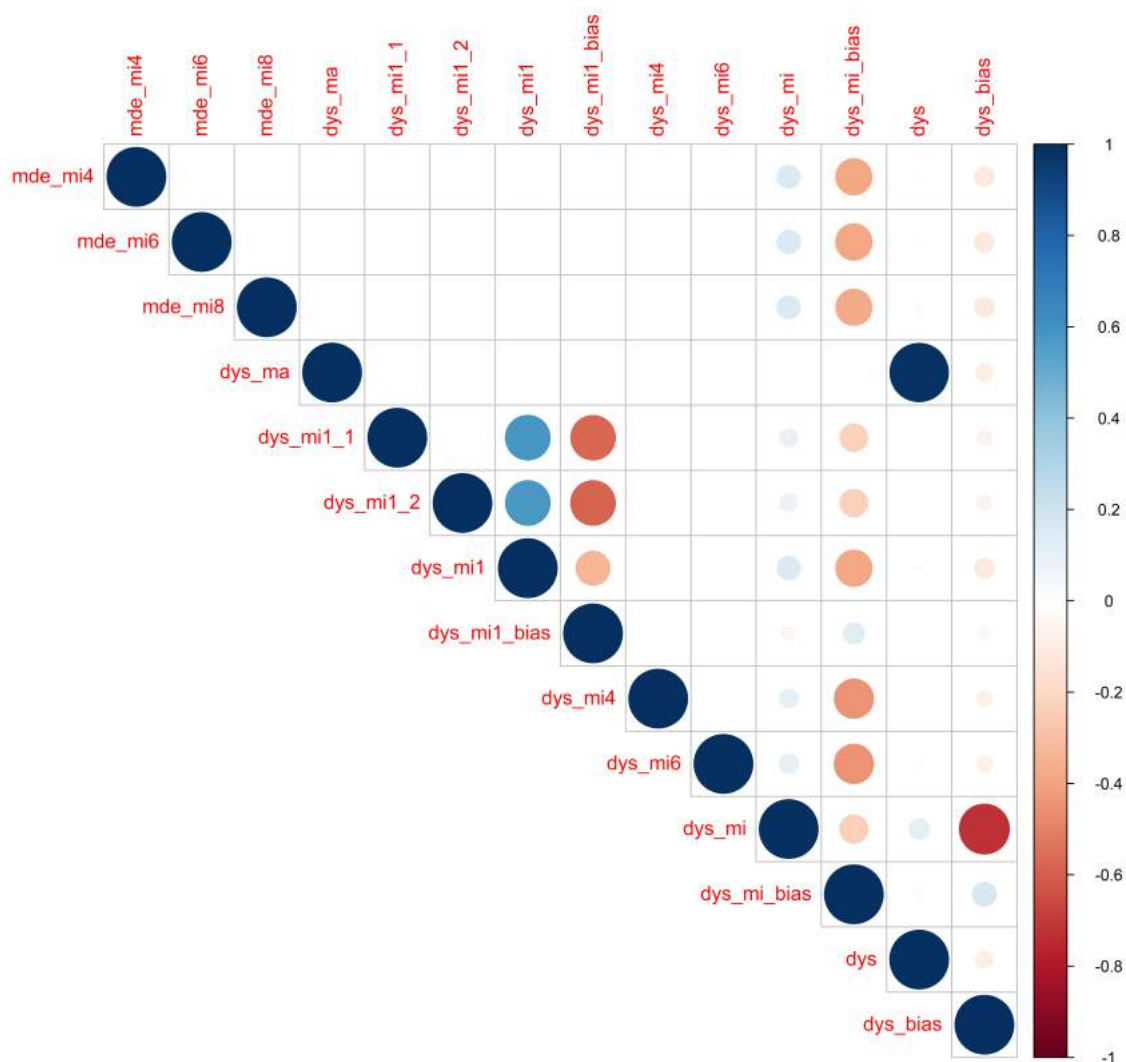


only

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

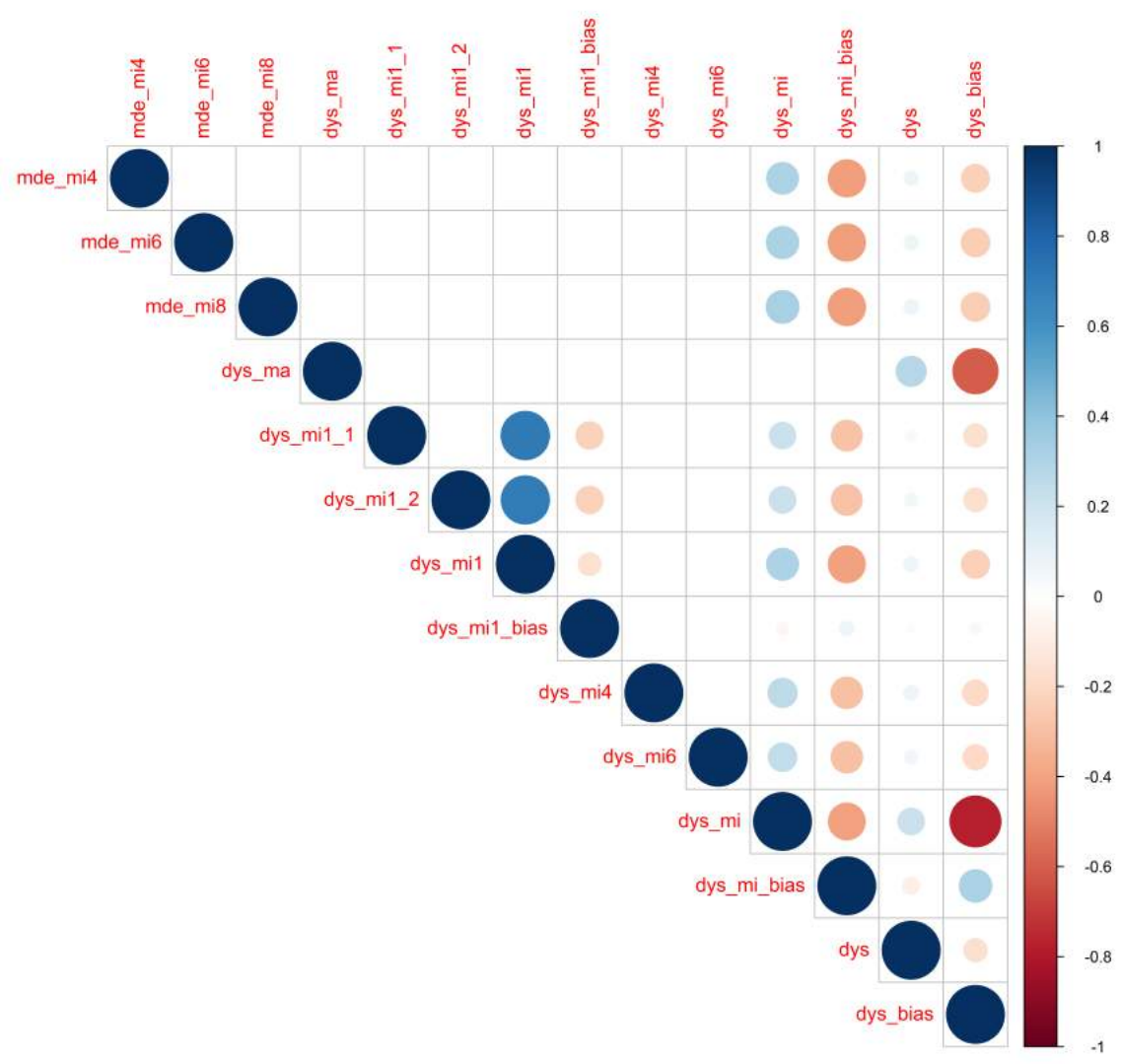


only

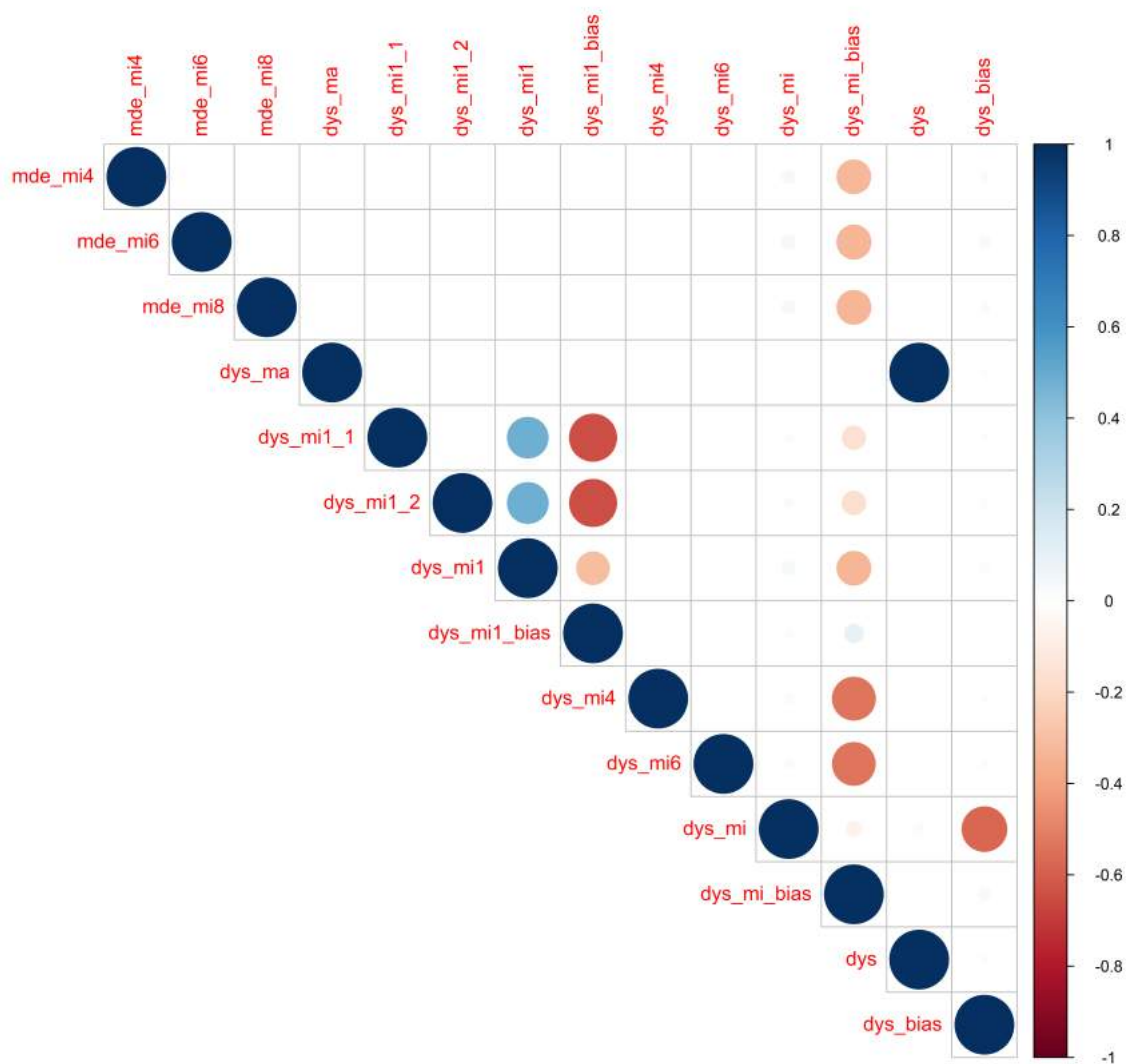


only

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

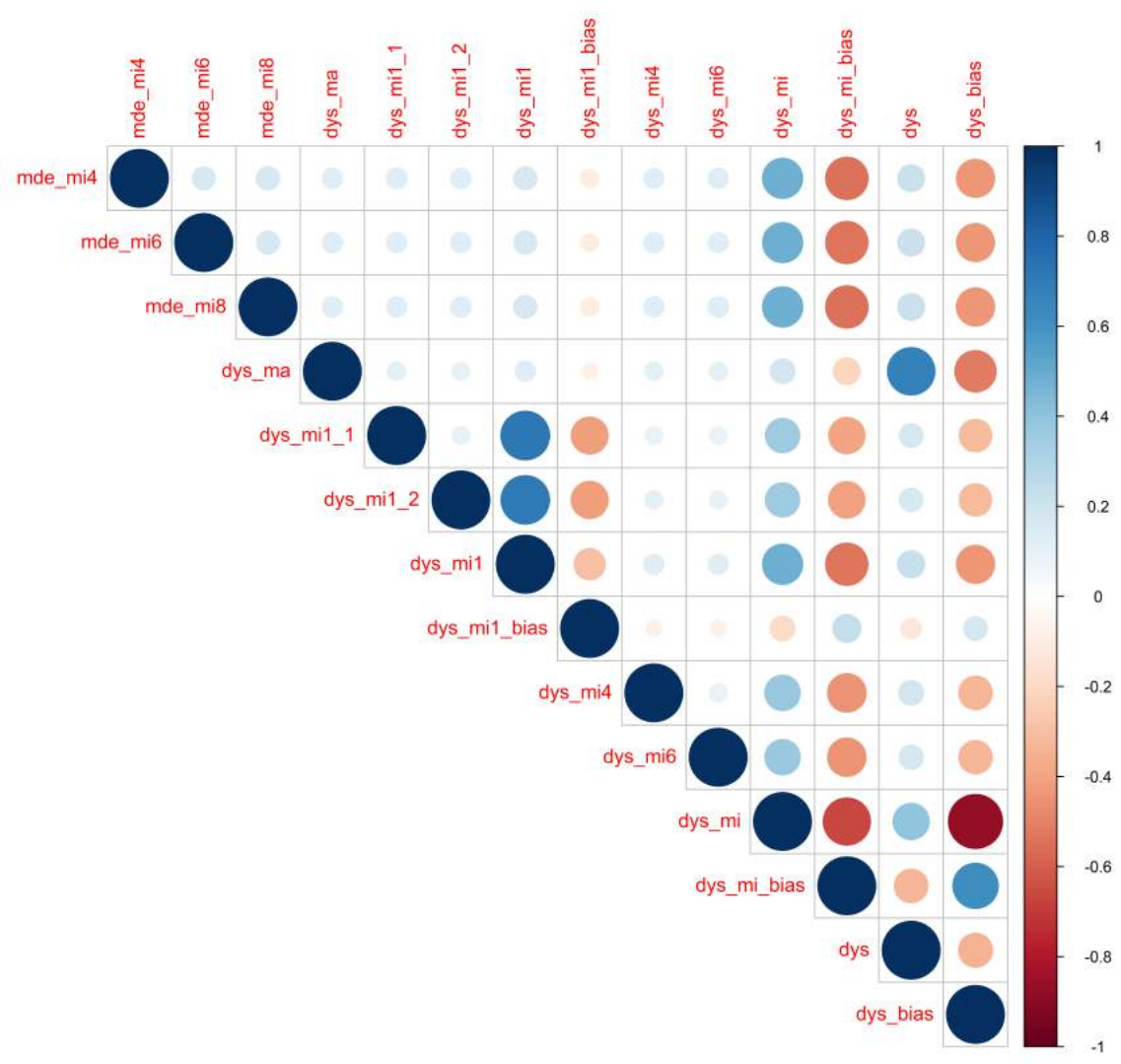


only

