1  Evaluating predictive biomarkers for a binary outcome with linear versus logistic regression –

2  Practical recommendations for the choice of the model

3

4  Damian Gola[1,2], Nicole Heßler[1,2], Markus Schwaninger[2,3], Andreas Ziegler[1,2,4], Inke R.

5  König[1,2,4,5*]

6

7  [1]Institut für Medizinische Biometrie und Statistik, Universität zu Lübeck, Universitätsklinikum

8  Schleswig-Holstein, Campus Lübeck, Lübeck, Germany.

9  [2]German Centre for Cardiovascular Research (DZHK), partner site Hamburg/Kiel/Lübeck,

10  Lübeck, Germany.

11  [3]Institut für Experimentelle und Klinische Pharmakologie und Toxikologie, Universität zu

12  Lübeck, Lübeck, Germany

13  [4] School of Mathematics, Statistics and Computer Science, University of KwaZulu-Natal,

14  Pietermaritzburg, South Africa.

15  [5]Airway Research Center North (ARCN), Member of the German Center for Lung Research

16  (DZL).

17

18  E-Mail: inke.koenig@imbs.uni-luebeck.de (IRK)

19

# Abstract

A predictive biomarker can forecast whether a patient benefits from a specific treatment under study. To establish predictiveness of a biomarker, a statistical interaction between the biomarker status and the treatment group concerning the clinical outcome needs to be shown. In clinical trials looking at a binary outcome, linear or logistic regression models may be used to evaluate the interaction, but the effects in the two models are different and differently interpreted. Specifically, the effects are estimated as absolute risk reductions (ARRs) and odds ratios (ORs) in the linear and logistic model, thus measuring the effect on an additive and multiplicative scale, respectively.

We derived the relationship between the effects of the linear and the logistic regression model allowing for translations between the effect estimates between both models. In addition, we performed a comprehensive simulation study to compare the power of the two models under a variety of scenarios in different study designs. In general, the differences in power to detect interaction were minor, and visible differences were detected in rather unrealistic scenarios of effect size combinations and were usually in favor of the logistic model.

Based on our results and theoretical considerations, we recommend to 1) estimate logistic regression models because of their statistical properties, 2) test for interaction effects and 3) calculate and report both ARRs and ORs from these using the formulae provided.

## Introduction

41

42 Novel technologies and increased accumulated knowledge on the functional background of

43 diseases have made the application of biomarkers in clinical studies increasingly popular.

44 Their use is extremely diverse and includes serving as a tool for diagnosis, for staging the

45 disease, for forecasting disease prognosis or for monitoring and predicting clinical response

46 [1]. For many instances, it is most helpful to distinguish between prognostic biomarkers and

47 predictive biomarkers [2].

48 Prognostic biomarkers can forecast the development of the disease. In a randomized clinical

49 trial, this would usually be the outcome of the study such as remission. Importantly, this

50 forecast is independent of the intervention but an overall prognosis. Put differently, patients

51 with different prognostic biomarker profiles would have a different course of disease,

52 regardless of the intervention group. For example, the epidermal growth factor receptor

53 tyrosine kinase status is a prognostic factor for survival in patients with non-small cell lung

54 cancer [3], irrespective of the treatment. Predictive biomarkers, in contrast, predict the

55 effect of the intervention itself and therefore serve as companion diagnostic tests [4]. Thus,

56 patients with different predictive biomarker values would differ in how likely they are to

57 benefit from the specific therapy under study or to suffer from side effects. For instance,

58 several studies have shown that eosinophil counts in peripheral blood are predictors for

59 treatment response to Anti-IL-5 in patients with severe asthma [5-7].

60 Biomarkers are considered in clinical trials using different study designs, and these are

61 described in detail in the literature [2, 4, 8]. Which design should be used depends, among

62 other aspects, mostly on what is already known about the biomarker and the overall aim of

63 the study. If the aim is to prove the predictiveness of a biomarker, all patients regardless of

64 their biomarker status need to be randomized to the treatment groups. This is integrated in

65     the so-called "randomize-all" or "biomarker-stratified" design. Specifically, in the

66     "randomize-all" design, eligible patients are randomized into the treatment groups before

67     their biomarker status is assessed (Fig 1A). In the "biomarker-stratified randomization"

68     design, the biomarker status is assessed first. Then, patients with positive and negative

69     biomarker status are randomized separately (Fig 1B).

70

71     **Fig 1. Trial designs used in the simulation study.** (A) In the "randomize-all" design $n$ patients

72     are assigned irrespectively of their biomarker status to one the treatment groups based on

73     the randomization factor $\gamma$. (B) In the "biomarker-stratified randomization" design, $n$

74     patients are assigned to two randomizations based on their biomarker status.