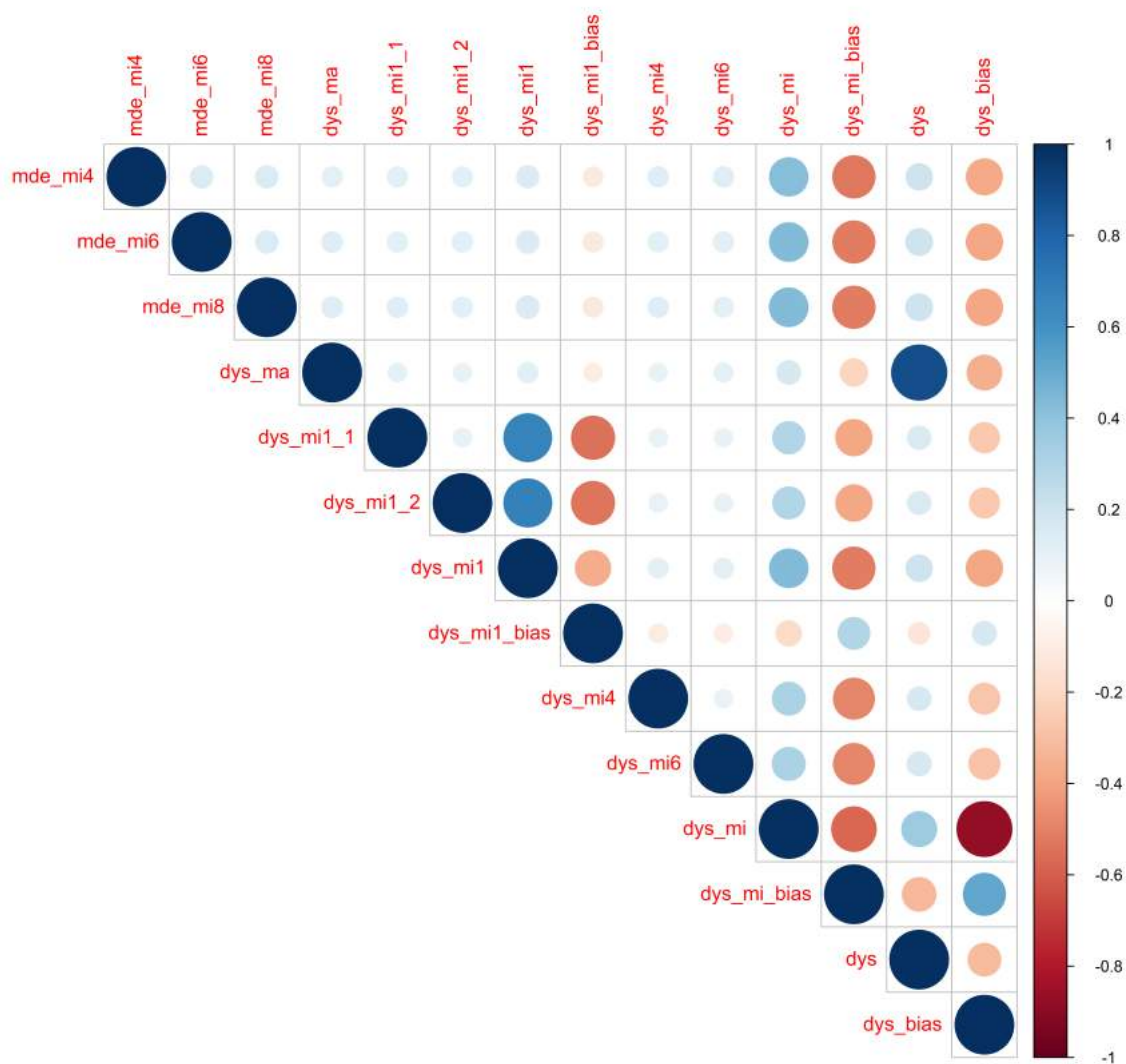


only

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

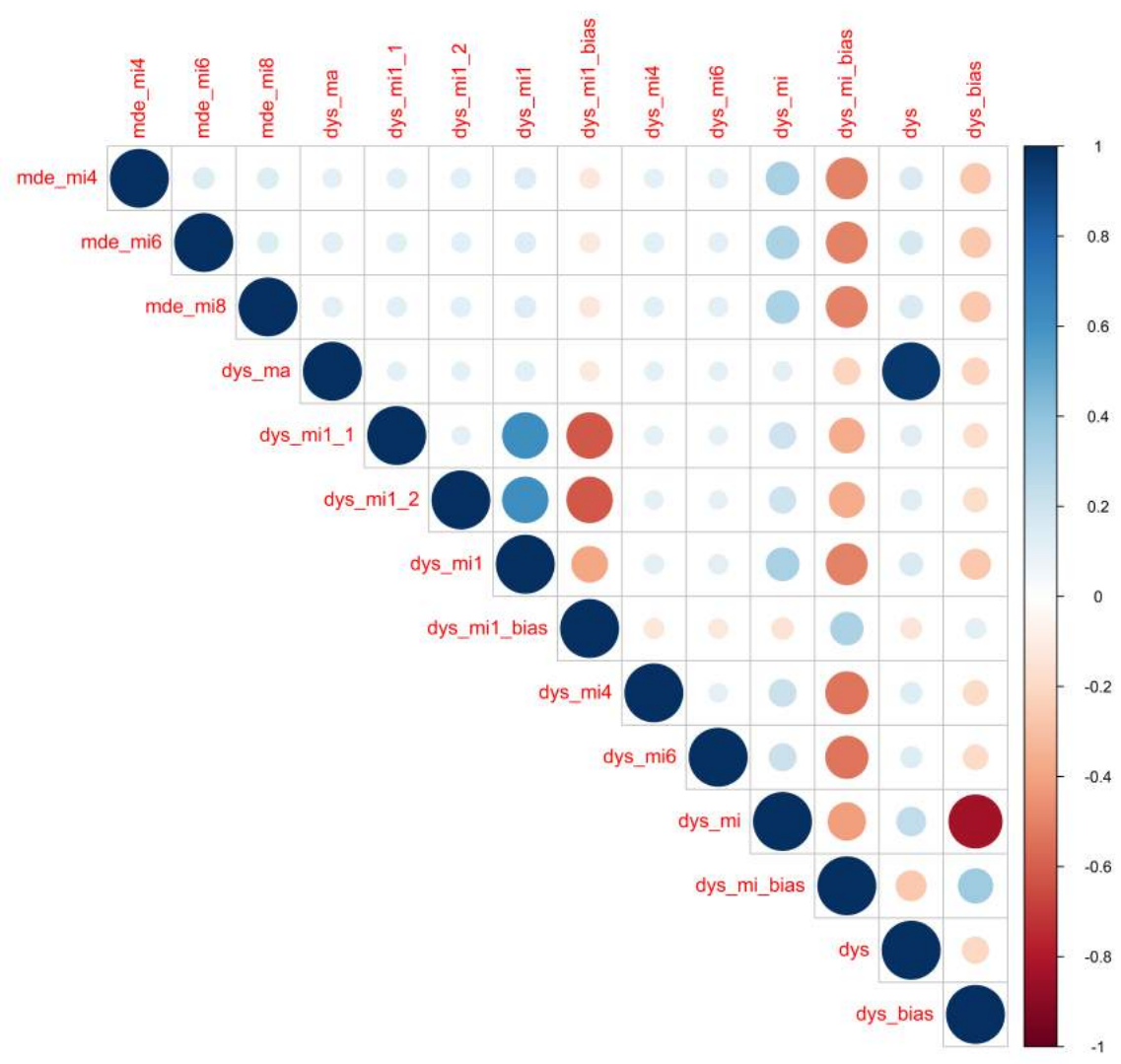


only

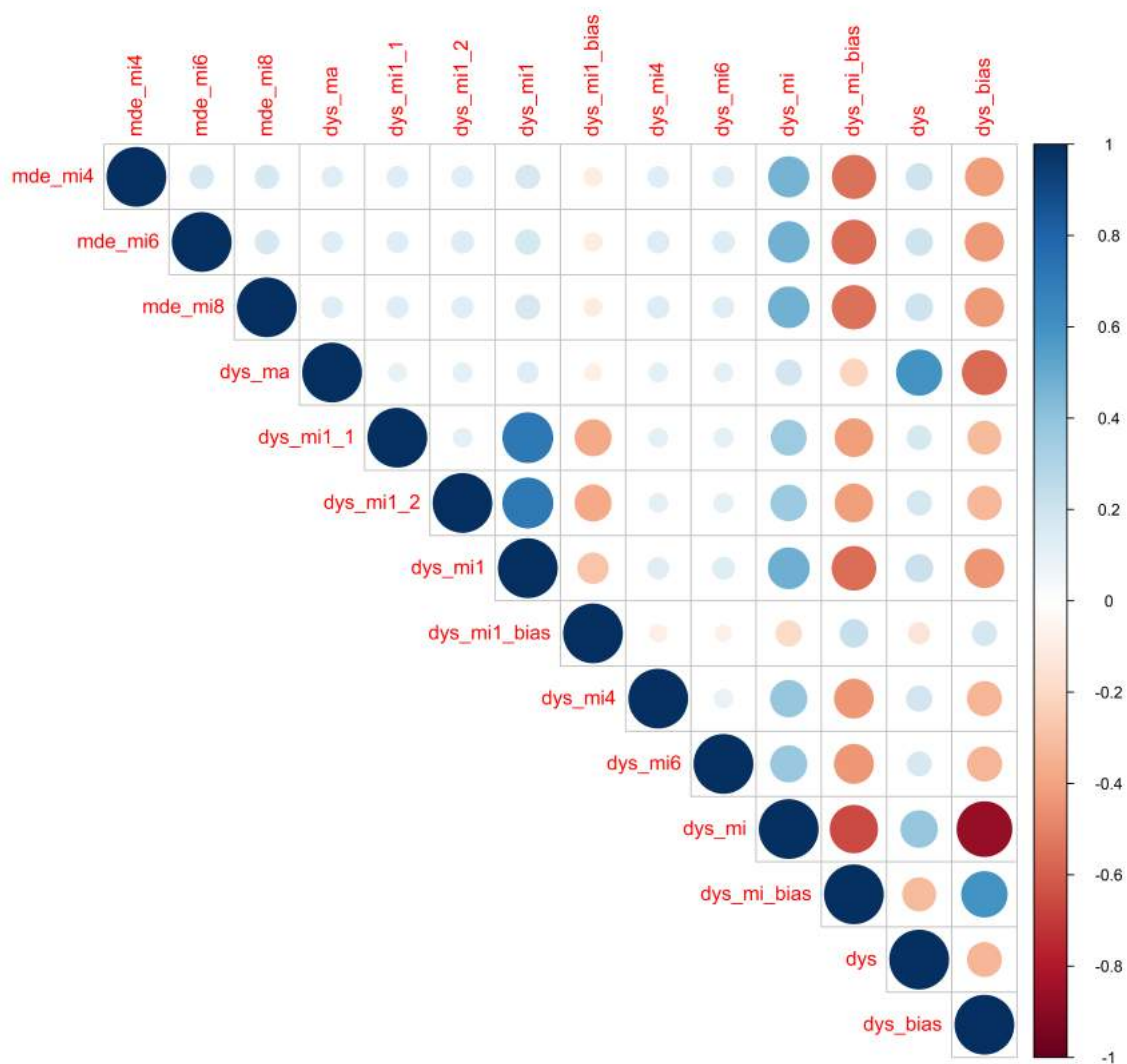


only

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

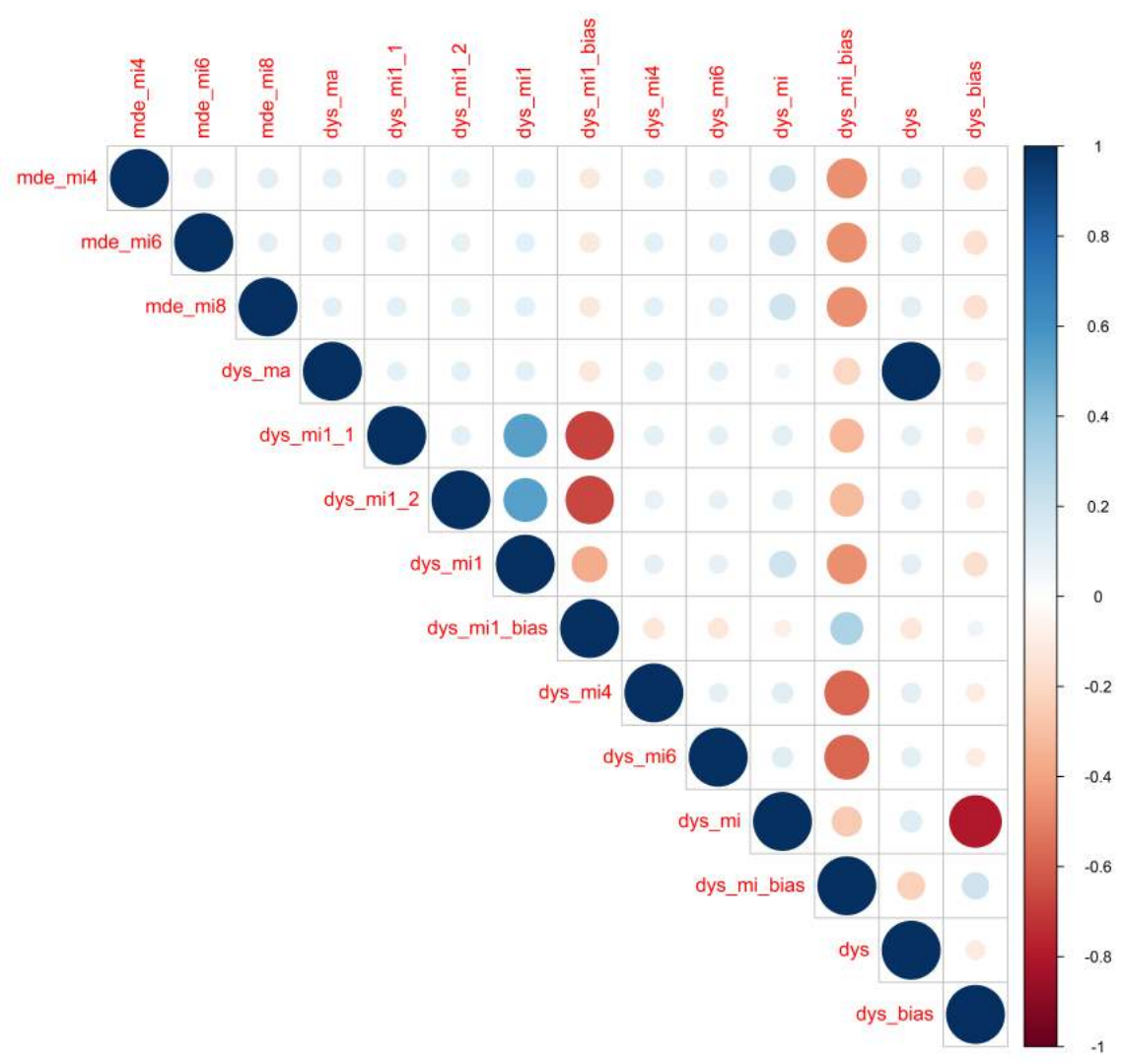


only

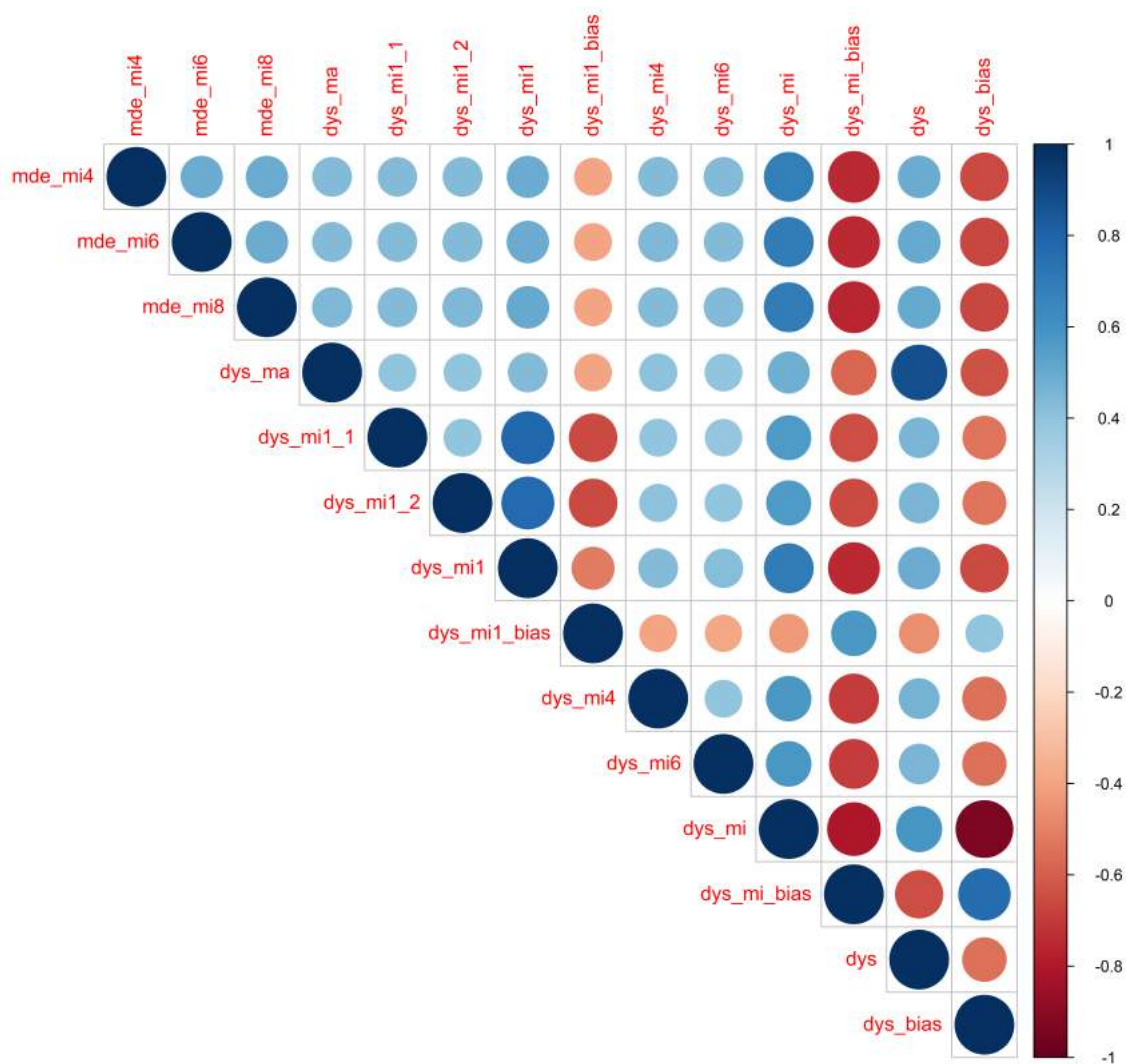


only

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

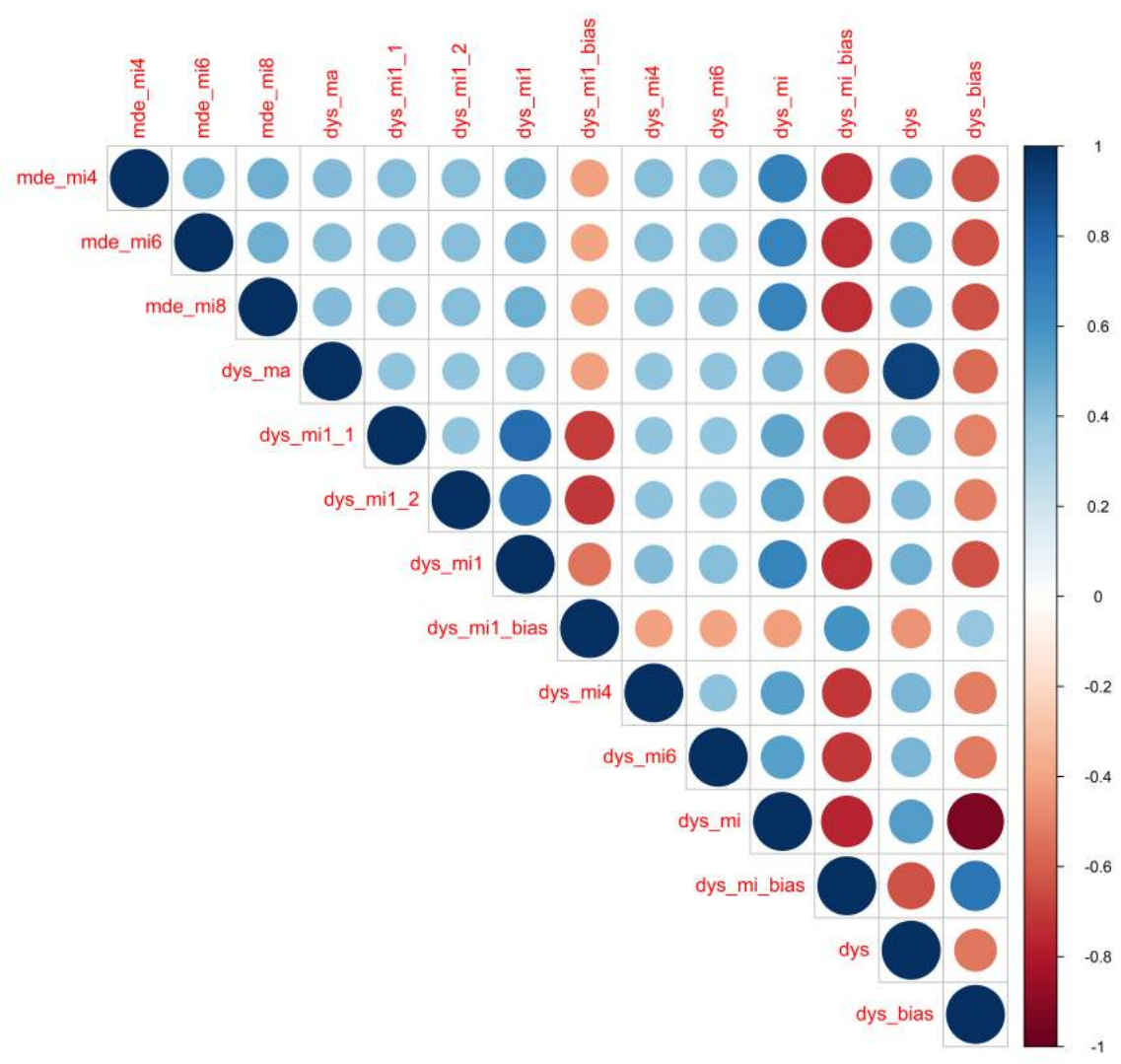


only

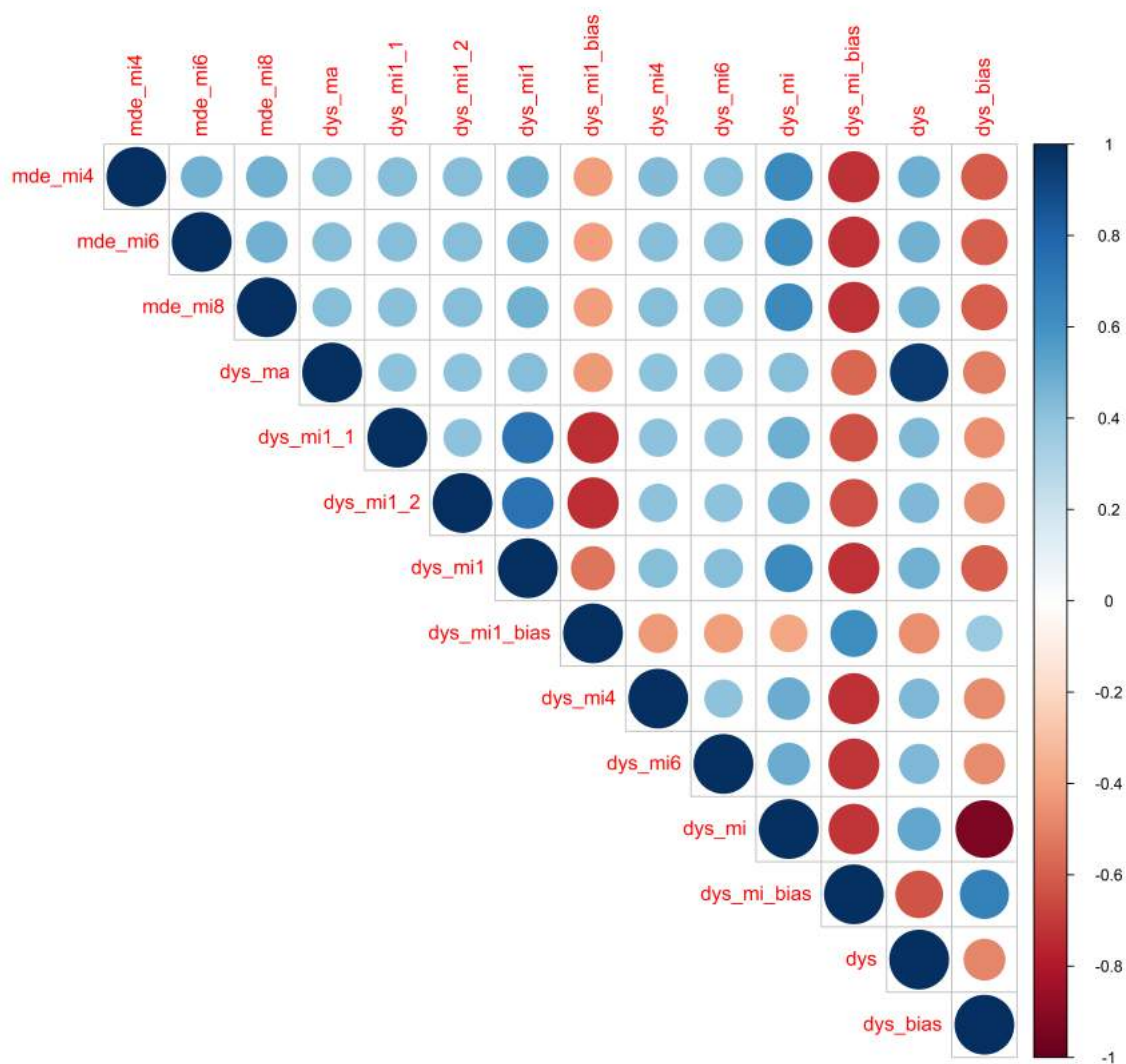


only

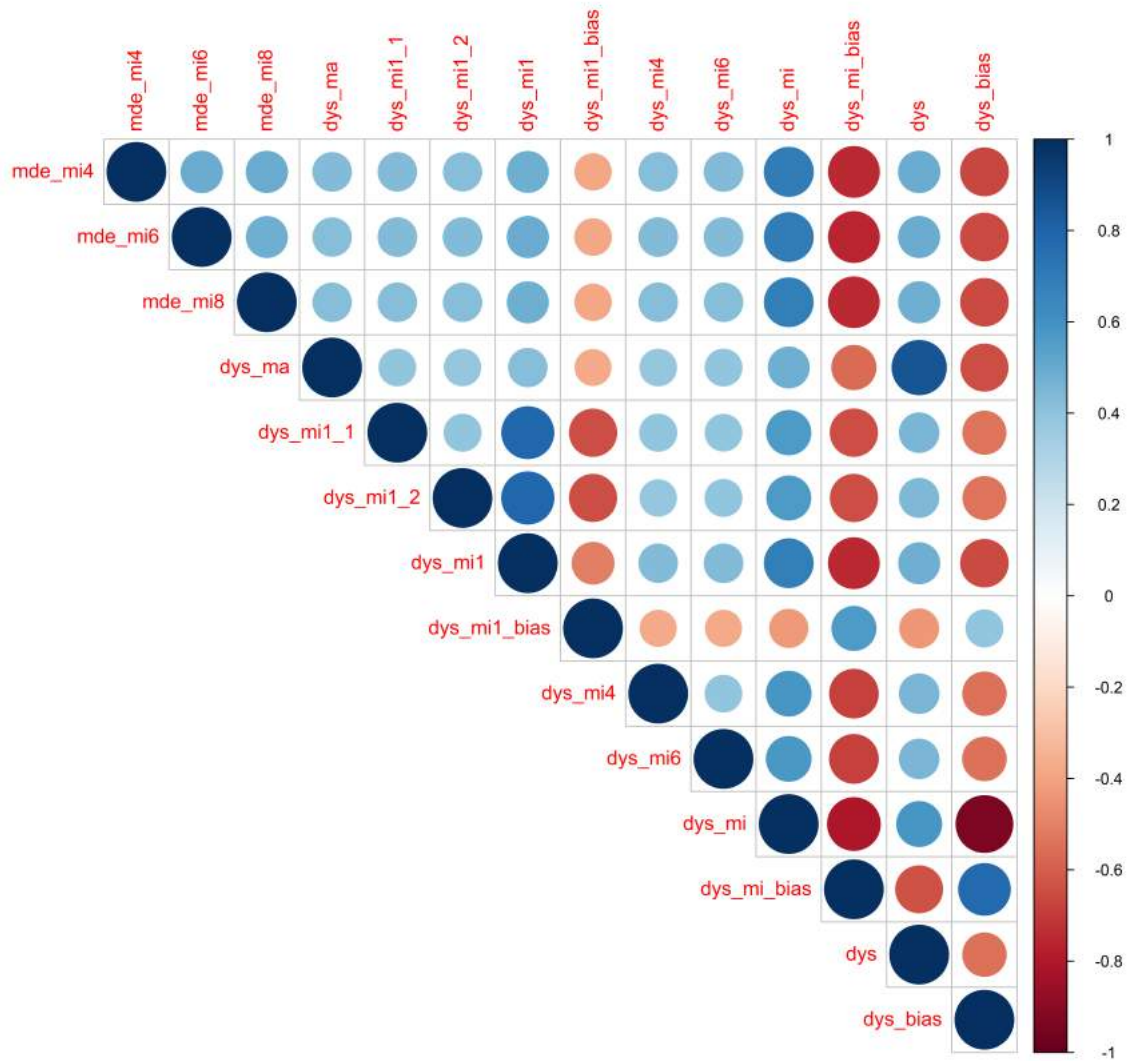
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