75

76     If, in contrast, only patients with a positive biomarker status are randomized as in the

77     "targeted" design, it can only be shown that there is a treatment effect in this group, which

78     does not rule out that also biomarker negative patients benefit from the intervention, who

79     were not investigated. Furthermore, for establishing a predictive biomarker the trial needs

80     to show statistically that the treatment effect depends on the biomarker status, i.e., the

81     interaction between treatment arm and biomarker status has to be established. However, it

82     does not suffice to analyze biomarker positive and negative subgroups in separate trials and

83     report an effect in one but not the other group [9]: Firstly, not finding the therapeutic effect

84     in one group might be due to a lack of power. For example, in the study by Pant et al. [10]

85     predictiveness of albumin for the treatment of advanced pancreatic cancer with

86     bevacizumab was claimed on the finding of a positive effect in patients with normal baseline

87     albumin but not in others. However, only 26 patients with non-normal serum albumin levels

88     were included in the study. Hence, the confidence interval of the effect is very wide in this

4

89    subgroup and indeed includes the effect observed in patients with normal serum albumin.

90    Consequently, it cannot be ruled out that the effect was only not detected in the smaller

91    group, and no interaction between the treatment and albumin can be observed. A second

92    reason against claiming predictiveness based on the analysis of subgroups only is that even if

93    there are effects in both subgroups, predictiveness of the biomarker cannot be excluded,

94    because the therapeutic effect might be weaker (quantitative interaction) or in the opposite

95    direction (qualitative interaction) in the second subgroup.

96    In the following, we will describe the statistical methods to evaluate the biomarker-by-

97    treatment interaction that needs to be shown for the predictiveness of a biomarker.

## Statistical evaluation of biomarker-by-treatment interaction

99    The statistical method of choice to evaluate the biomarker-by-treatment interaction

100    depends on the data, i.e., the scale of the outcome variable and additional covariables that

101    are to be included in the model. In the following, we will focus on the simple setting of a

102    dichotomous outcome without further covariables. As a first approach, a linear regression

103    framework can be used in which the risk or probability of the dichotomous outcome $y$ (e.g.

104    therapy success) is modeled as a function of the dichotomous variables treatment $T$,

105    biomarker status $B$, and treatment-by-biomarker interaction $TB$ with

$$P(y = 1|T,B) = \beta_0 + \beta_T T + \beta_B B + \beta_{TB} TB.$$

107    Here, $T = 0$ or $T = 1$ if a patient receives the control treatment or the experimental

108    treatment, $B = 0$ or $B = 1$ if a patient is biomarker negative or positive, and $TB = 1$ only if a

109    biomarker positive patient receives the experimental treatment. Through this, the

110    coefficients $\beta_T$ and $\beta_B$ can be interpreted as the increase in risk with a change in the

111    treatment group and the biomarker status, respectively. The interpretation of these effect

5

112    estimates as absolute risk reductions (ARRs) is beneficial since it can be directly related to

113    the number needed to treat (NNT=1/ARR) [11]. The coefficient $\beta_{TB}$ indicates whether the

114    influence of *T* and *B* on *y* is independent, in which case it would equal 0. If it deviates from 0,

115    there is a statistical interaction between *T* and *B* regarding the risk of the outcome on the

116    additive scale [12].

117    However, this model has some statistical disadvantages. For example, the predicted

118    probability might be out of the range of possible values between 0 and 1. The standard

119    statistical model for analyzing dichotomous outcome in the life sciences therefore is the

120    logistic regression model. Here, the log odds of the outcome *y* is modeled as a function of *T*

121    and *B* and their interaction $TB$ by

122    $$\text{logit}(P(y=1|T,B)) = b_0 + b_T T + b_B B + b_{TB} TB .$$

123    From this, the coefficients $b_T$ and $b_B$ can be exponentiated to be interpreted as the increase

124    in odds of the outcome with a change in the treatment group and the biomarker status,

125    respectively. The coefficient $b_{TB}$, when exponentiated, then measures the treatment-by-

126    biomarker interaction as the odds ratio (OR) on the multiplicative scale. One advantage of

127    this model is that the predicted outcome probability will be guaranteed to lie between 0 and

128    1. Furthermore, the logit link is the natural parameter from the linear exponential family

129    which provides excellent statistical properties.

130    The linear and the logistic models are different, they have different effect sizes. This can be

131    seen from S1 Appendix in which we have derived the relation between ARRs from the linear

132    probability model and ORs from the logistic regression model.

133    Concerning the interaction effect, it can be shown that the models lead to different results,

134    meaning that the evidence for interaction will differ in strength, and that interaction in one

135    model does not imply interaction in the other. For example, in the study by Bokemeyer et al.

6

136    [13], patients with metastatic colorectal cancer had been randomized to receive FOLFOX-4

137    with or without cetuximab and were screened for *K-ras* mutations. A randomize-all design

138    was used, and, amongst other criteria, the best overall response in both *K-ras* positive and

139    negative patients was analyzed separately. We re-analyzed the data presented in the paper

140    and derived that the relative risk of response from a linear regression model under

141    cetuximab plus FOLFOX-4 versus FOLFOX-4 only was 1.68 in the wild type and 0.64 in the

142    mutation group, respectively. The corresponding p-value for the interaction was 0.00019. In

143    the logistic regression model, the odds ratio of response was 2.60 in the wild type and 0.46

144    in the mutation group, respectively, with an interaction p-value of 0.00023. Therefore, even

145    though interaction was established in both models, the p-values differ [13].

146    Therefore, given the statistical advantage of the logistic regression model over the linear

147    probability model, one may question the use of the linear regression model in this setting in

148    general. However, it has been shown that the statistical problems may not be as large as

149    anticipated [12, 14] and that, considering the interpretation of the effects, there are indeed

150    some merits to the linear model. As notional example, we consider the data in Table 1 (left),

151    showing the risk or probability of an outcome depending on the treatment and biomarker

152    status. In this example, changing the biomarker status from negative to positive always

153    increases the risk by 20%, and changing the treatment from control to experimental always

154    increases the risk by 40%. Thus, there is no additive biomarker by treatment interaction. We

155    now assume that we wish to select patients who will benefit most from the treatment. If

156    there were 100 patients each who were biomarker positive and negative, 10 and 30 would

157    reach a positive outcome, respectively, under control treatment (Fig 2A). Switching to the

158    experimental treatment instead, the numbers could be increased to 50 and 70, respectively.

159    This means that in either biomarker group, 20 patients would benefit from the experimental

7

160    treatment, indicating that the biomarker status does not need to be taken into account

161    when offering the treatment, which is mirrored by the lack of an additive interaction.

162    Consider now the data in Table 1 (right), where changing the biomarker status from negative

163    to positive increases the risk by 10% under control but by 30% under the experimental

164    therapy, and changing the treatment from control to experimental increases the risk by 20%

165    for biomarker negative and by 40% for biomarker positive patients. Phrased differently,

166    changing the biomarker status is always associated with doubling the risk, and changing the

167    therapy regimen with a 3-fold increase. In this case, there is therefore no multiplicative

168    interaction. Translating these risks into patient numbers who will benefit from the treatment

169    (Fig 2B) now shows that by switching the treatment from control to experimental would

170    benefit 20 biomarker-negative but 40 biomarker-positive patients. Given limited resources,

171    it might therefore be reasonable to offer the experimental treatment preferably to

172    biomarker positive patients, even though there is no biomarker by treatment interaction on

173    the multiplicative scale. From a health economic point of view, it can therefore be argued

174    that interaction on the additive scale, thus use of the linear regression model, should at least

175    complement the logistic regression model.

176

177    **Table 1. Notional risk of outcome.**

| Treatment | No additive interaction | | No multiplicative interaction | |
|-----------|-------|-------|-------|-------|
|           | B = 0 | B = 1 | B = 0 | B = 1 |
| T = 0     | 0.1   | 0.3   | 0.1   | 0.2   |
| T = 1     | 0.5   | 0.7   | 0.3   | 0.6   |

178    Notional risk of outcome in biomarker negative (B = 0) and biomarker positive (B = 1)
179    patients in the control (T = 0) and experimental treatment group (T = 1) in the scenario of no
180    additive interaction (left) and no multiplicative interaction (right).
181

8

182 **Fig 2. Number of patients with a positive outcome.** Based on a sample size of 100 in every

183 constellation in the scenario of (A) no additive interaction and (B) no multiplicative

184 interaction as specified in Table 1. Solid line: biomarker negative, dashed line: biomarker

185 positive.

186

187 Given that interactions on both scales can occur, are relevant and should be analyzed, we

188 need to know how powerful the statistical analyses will be. More specifically, if there is an

189 additive interaction, how likely will this be detected using the "false" model, i.e., the logistic

190 regression? Vice versa, how likely is it to detect a multiplicative interaction when using the

191 linear regression? To answer these questions, we performed a simulation study that will be

192 described in the following.

# Methods

193

## Simulation framework

194

195 In our simulation we start from a population with individuals affected and unaffected by the

196 disease under study, which is indicated by the disease status $D \in \{1, 0\}$. Additional to the

197 general probability of developing the disease, the probability might be influenced by having

198 or having not a certain biomarker status $B \in \{1,0\}$. A random sample $R$ of the diseased

199 individuals is recruited to a clinical trial, comparing an experimental treatment with the

200 control treatment, denoted by $T \in \{1,0\}$. The trial aims to answer the research question

201 whether the biomarker $B$ is predictive, i.e., whether it modifies the probability of treatment

202 success $y \in \{1,0\}$.