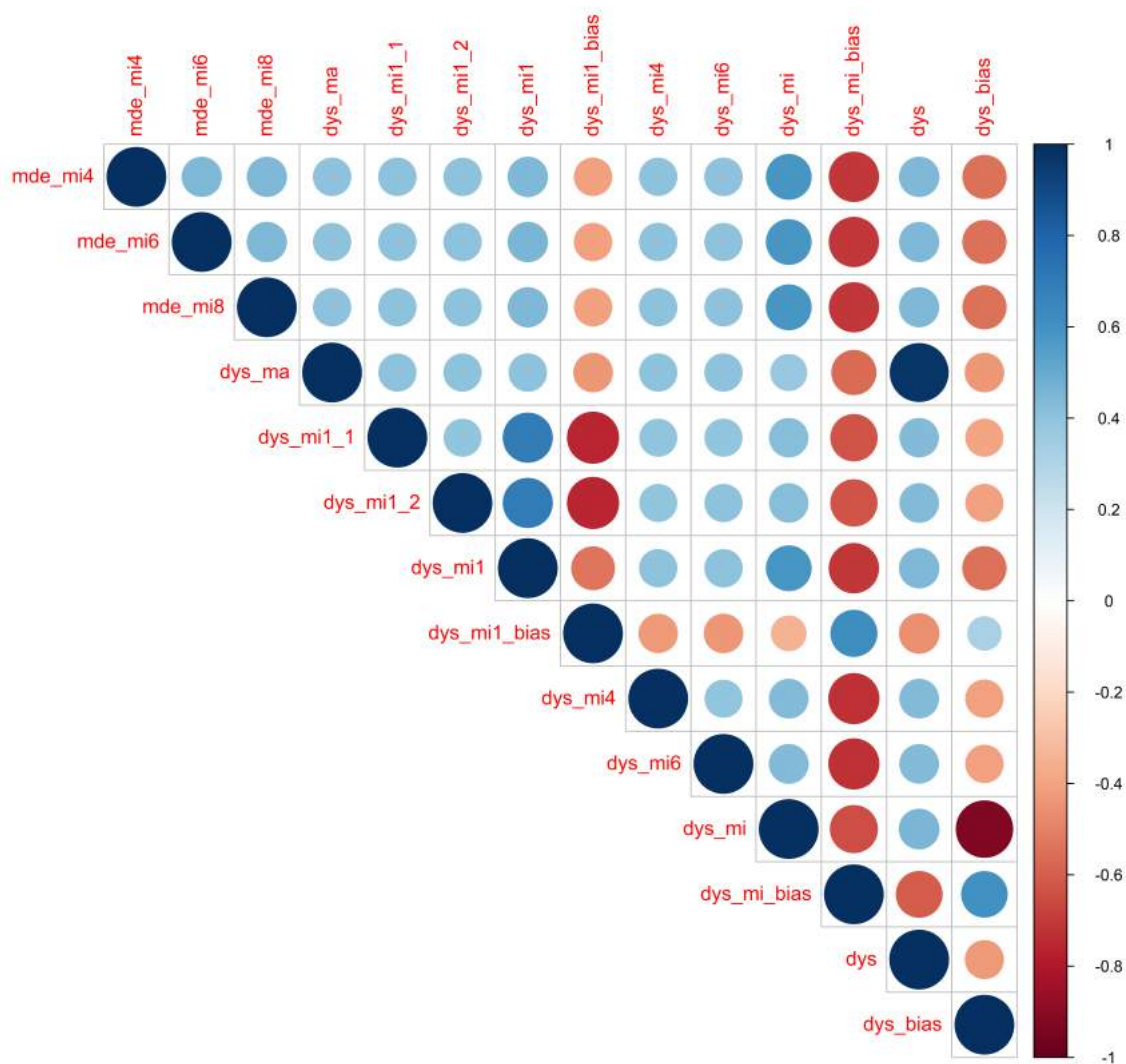
only



only

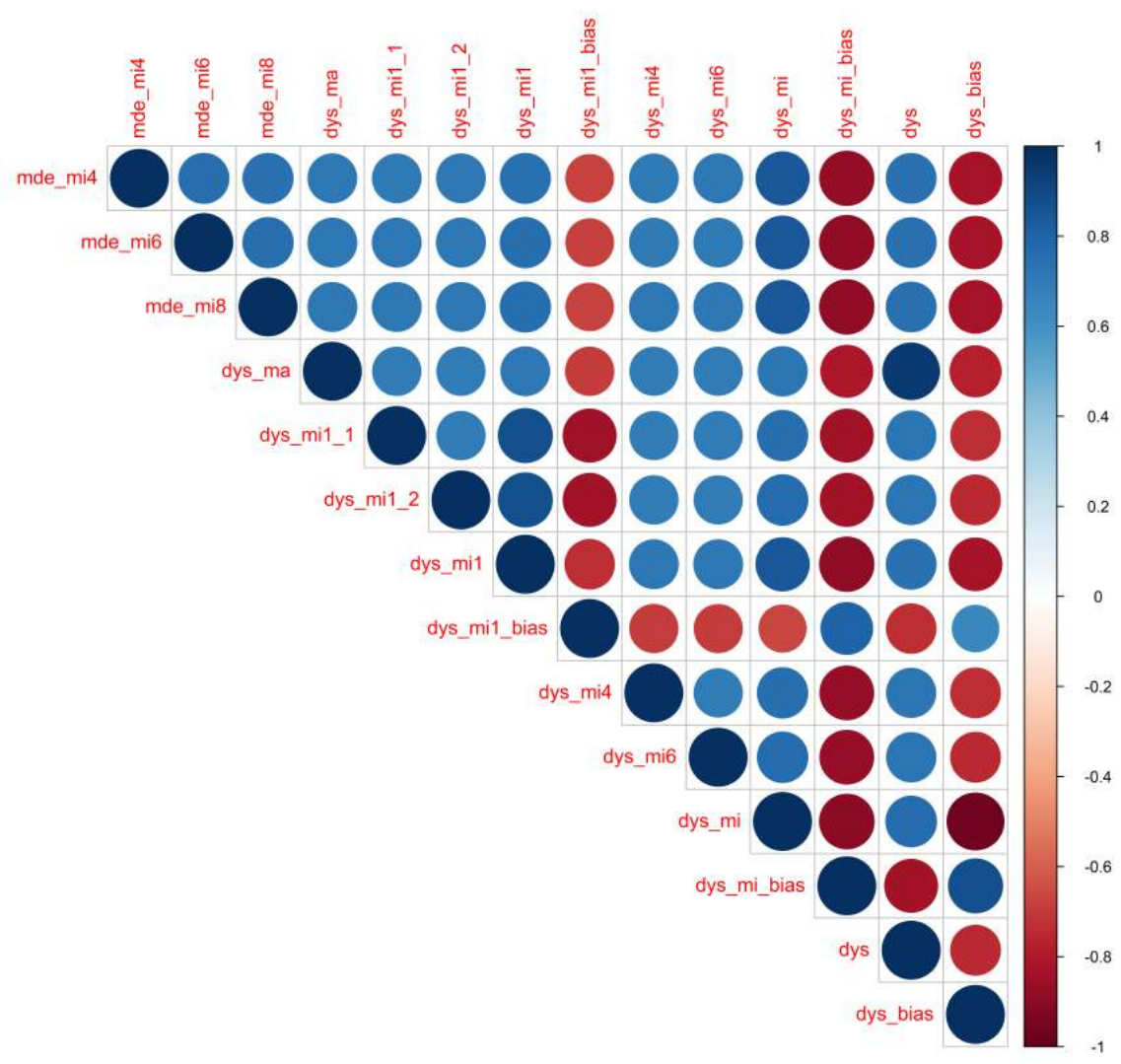


only

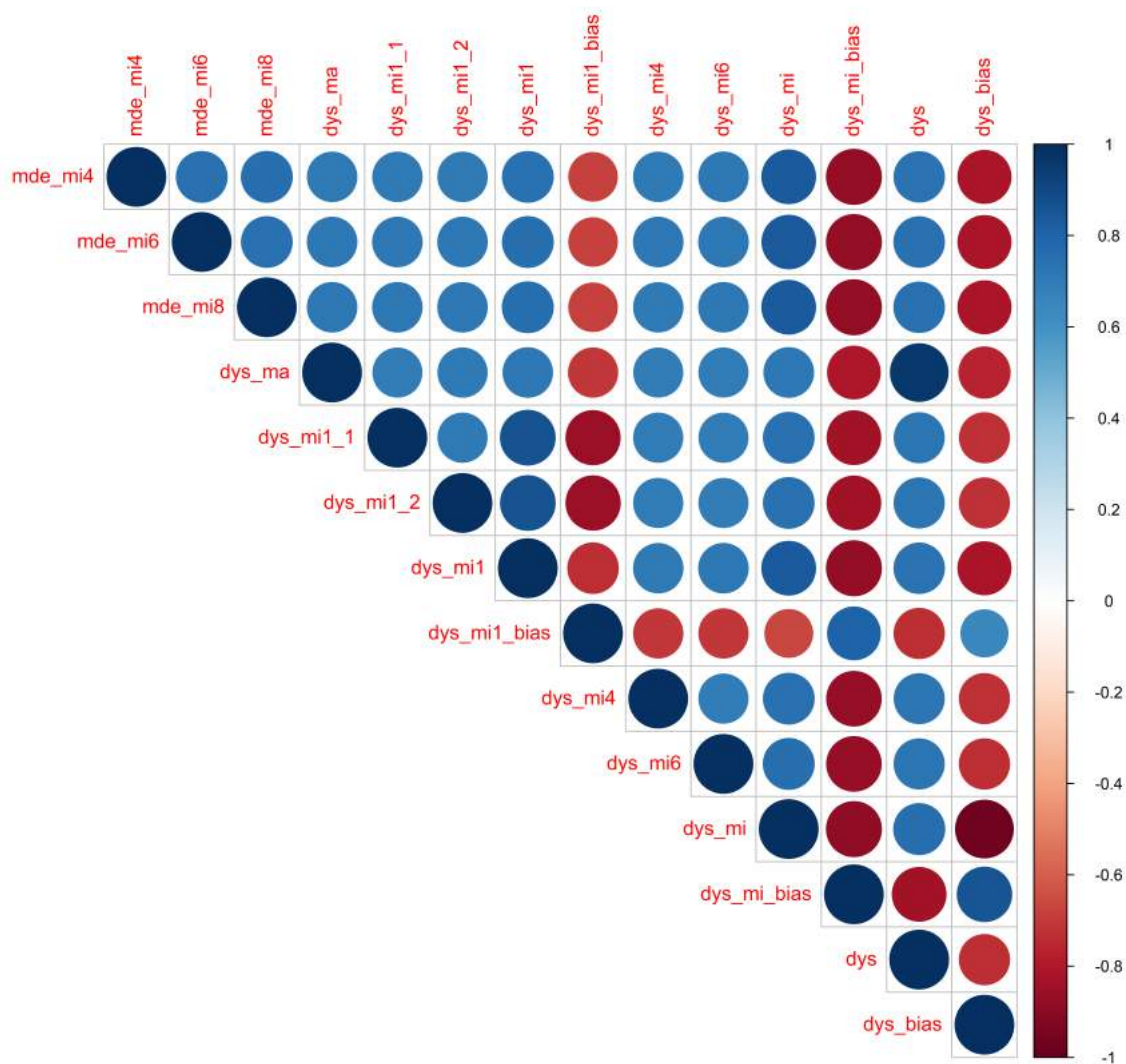


only

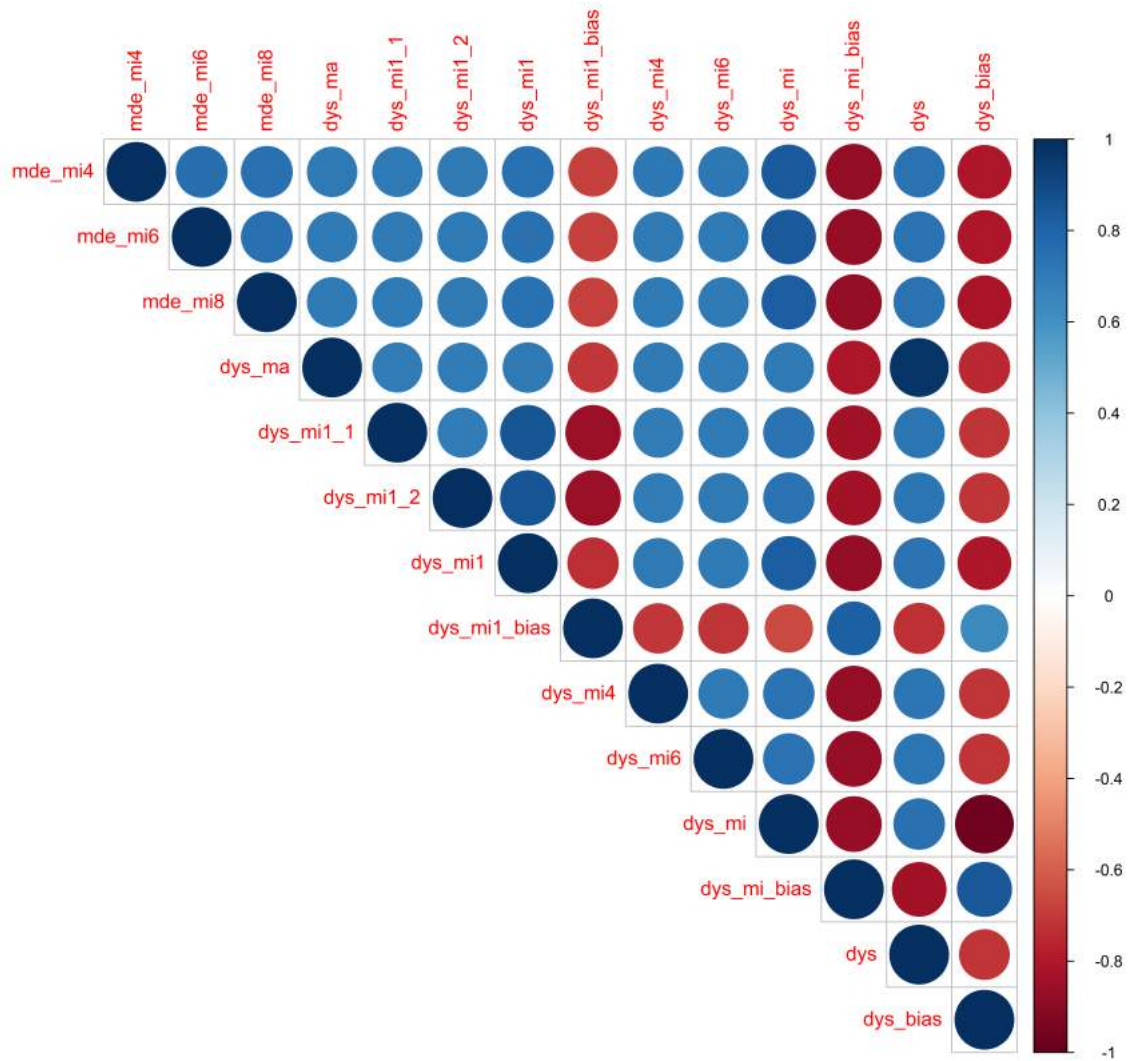
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