## Population simulation

203

9

204    We define the prevalence of a dichotomous biomarker $B$ by $P(B = 1) = \phi$. Populations are

205    simulated by modelling the disease probability by

$$\text{logit}(P(D = 1 \mid B)) = b_0^D + b_B^D B \qquad (1)$$

206    and sampling the disease status $D$ from a Bernoulli distribution with probability $P($

207    $(D = 1 \mid B))$. Here, $b_0^D$ is the baseline $\log(odds)$ of the disease and $b_B^D$ is a prognostic effect

208    of the biomarker $B$.

## Trial designs

210    As illustrated in Fig 1, in the "randomize-all" design $n$ patients are drawn randomly from a

211    simulated population. Based on the randomization factor $\gamma \in (0,1)$, $\gamma n$ randomly chosen

212    patients receive the biomarker guided treatment ($T = 1$) and $(1 - \gamma)n$ randomly chosen

213    patients receive the control treatment ($T = 0$). After the assignment to a treatment arm the

214    biomarker status is revealed. Thus, the numbers of biomarker positive ($n_+$) and biomarker

215    negative ($n_-$) patients in each treatment group are determined by the biomarker

216    prevalence $\phi$. In the "biomarker-stratified randomization" design the biomarker status is

217    revealed before randomization. This enables to draw $n_+$ biomarker positive and $n_-$

218    biomarker negative, $n = n_+ + n_-$ in total, patients from a simulated population. By

219    specifying $n_-$ and $n_+$, the prevalence of the biomarker under consideration is not reflected

220    in this design. In each biomarker stratum, the randomization factors $\gamma_+ \in (0,1)$ and $\gamma_-$

221    $\in (0,1)$ determine the proportion of patients receiving control or biomarker guided

222    treatment.

## Data simulation

224     In the present simulation study, treatment success is simulated on both the linear and

225     logistic scale in both trial designs for varying parameters. The procedure to simulate this

226     data is as follows:

     1. Draw $n$ patients from a population based on formula ( 1 ).

227

     2. Assign patients to treatment arms based on $\gamma$ or $\gamma_+$ and $\gamma_-$, depending on the trial

228

229         design.

230      3. Calculate the treatment success probability $P(y = 1)$ by applying either

$$P(y = 1 \mid T, B) = \text{expit}(b_0 + b_T T + b_B B + b_{TB} TB) \qquad (2)$$

231         or

$$P(y = 1 \mid T, B) = \mu + \beta_T T + \beta_B B + \beta_{TB} TB \qquad (3)$$

232         for every patient with $\text{expit}(c) = \frac{\exp(c)}{1 + \exp(c)}$, and $T$ and $B$ denote the treatment and

233         biomarker status, respectively.

234      4. Sample the treatment success from a Bernoulli distribution using the probability from

235         formula ( 2 ) or ( 3 ).

236     We consider $\phi \in \{0.1, 0.25, 0.5\}$ as prevalence for the biomarker, and we use $b_0^D = 0$ and $b_B^D$

237     $= 0$ to simulate populations, i.e., there is no prognostic effect of the biomarker. We create

238     study populations of sizes $n \in \{100, 200, 500, 1000\}$. In case of the "biomarker-stratified

239     randomization" trial either half of the study population is biomarker positive and the other

240     half is biomarker negative; alternatively, the proportion of biomarker positive patients is

241     determined by the biomarker prevalence in the respective simulated population, i.e.

242     specifying $n_+ = n_- = \frac{n}{2}$ explicitly or specifying only $n$, and from this follows $n_+ \approx \phi n$. We

243     use $\gamma, \gamma_+, \gamma_- \in \{0.25, 0.5, 0.75\}$ as randomization factors, and in the "biomarker-stratified

244     randomization" trial all combinations of the values of $\gamma_+$ and $\gamma_-$ are considered. The effect

245      sizes to determine the treatment success probability are the cross-product of a range of

246      possible values. On the linear scale we use

247         •   $\beta_0 = 0.5$,

248         •   $\beta_T \in \{0, 0.1, 0.2, 0.3, 0.4\}$,

249         •   $\beta_B \in \{-0.4, -0.3, -0.2, -0.1, 0, 0.1, 0.2, 0.3, 0.4\}$ and

250         •   $\beta_{TB} \in \{-0.4, -0.3, -0.2, -0.1, 0, 0.1, 0.2, 0.3, 0.4\}$.