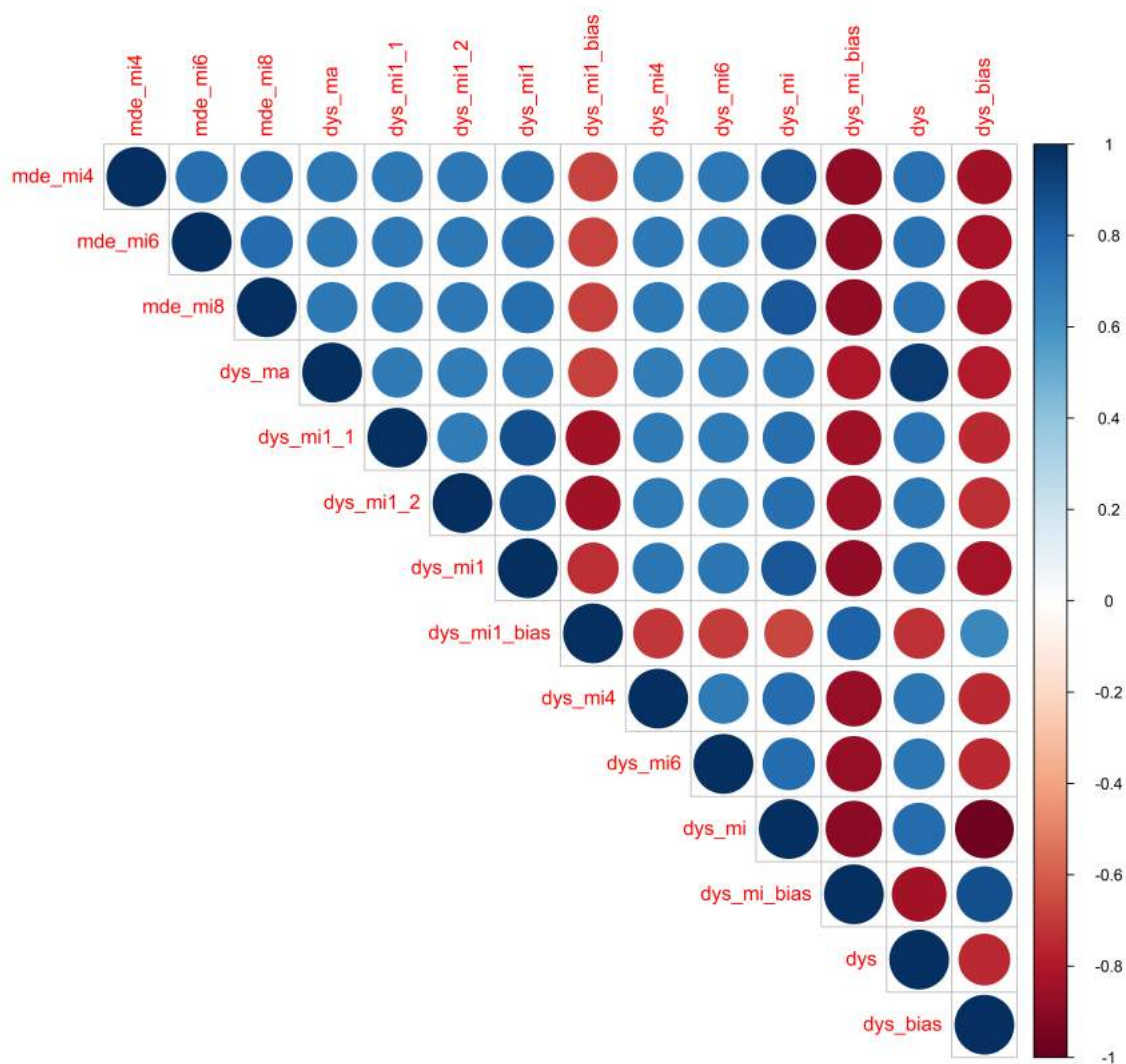
only



only

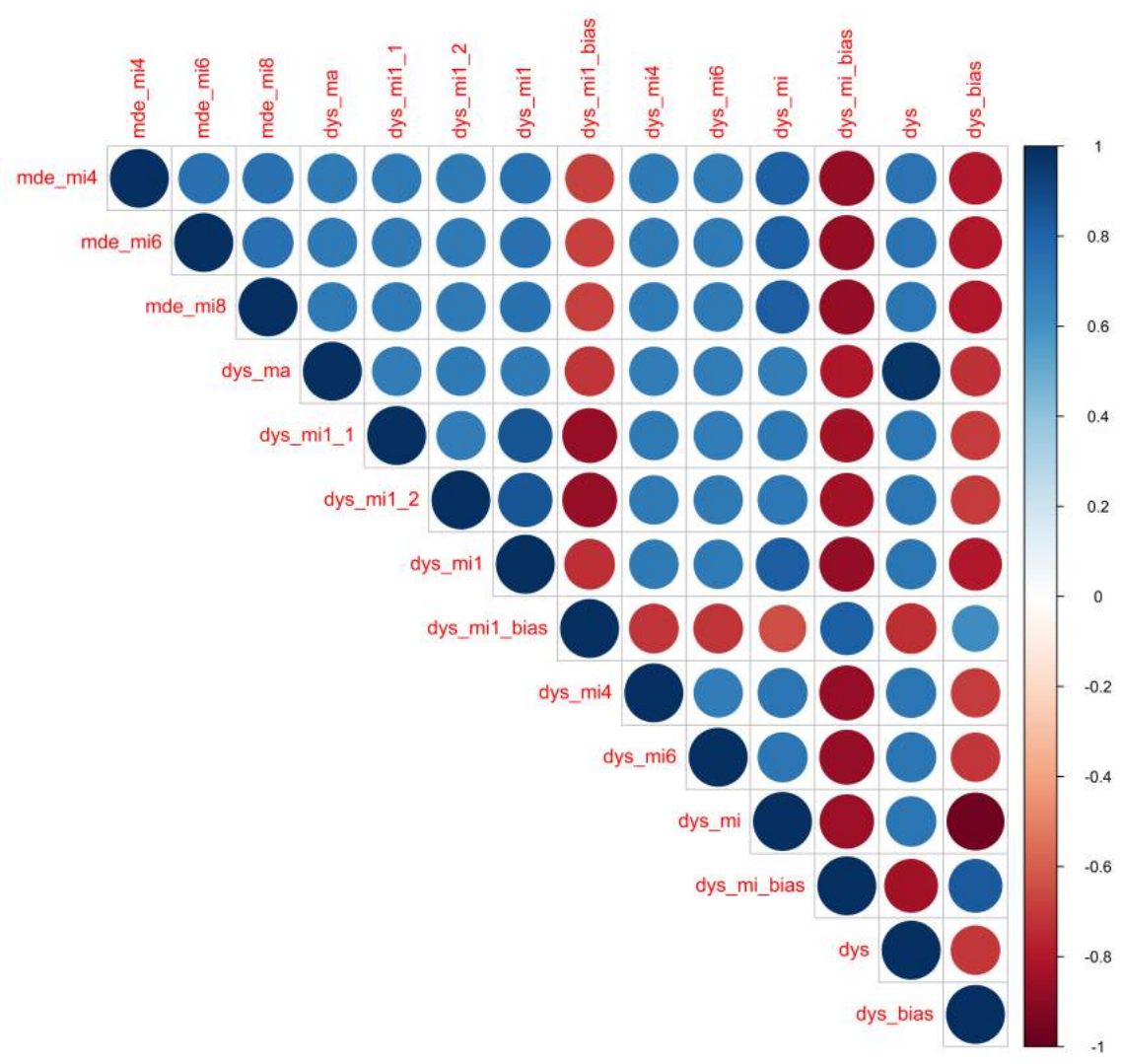


only

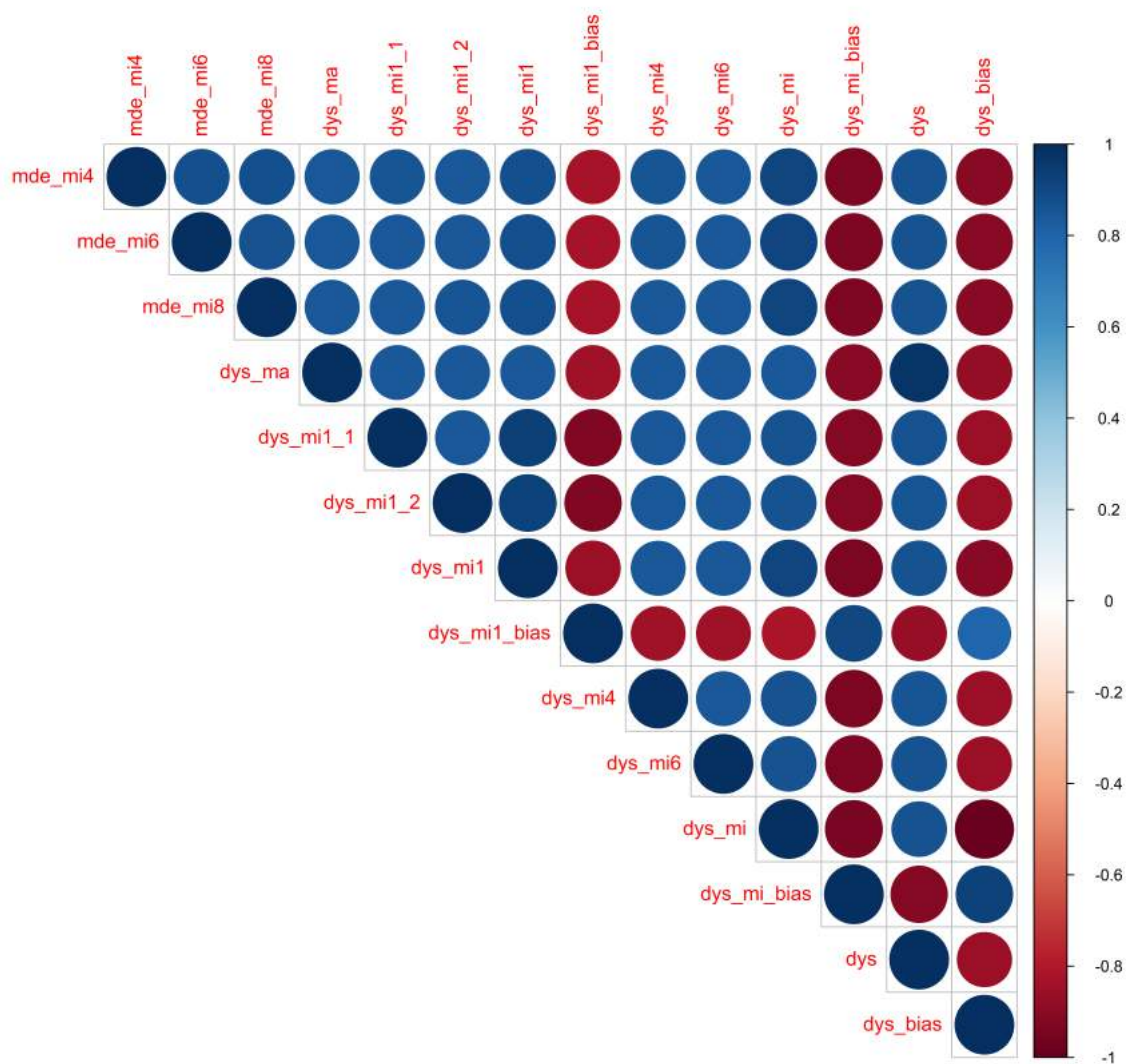


only

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

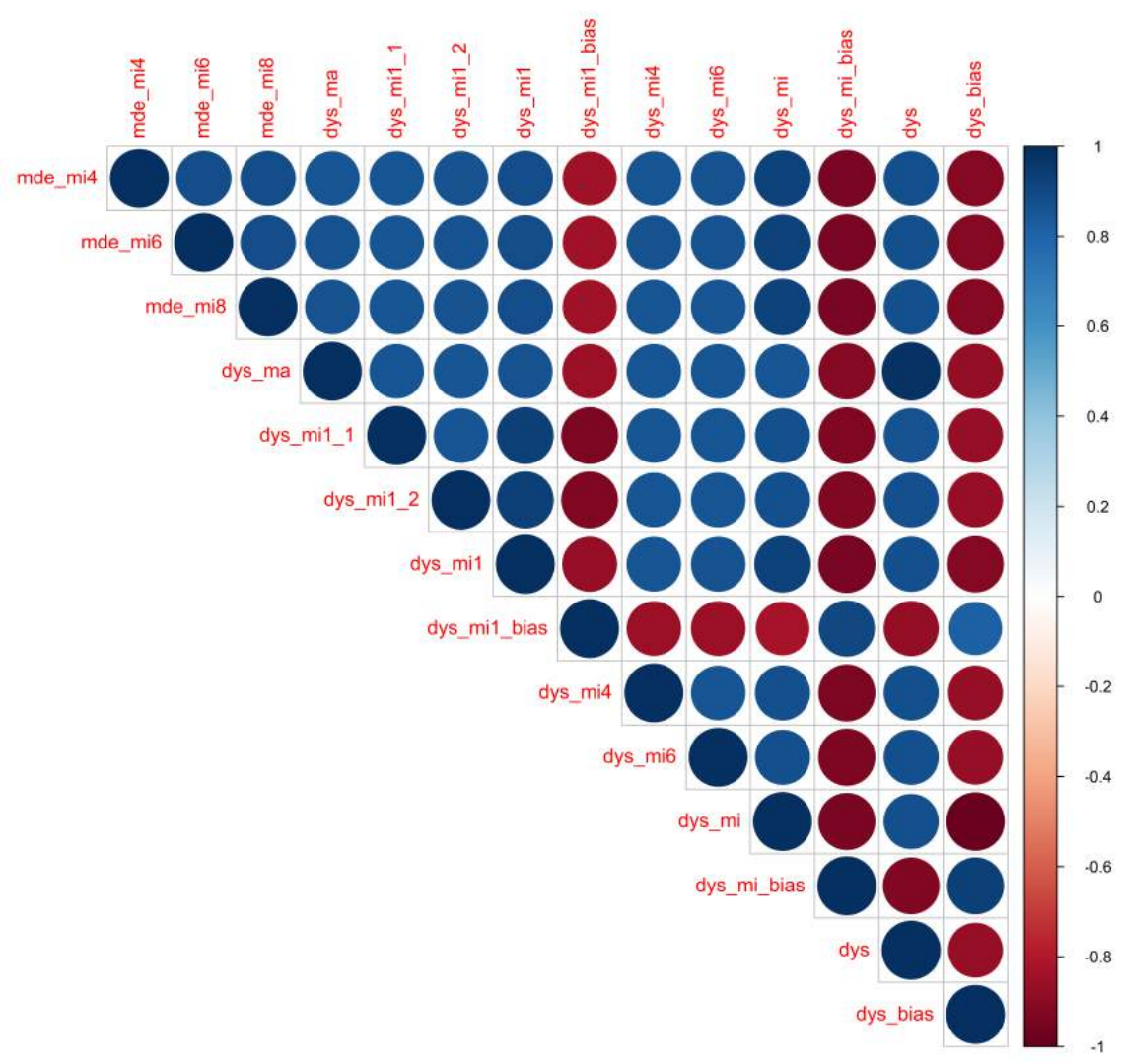


only

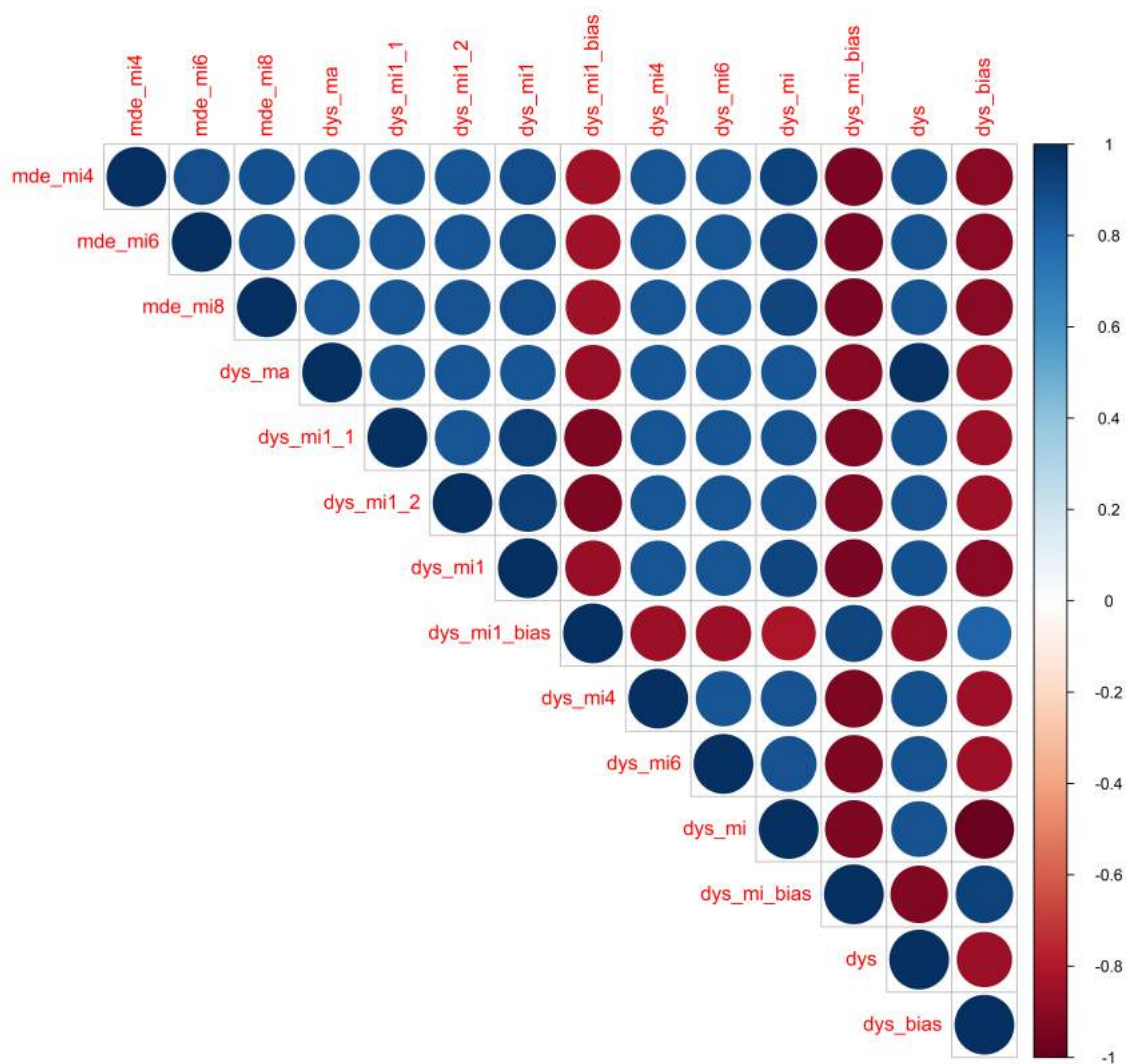


only

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

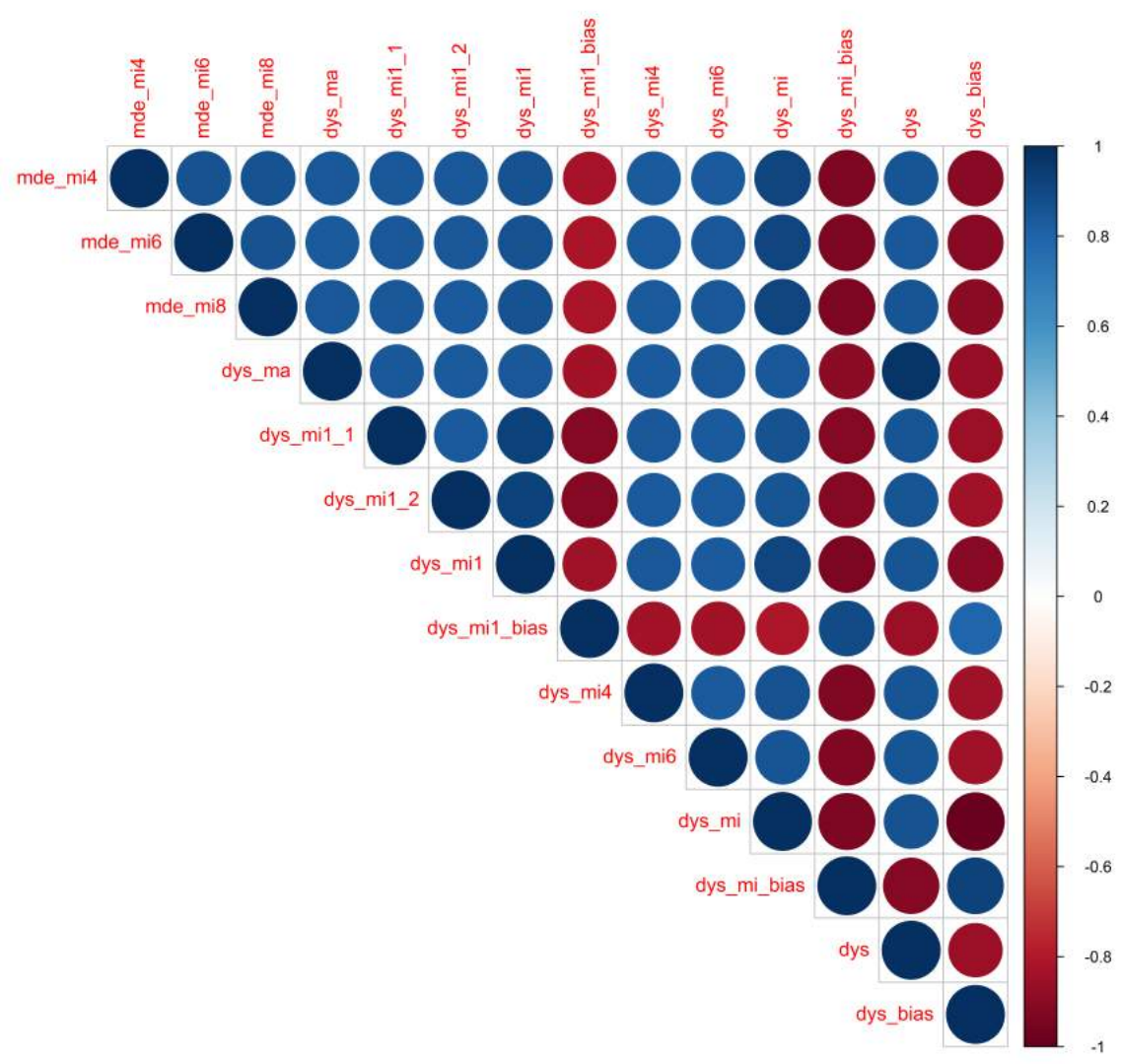


only