251      Combinations of effect sizes leading to a probability of therapy success less than 0 or greater

252      than 1 are excluded, e.g. $\beta_0 = 0.5$, $\beta_T = 0$, $\beta_B = -0.4$, $\beta_{TB} = -0.4$ is not valid.

253      On the logistic scale we use

254         •   $b_0 = 0$,

255         •   $b_T \in \{0, 0.2231, 0.4055, 0.5596, 0.6931\}$ corresponding to OR

256           $\in \{1, 1.25, 1.50, 1.75, 2\}$,

257         •   $b_B \in \{-0.6931, -0.5596, -0.4055, -0.2231, 0, 0.2231, 0.4055, 0.5596, 0.6931\}$
258           corresponding to $OR \in \{0.5, 0.5713, 0.6667, 0.8, 1, 1.25, 1.5, 1.75, 2\}$

259         •   $b_{TB}$
260           $\in \{-0.6931, -0.5596, -0.4055, -0.2231, 0, 0.2231, 0.4055, 0.5596, 0.6931\}$

261           corresponding to $OR \in \{0.5, 0.5713, 0.6667, 0.8, 1, 1.25, 1.5, 1.75, 2\}$.

262      In total, we use 680 unique effect size combinations for our simulations. Note that effect size

263      combinations having $\beta_{TB} = 0$ or $b_{TB} = 0$ act as null models for the respective regression

264      model analysis.


265      ## Analyses

266      All simulated data sets are analyzed using both linear and logistic models. Following Kraft et

267      al. [15], the likelihood ratio-based deviance test between the saturated model

$$\text{logit}(\hat{\pi}) = \hat{b}_0 + \hat{b}_T T + \hat{b}_B B + \hat{b}_{TB} TB \qquad (4)$$

268   or

$$\hat{\pi} = \hat{\mu} + \hat{\beta}_T T + \hat{\beta}_B B + \hat{\beta}_{TB} TB, \qquad (5)$$

269   where $\pi = P(D = 1 \,|\, B)$, and a model considering both main effects of treatment and

270   biomarker but no interaction effect (restricted deviance test) is calculated. In addition, a

271   Wald-like test on the null hypotheses $H_0{:}b_{TB} = 0$ (logistic regression model) or $H_0{:}\beta_{TB} = 0$

272   (linear regression model) in the respective saturated models ( 4 ) and ( 5 ) is performed. To

273   obtain reliable estimates for the power to detect an interaction between treatment and

274   biomarker effect, 1000 replicates are run. For each replicate it is noted whether the two-

275   sided p-value of the respective test is less than $\alpha = 0.05$.

276   All simulations and analyses are done in R 3.3.1 [16] utilizing the R package batchtools [17].

277   The code is available in the supplement (S2 Appendix).

## Results

279   Table 2 shows the estimated frequency of type I errors of the interaction test, i.e., the

280   restricted deviance test, in logistic and linear regression models to detect a interaction effect

281   simulated via the linear (upper part) or logistic (lower part) model. Given are the frequencies

282   in the "randomize-all" trial design with biomarker prevalence $\phi = 0.1$ and randomization

283   factor $\gamma = 0.5$ for some selected effect size combinations with no ($b_{TB} = \log (1)$ and $\beta_{TB}$

284    $= 0$), moderate ($b_{TB} = \log (1.5)$ or $b_{TB} = \log (\frac{2}{3})$ and $\beta = \pm\, 0.2$) and strong ($b_{TB} = \log (0.5)$

285   or $b_{TB} = \log (2)$ and $\beta = \pm\, 0.4$) effects. The effect sizes are given on both the linear and

286   logistic scale for sample sizes $n = 200$ and $n = 500$, sorted by the biomarker main effects

287   (Table 2). Other scenarios meeting these restrictions but not displayed are redundant such

288   that the effects $\beta_T, \beta_B, b_T$ or $b_B$ have opposite signs or are permuted.

289

**Table 2. Estimated type I error frequency at the nominal two-sided 0.05 test-level in the "randomize-all" design.**

| | | | | | | | n = 200 | | n = 500 | |
|---|---|---|---|---|---|---|---|---|---|---|
| Scen | $\beta_T$ | $\beta_B$ | $\beta_{TB}$ | $b_T$ | $b_B$ | $b_{TB}$ | logistic | linear | logistic | linear |
| 1 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.060 | 0.046 | 0.057 | 0.053 |
| 2 | 0.0000 | -0.1000 | 0.0000 | 0.0000 | -0.4055 | 0.0000 | 0.054 | 0.040 | 0.053 | 0.046 |
| 3 | 0.0000 | -0.2000 | 0.0000 | 0.0000 | -0.8473 | 0.0000 | 0.064 | 0.026 | 0.058 | 0.038 |
| 4 | 0.0000 | -0.4000 | 0.0000 | 0.0000 | -2.1972 | 0.0000 | 0.043 | 0.004 | 0.086 | 0.006 |
| 5 | 0.4000 | -0.4000 | 0.0000 | 2.1972 | -2.1972 | 0.0000 | 0.039 | 0.045 | 0.068 | 0.051 |
| 6 | 0.0000 | -0.1667 | 0.0000 | 0.0000 | -0.6931 | 0.0000 | 0.062 | 0.029 | 0.048 | 0.035 |
| 7 | 0.1667 | -0.1667 | 0.0000 | 0.6931 | -0.6931 | 0.0000 | 0.062 | 0.044 | 0.052 | 0.050 |