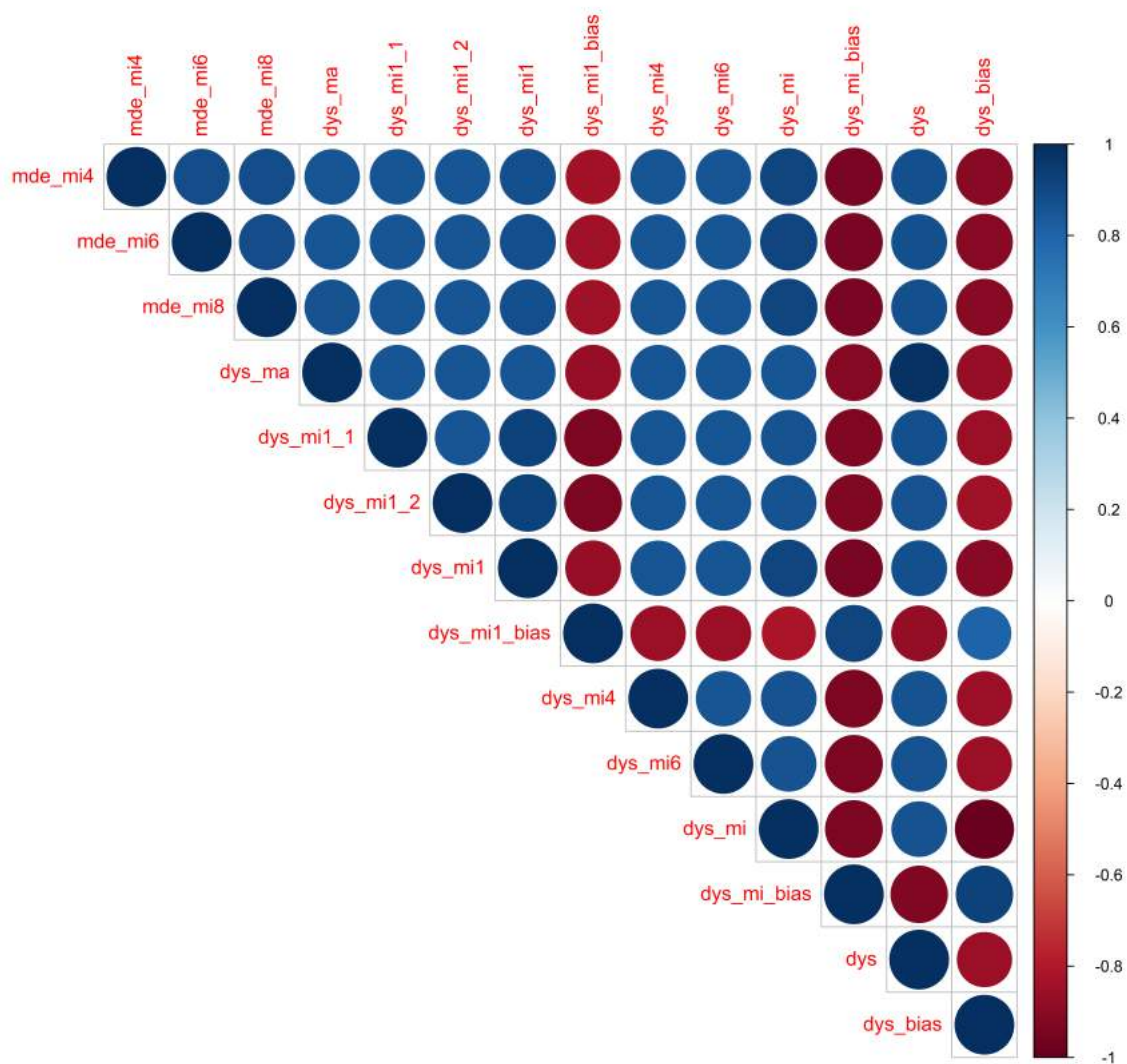


only

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

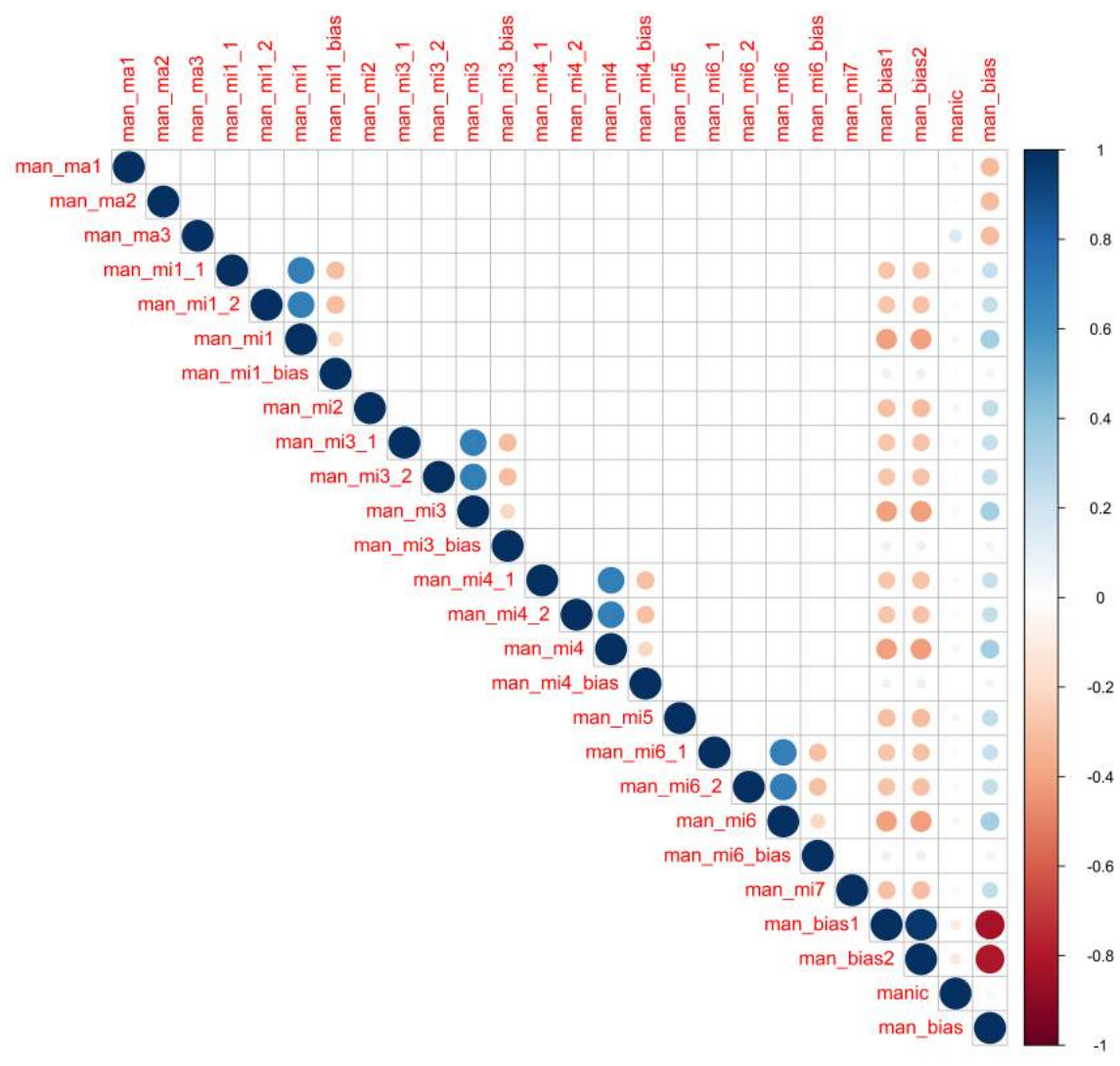


only



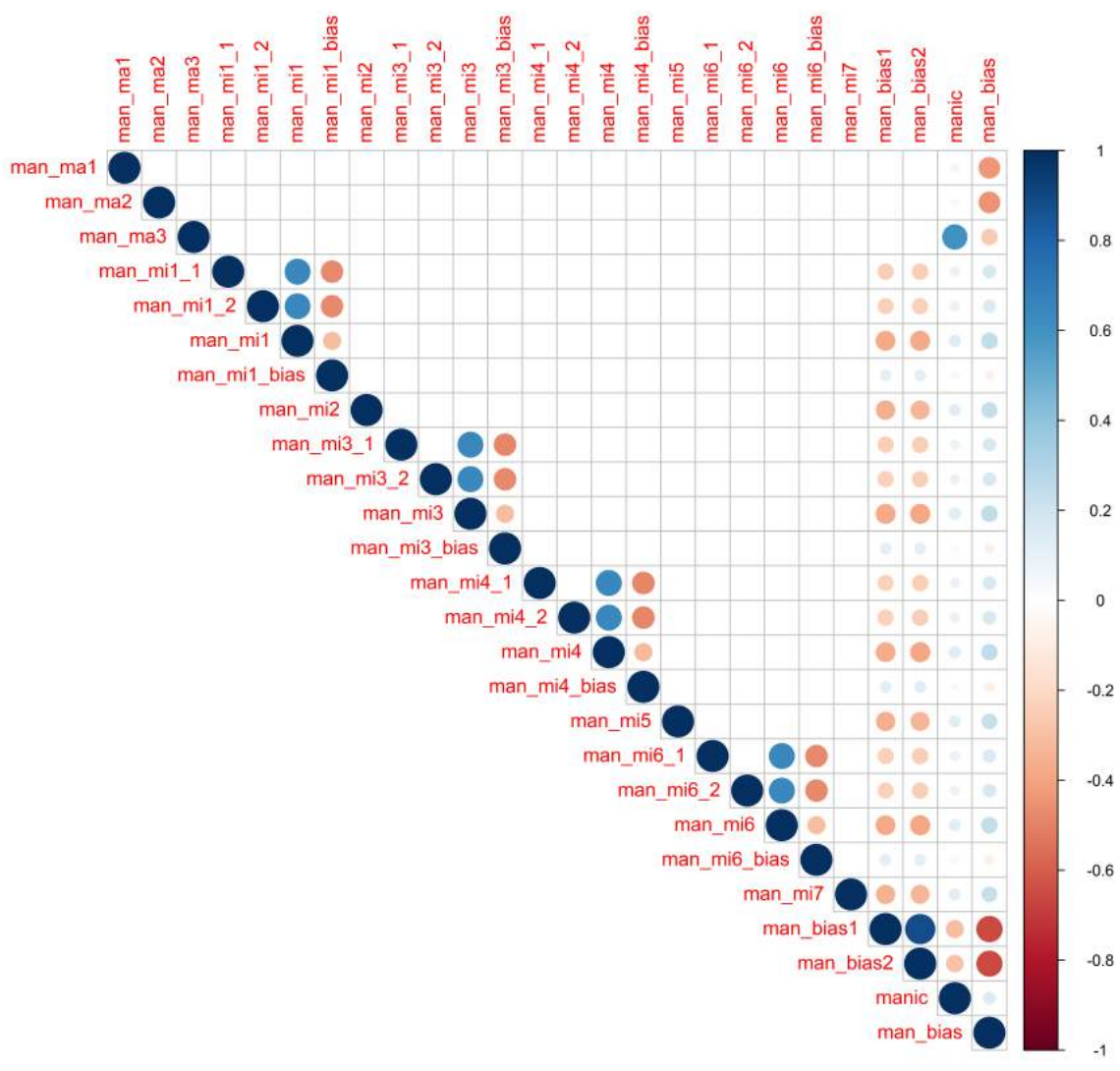
only

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



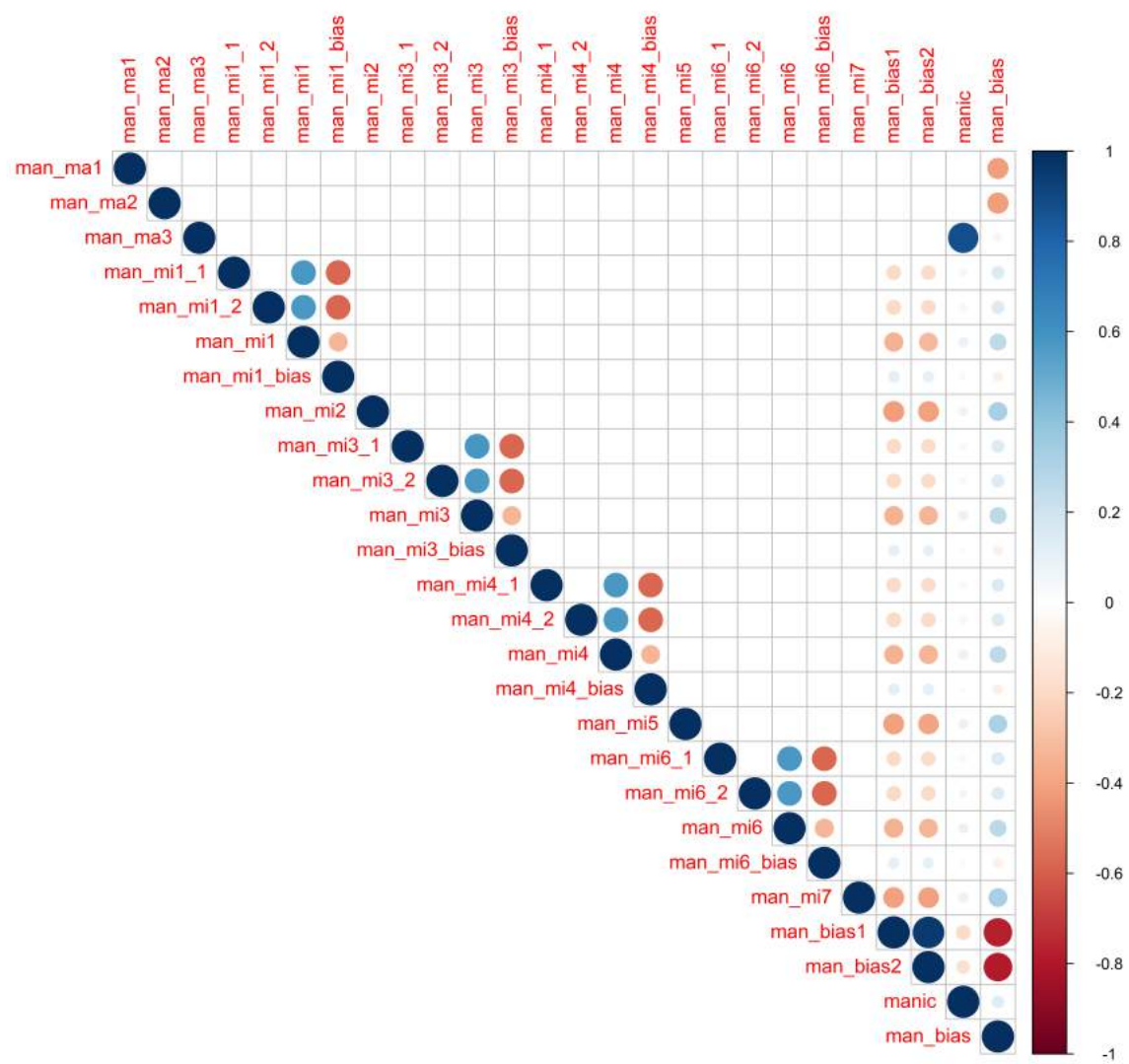
only

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

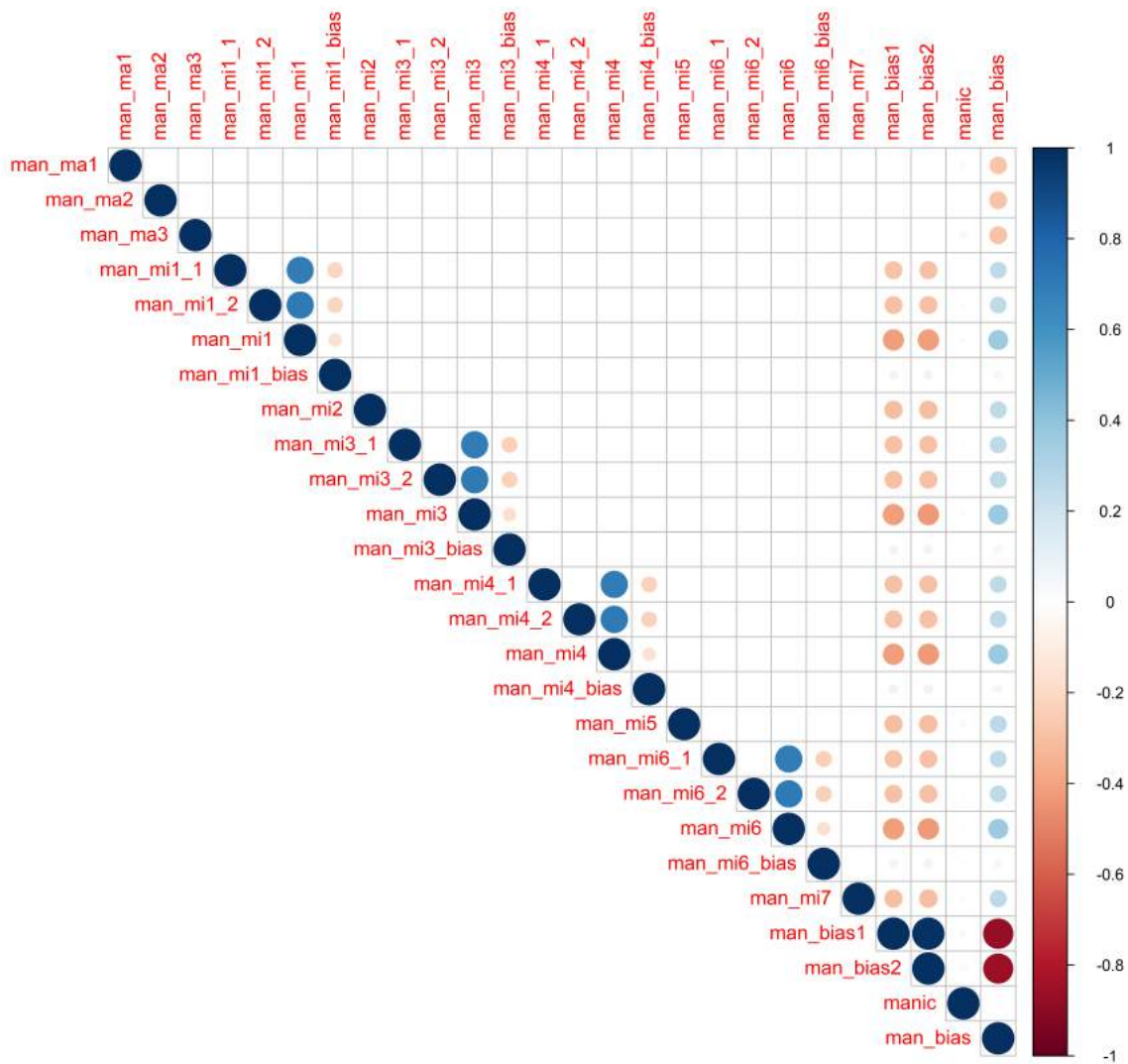