292 Frequency estimates are based on the likelihood ratio-based restricted deviance test in the

293 "randomize-all" trial design with biomarker prevalence $\phi = 0.1$ and randomization factor

294 $\gamma = 0.5$. $\beta_0 = 0.5$ and $b_0 = 0$. Scen = Number of scenario with respective effect size

295 combination $\beta_T$, $\beta_B$, $\beta_{TB}$ or $b_T$, $b_B$, $b_{TB}$. Logistic and linear refer to the type I error frequency in

296 the logistic and linear regression model, respectively.

297

298 Table 2 shows that the frequency of type I errors for the restricted deviance test in both

299 regression models mainly is near to 0.05, as expected, and thus in line with the specified

300 significance level of α = 0.05. However, in some scenarios the linear and logistic model

301 deviate from the specified significance level. Based on Bradley's liberal criterion of

302 robustness [18], the type I error frequency should be between 0.025 and 0.075. Both the

303 logistic and the linear model fail to fall into this range in scenario 4, which is characterized by

304 a single strong main effect. The total number and percentage of scenarios violating Bradley's

14

305   criterion in the "randomize-all" design is shown in Table 3. In total, 54 times (5% of all

306   scenarios) the logistic model has a type I error outside Bradley's bounds, whereas the linear

307   model violates this criterion 123 times (11% of all scenarios). Comparing the numbers per

308   model and criterion bound, it is of special interest that the logistic model tends to violate the

309   upper bound (liberal) whereas the linear model tends to violate the lower bound

310   (conservative).

311

312   **Table 3. Number of scenarios in which type I error frequencies deviate from Bradley's**

313   **criterion [18] in the "randomize-all" design.**

|          |           | n = 100   | n = 200   | n = 500   | n = 1000  | Σ         |
|----------|-----------|-----------|-----------|-----------|-----------|-----------|
|          | > 0.075   | 23 (8%)   | 22 (8%)   | 5 (2%)    | 2 (1%)    | 52 (5%)   |
| logistic | < 0.025   | 0 (0%)    | 2 (1%)    | 0 (0%)    | 0 (0%)    | 0 (0%)    |
|          | Σ         | 23 (8%)   | 24 (8%)   | 5 (2%)    | 2 (1%)    | 54 (5%)   |
|          | > 0.075   | 5 (2%)    | 6 (2%)    | 5 (2%)    | 5 (2%)    | 21 (2%)   |
| linear   | < 0.025   | 32 (11%)  | 25 (9%)   | 23 (8%)   | 22 (8%)   | 102 (9%)  |
|          | Σ         | 37 (13%)  | 31 (11%)  | 28 (10%)  | 27 (9%)   | 123 (11%) |

314   Based on the likelihood ratio-based restricted deviance test in the "biomarker-stratified"

315   trial design. All 1152 scenarios with $\beta_{TB} = b_{TB} = 0$ are considered.

316

317   We next look at the power of the restricted deviance test to detect an interaction effect

318   simulated via the linear (Table 4, upper part) or logistic (Table 4, lower part) model in the

319   same setting, i.e., the "randomize-all" trial design with the same effect specifications as

320   before. Results are sorted by the interaction effects.

321

**Table 4. Estimated power at the nominal two-sided 0.05 test-level in the "randomize-all"**

**design.**

| | | | | | | | n = 200 | | n = 500 | |
|---|---|---|---|---|---|---|---|---|---|---|
| Scen | $\beta_T$ | $\beta_B$ | $\beta_{TB}$ | $b_T$ | $b_B$ | $b_{TB}$ | logistic | linear | logistic | linear |
| 8 | 0.2000 | -0.4000 | 0.0000 | 0.8473 | -2.1972 | 0.5026 | 0.077 | 0.015 | 0.108 | 0.017 |
| 9 | 0.0000 | -0.1000 | 0.1000 | 0.0000 | -0.4055 | 0.4055 | 0.084 | 0.065 | 0.107 | 0.102 |
| 10 | 0.0000 | 0.0000 | -0.1000 | 0.0000 | 0.0000 | -0.4055 | 0.080 | 0.064 | 0.105 | 0.100 |
| 11 | 0.1000 | -0.1000 | -0.1000 | 0.4055 | -0.4055 | -0.4055 | 0.086 | 0.072 | 0.105 | 0.100 |
| 12 | 0.1000 | -0.1667 | -0.1000 | 0.4055 | -0.6931 | -0.4055 | 0.085 | 0.062 | 0.103 | 0.097 |
| 13 | 0.1667 | -0.1667 | -0.1000 | 0.6931 | -0.6931 | -0.4055 | 0.076 | 0.059 | 0.113 | 0.115 |
| 14 | 0.0000 | -0.4000 | 0.2000 | 0.0000 | -2.1972 | 1.3499 | 0.218 | 0.084 | 0.423 | 0.239 |
| 15 | 0.0000 | 0.0000 | -0.2000 | 0.0000 | 0.0000 | -0.8473 | 0.144 | 0.113 | 0.282 | 0.258 |
| 16 | 0.0000 | -0.2000 | -0.2000 | 0.0000 | -0.8473 | -1.3499 | 0.204 | 0.071 | 0.436 | 0.227 |
| 17 | 0.2000 | -0.4000 | -0.2000 | 0.8473 | -2.1972 | -0.8473 | 0.077 | 0.054 | 0.156 | 0.223 |
| 18 | 0.4000 | -0.4000 | -0.2000 | 2.1972 | -2.1972 | -0.8473 | 0.088 | 0.148 | 0.160 | 0.376 |
| 19 | 0.0000 | -0.2000 | 0.4000 | 0.0000 | -0.8473 | 1.6946 | 0.437 | 0.401 | 0.770 | 0.764 |
| 20 | 0.0000 | -0.4000 | 0.4000 | 0.0000 | -2.1972 | 2.1972 | 0.556 | 0.404 | 0.881 | 0.799 |
| 21 | 0.2000 | -0.4000 | 0.4000 | 0.8473 | -2.1972 | 2.1972 | 0.513 | 0.396 | 0.876 | 0.827 |
| 22 | 0.1000 | 0.1000 | -0.0077 | 0.4055 | 0.4055 | 0.0000 | 0.067 | 0.038 | 0.052 | 0.040 |
| 23 | 0.1000 | 0.1667 | -0.0167 | 0.4055 | 0.6931 | 0.0000 | 0.070 | 0.040 | 0.063 | 0.042 |
| 24 | 0.1667 | 0.1667 | -0.0333 | 0.6931 | 0.6931 | 0.0000 | 0.063 | 0.035 | 0.059 | 0.036 |
| 25 | 0.1000 | -0.1667 | -0.0048 | 0.4055 | -0.6931 | 0.0000 | 0.065 | 0.050 | 0.059 | 0.050 |
| 26 | 0.1000 | 0.1000 | 0.0714 | 0.4055 | 0.4055 | 0.4055 | 0.095 | 0.042 | 0.096 | 0.056 |
| 27 | 0.1000 | 0.1667 | 0.0515 | 0.4055 | 0.6931 | 0.4055 | 0.080 | 0.030 | 0.107 | 0.042 |
| 28 | 0.1667 | 0.1667 | 0.0238 | 0.6931 | 0.6931 | 0.4055 | 0.073 | 0.028 | 0.096 | 0.026 |
| 29 | 0.0000 | -0.1667 | 0.0952 | 0.0000 | -0.6931 | 0.4055 | 0.090 | 0.061 | 0.102 | 0.090 |