only

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

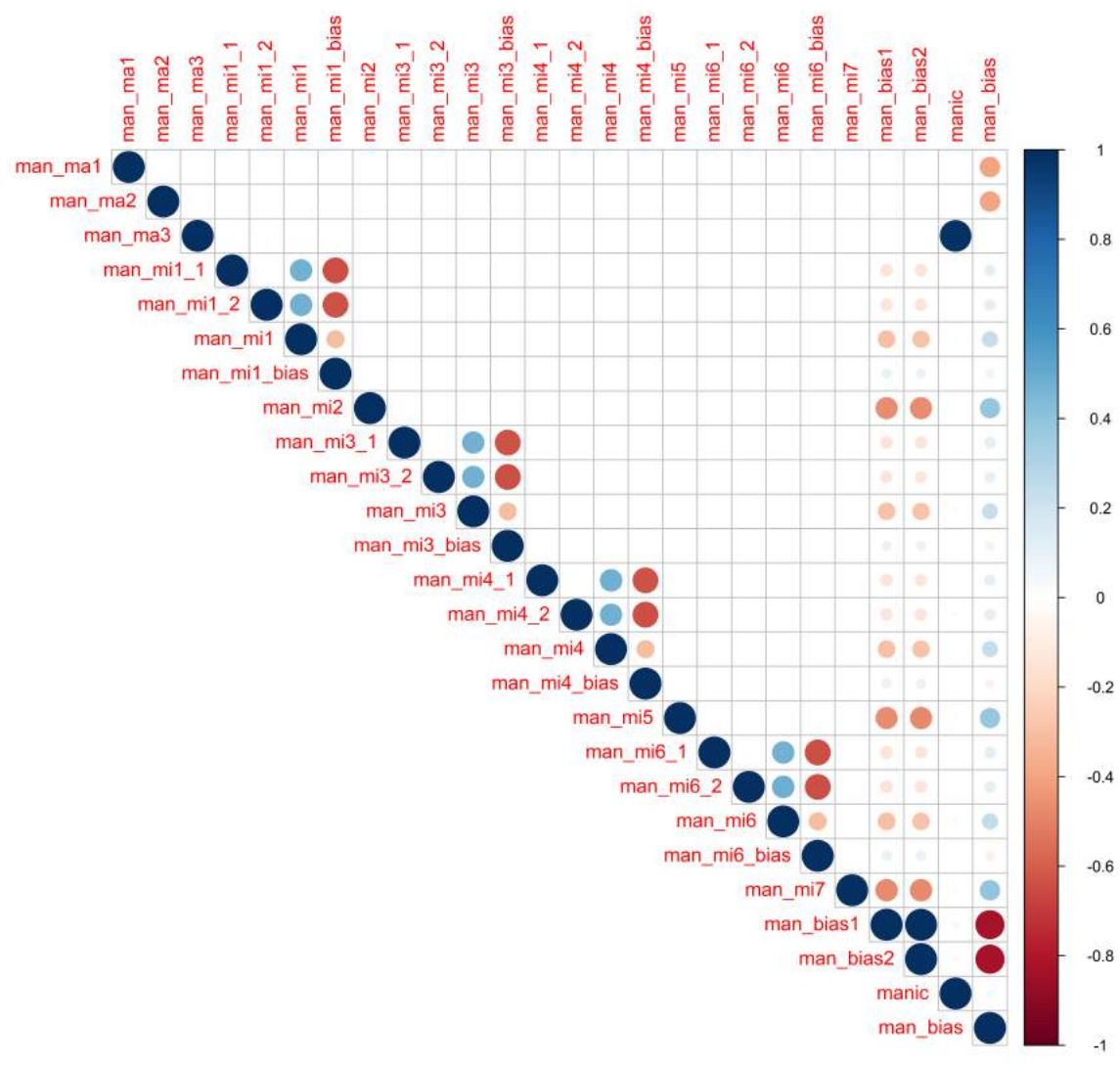


only

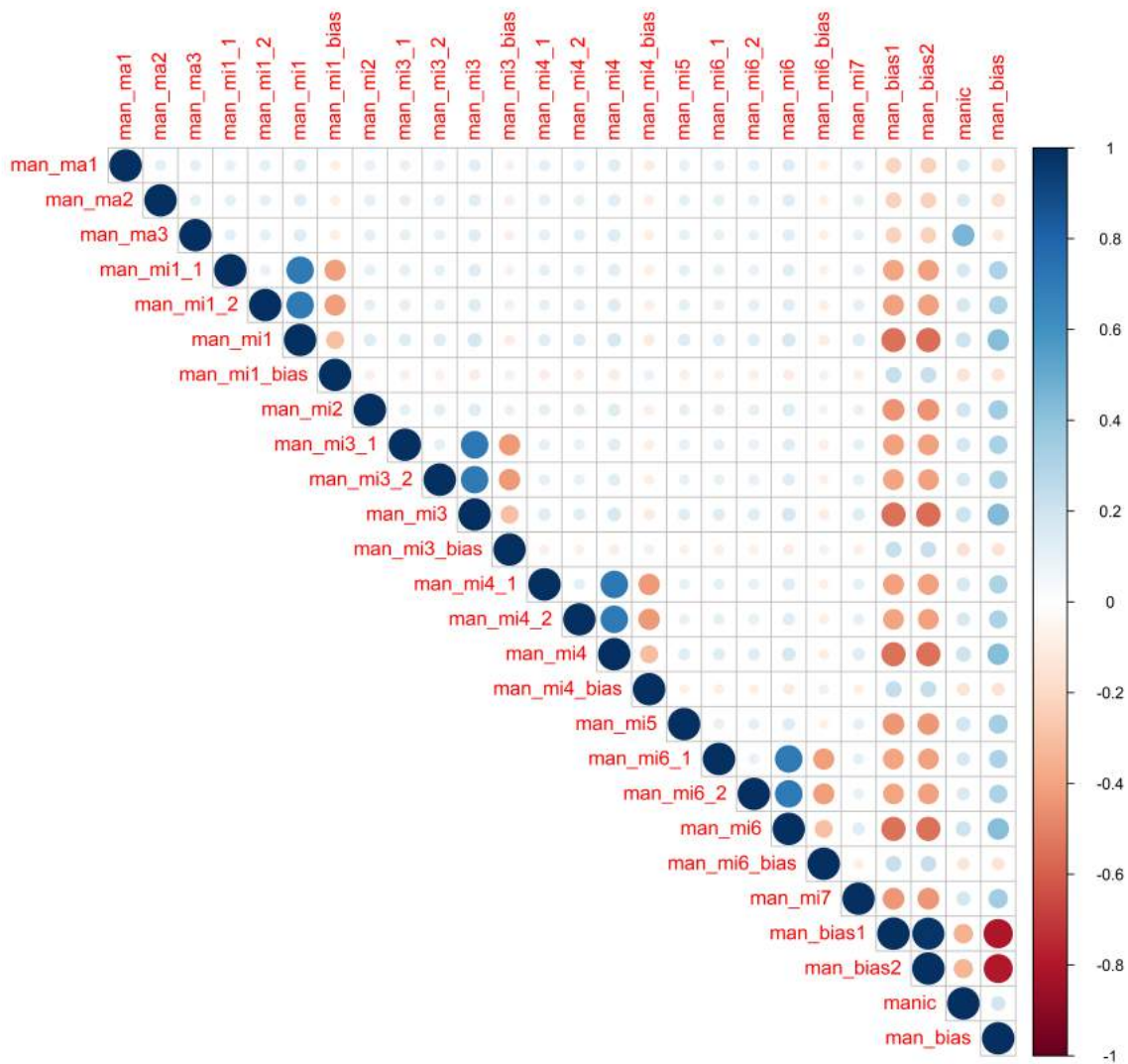


only

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

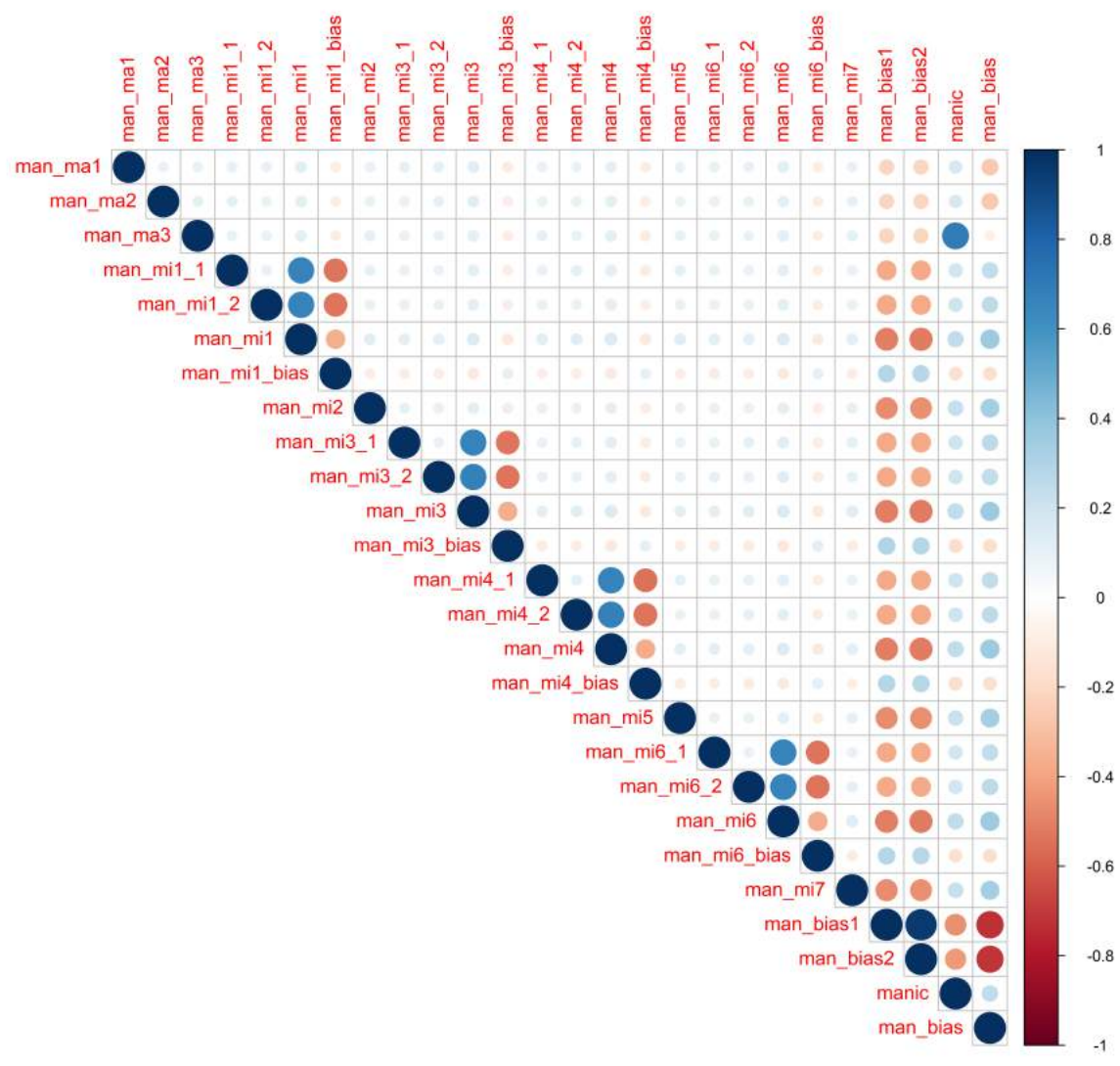


only

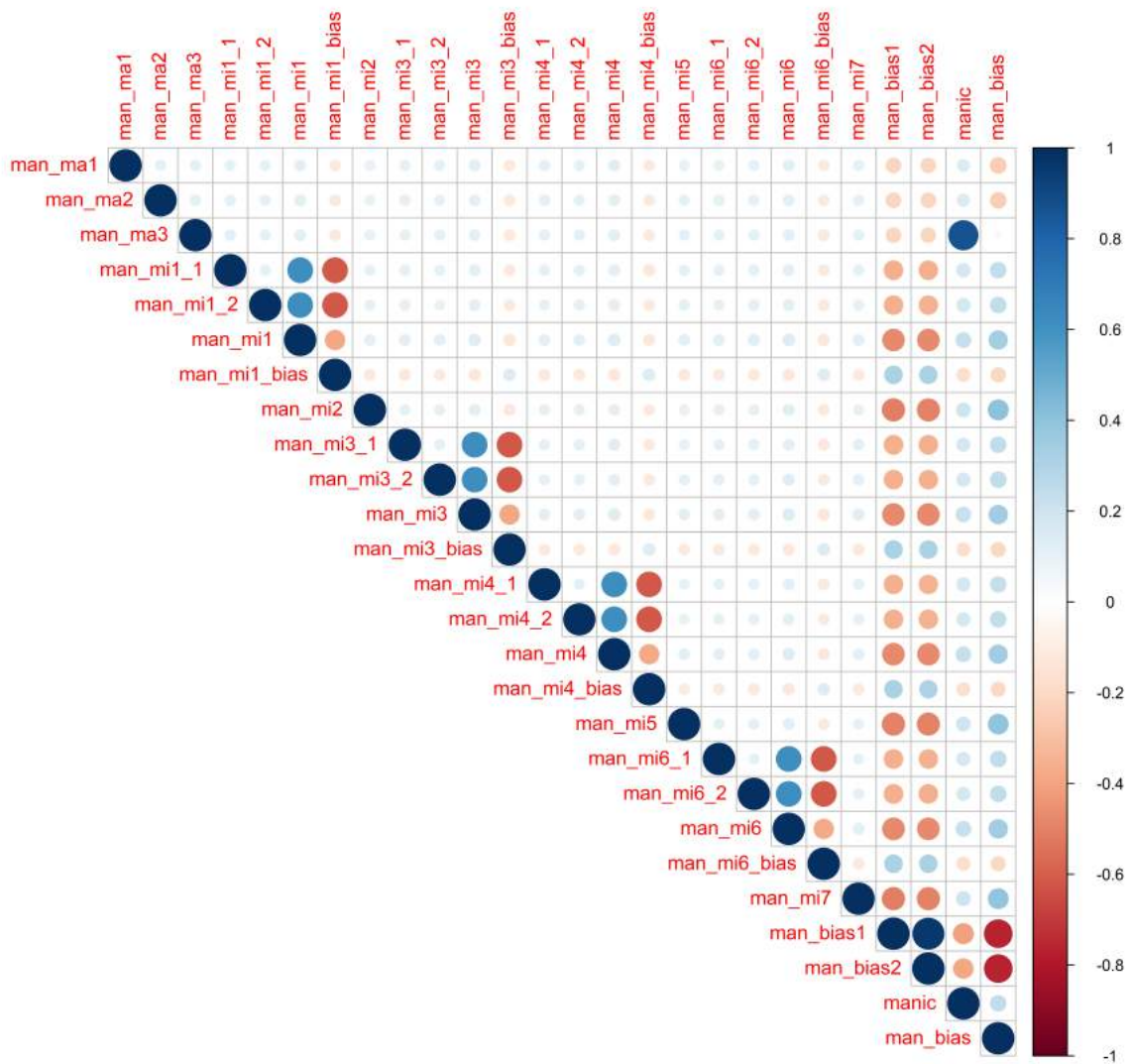


only

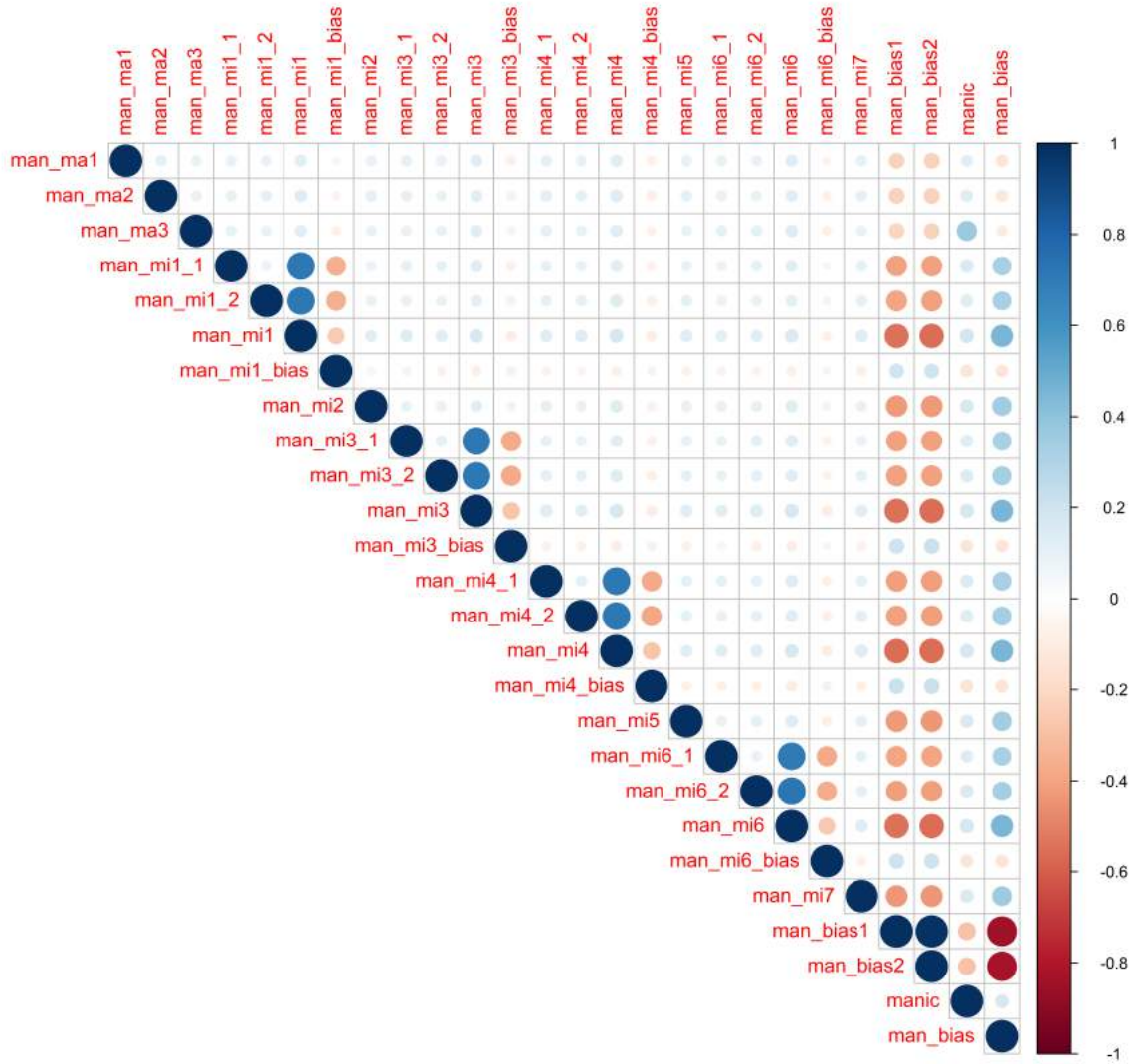
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



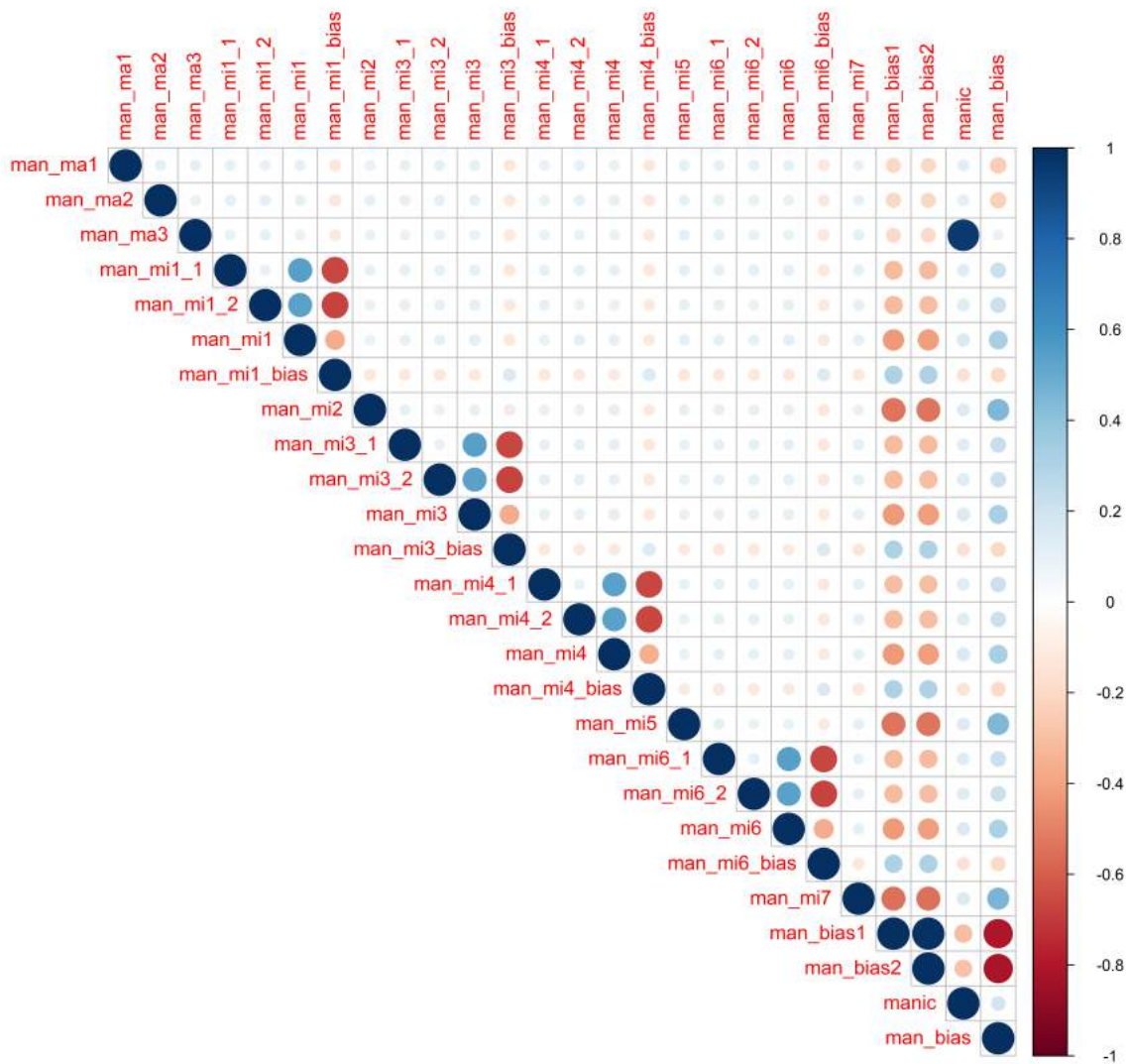
only



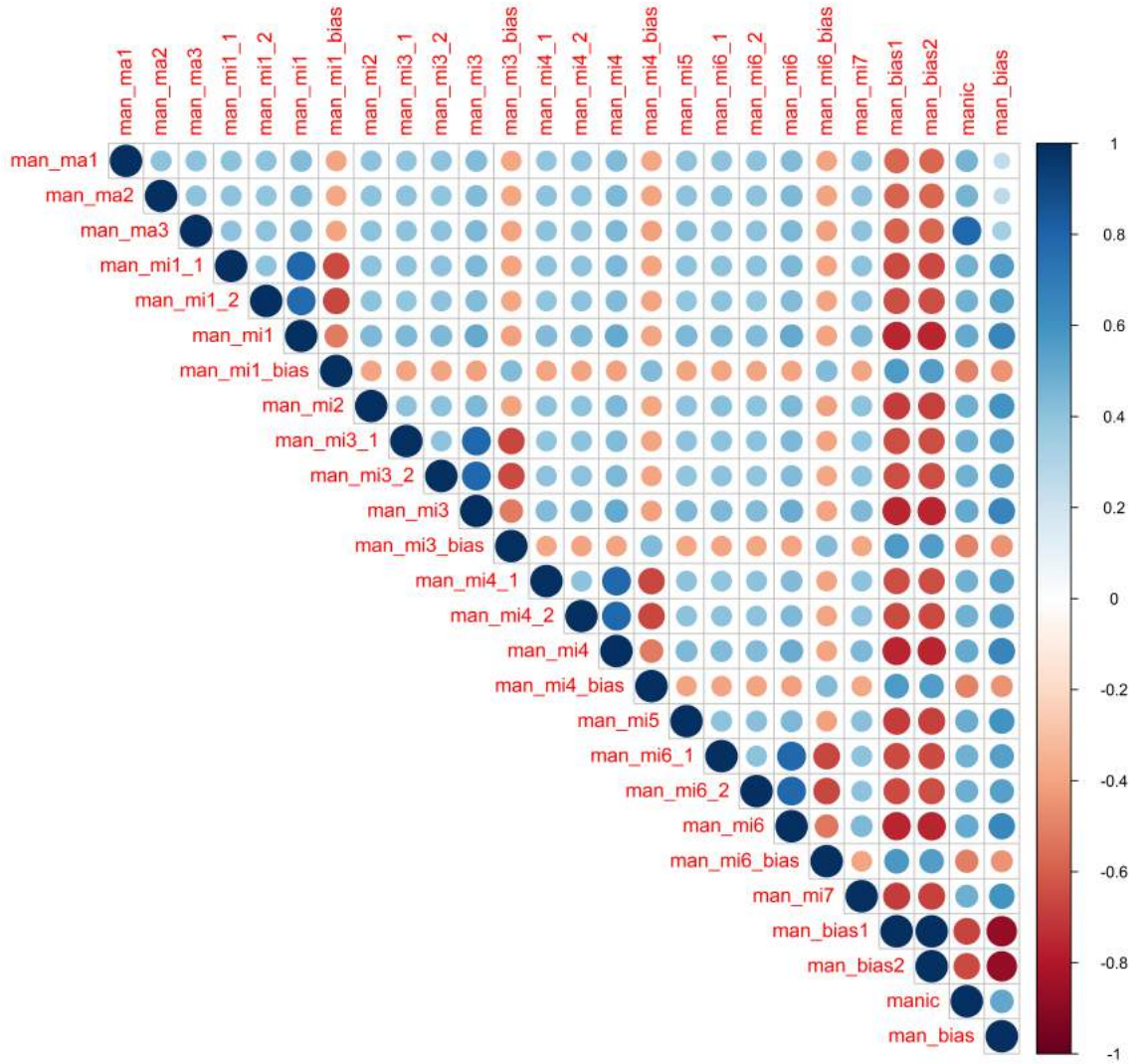
only



only

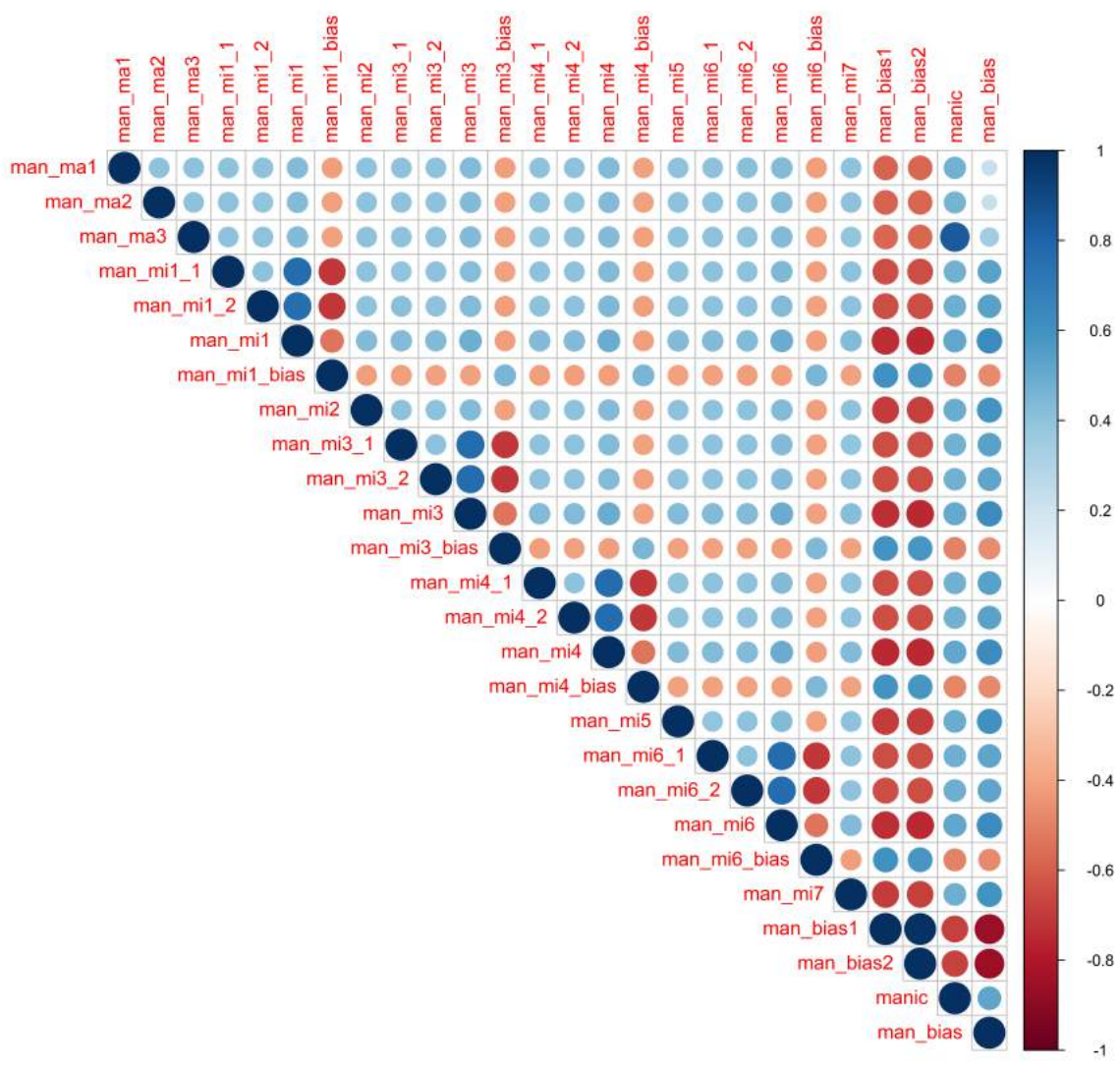


only



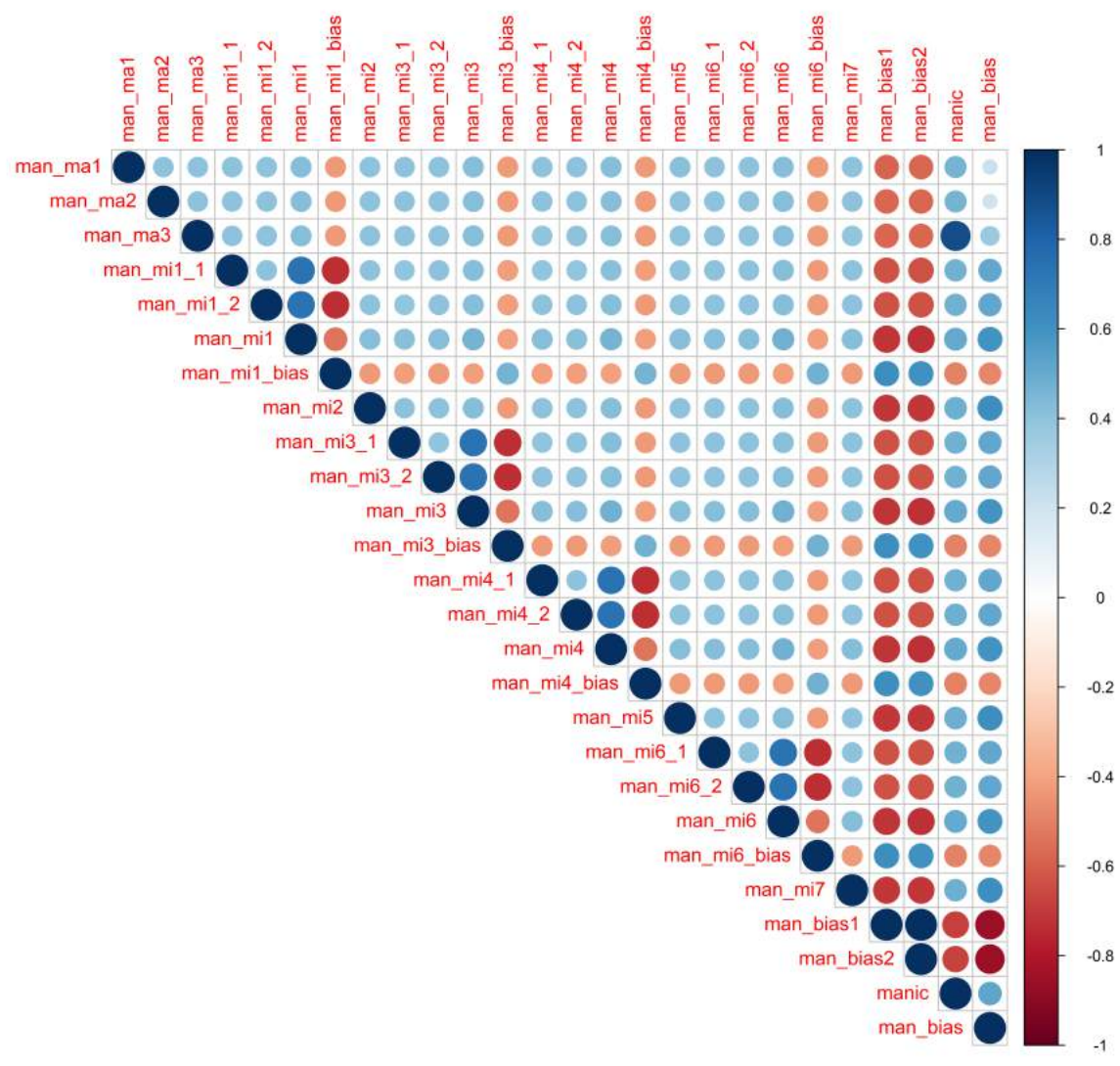
only

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

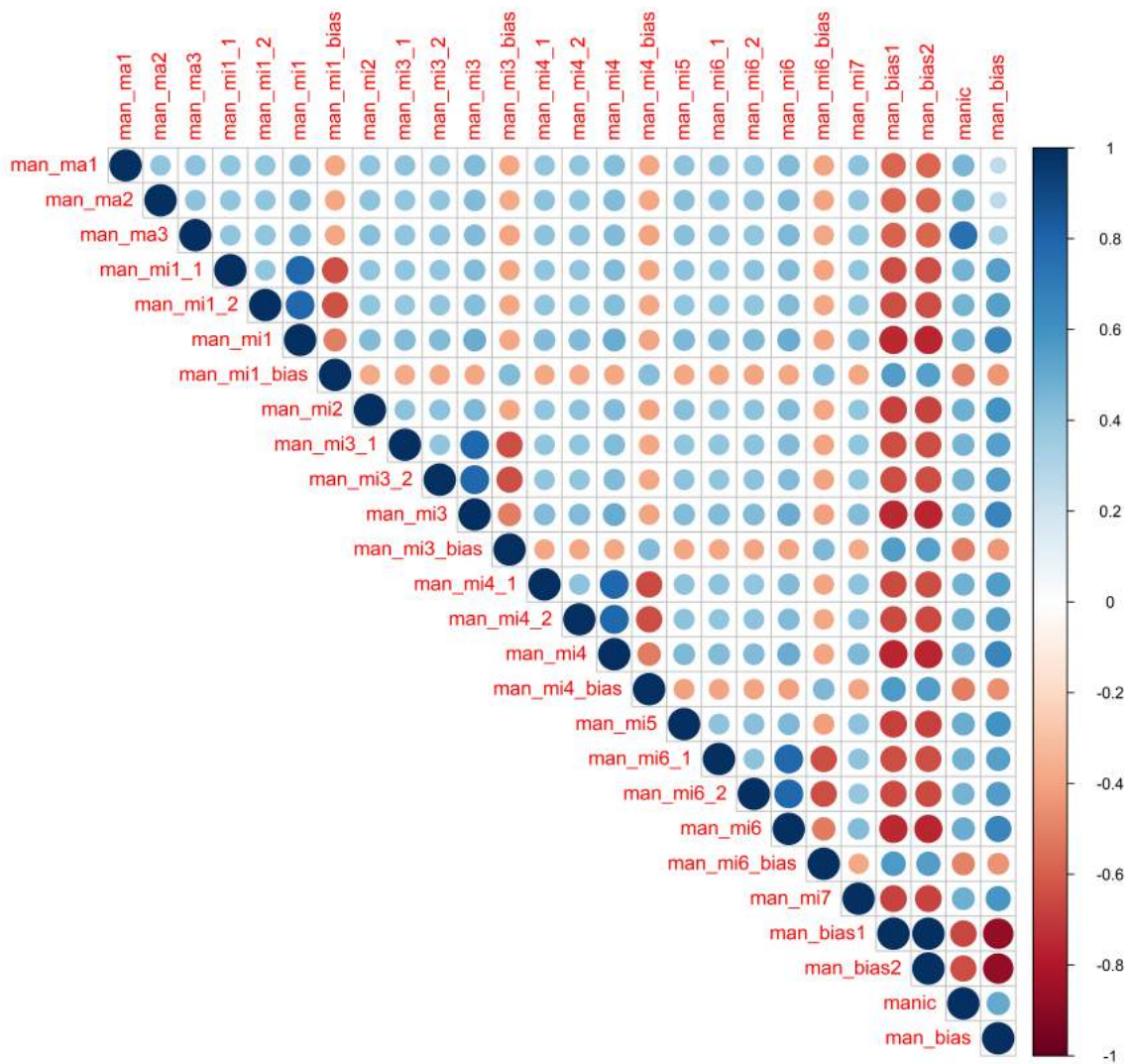


only

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

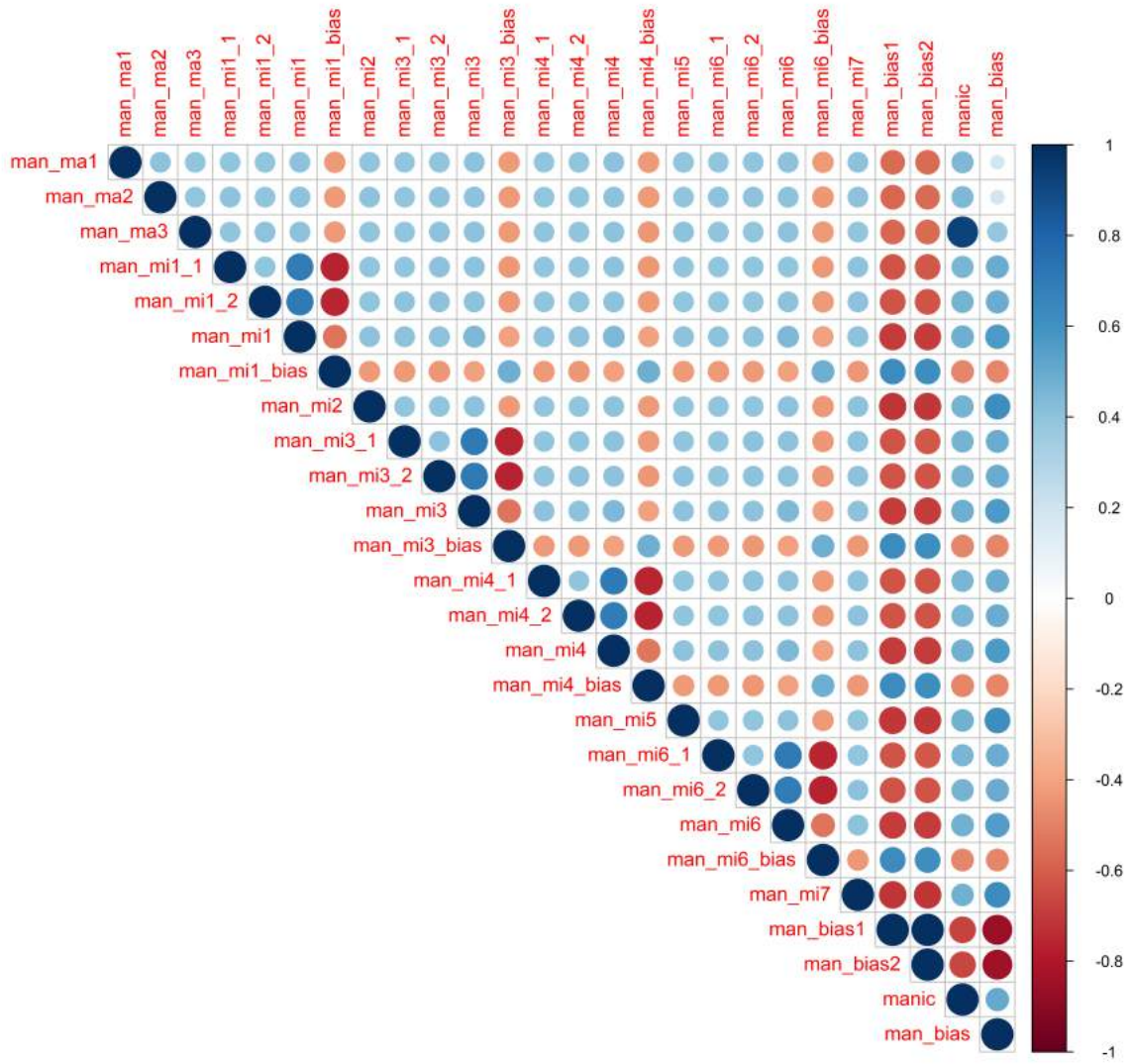


only



only

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



only