| Scen | $\beta_T$ | $\beta_B$ | $\beta_{TB}$ | $b_T$ | $b_B$ | $b_{TB}$ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 30 | 0.1000 | -0.1667 | 0.0961 | 0.4055 | -0.6931 | 0.4055 | 0.089 | 0.068 | 0.100 | 0.092 |
| 31 | 0.0000 | -0.1000 | -0.0923 | 0.0000 | -0.4055 | -0.4055 | 0.080 | 0.044 | 0.093 | 0.082 |
| 32 | 0.0000 | -0.1667 | -0.0833 | 0.0000 | -0.6931 | -0.4055 | 0.081 | 0.037 | 0.108 | 0.067 |
| 33 | 0.1667 | 0.1667 | -0.1061 | 0.6931 | 0.6931 | -0.4055 | 0.089 | 0.074 | 0.088 | 0.097 |
| 34 | 0.1000 | 0.1000 | 0.1182 | 0.4055 | 0.4055 | 0.6931 | 0.126 | 0.048 | 0.181 | 0.095 |
| 35 | 0.1000 | 0.1667 | 0.0905 | 0.4055 | 0.6931 | 0.6931 | 0.106 | 0.036 | 0.168 | 0.058 |
| 36 | 0.1667 | 0.1667 | 0.0556 | 0.6931 | 0.6931 | 0.6931 | 0.098 | 0.030 | 0.174 | 0.027 |
| 37 | 0.0000 | -0.1000 | 0.1714 | 0.0000 | -0.4055 | 0.6931 | 0.132 | 0.110 | 0.235 | 0.227 |
| 38 | 0.0000 | -0.1667 | 0.1667 | 0.0000 | -0.6931 | 0.6931 | 0.141 | 0.108 | 0.207 | 0.194 |
| 39 | 0.1000 | -0.1667 | 0.1667 | 0.4055 | -0.6931 | 0.6931 | 0.141 | 0.113 | 0.196 | 0.184 |
| 40 | 0.0000 | 0.0000 | -0.1667 | 0.0000 | 0.0000 | -0.6931 | 0.123 | 0.097 | 0.202 | 0.191 |
| 41 | 0.0000 | -0.1000 | -0.1500 | 0.0000 | -0.4055 | -0.6931 | 0.112 | 0.065 | 0.184 | 0.145 |
| 42 | 0.0000 | -0.1667 | -0.1333 | 0.0000 | -0.6931 | -0.6931 | 0.116 | 0.049 | 0.196 | 0.121 |
| 43 | 0.1000 | -0.1000 | -0.1667 | 0.4055 | -0.4055 | -0.6931 | 0.127 | 0.107 | 0.211 | 0.203 |
| 44 | 0.1000 | -0.1667 | -0.1606 | 0.4055 | -0.6931 | -0.6931 | 0.124 | 0.096 | 0.197 | 0.183 |
| 45 | 0.1667 | -0.1667 | -0.1667 | 0.6931 | -0.6931 | -0.6931 | 0.122 | 0.108 | 0.202 | 0.211 |
| 46 | 0.1000 | 0.1000 | -0.1706 | 0.4055 | 0.4055 | -0.6931 | 0.139 | 0.123 | 0.190 | 0.188 |
| 47 | 0.0000 | 0.0000 | -0.4000 | 0.0000 | 0.0000 | -2.1972 | 0.512 | 0.355 | 0.881 | 0.803 |
| 48 | 0.4000 | -0.4000 | -0.4000 | 2.1972 | -2.1972 | -2.1972 | 0.283 | 0.533 | 0.564 | 0.963 |

324 Power estimates are based on the likelihood ratio-based restricted deviance test in the

325 "randomize-all" trial design with biomarker prevalence $\phi = 0.1$ and randomization factor

326 $\gamma = 0.5$. $\beta_0 = 0.5$ and $b_0 = 0$. Scen = Number of scenario with respective effect size

327 combination $\beta_T$, $\beta_B$, $\beta_{TB}$ or $b_T$, $b_B$, $b_{TB}$. Logistic and linear refer to the power in the logistic and

328 linear regression model, respectively.

329

17

330    In some effect size combinations, an interaction effect is present only on one scale. In

331    scenario 8 an interaction effect is present only on the logistic scale. The interaction effect

332    size is rather small compared to the other effect sizes simulated, namely $b_{TB} = 0.5026$,

333    rendering an odds ratio of $1.6530$. Correspondingly, the power in the logistic regression

334    model to detect the interaction effect is very low at 0.077 (n=200) or 0.108 (n=500).

335    Conversely, scenarios 22 to 25 (Table 4, lower part) reflect the situation of no interaction

336    effect on the logistic scale but only on the linear scale. As in scenario 8 on the logistic scale,

337    the interaction effect sizes are rather small on the linear scale and the power in the linear

338    regression model is very low at $0.035 - 0.05$ (n=200) or $0.036 - 0.05$ (n=500).

339    The biggest differences in terms of power between the logistic and linear regression models

340    can be seen if the interaction effect sizes are most extreme and either no or main effects

341    with opposite signs are present. For example, in scenario 48, the restricted deviance test in

342    the linear regression model achieves a power of $0.533$, whereas the restricted deviance test

343    in the logistic regression model achieves a power of $0.283$ for sample size $n = 200$. This

344    scenario is characterized by a strong negative predictive effect of the biomarker, a positive

345    treatment effect and a strong negative interaction as illustrated in Fig 3A. In other scenarios,

346    the deviance test in the logistic regression model achieves a higher power than in the linear

347    regression model, for example, in scenarios 14, 16, and 20. Here the difference is between

348    $\sim$0.13 and $\sim$0.15, which is illustrated in Fig 3B for scenario 20. These are described by no

349    treatment effects and a negative predictive effect of the biomarker with an additional

350    interaction effect. For all other effect size combinations the differences in terms of power

351    are negligible.

352    S1 and S2 Tables list the corresponding type I error frequency and estimated power for the

353    same effect size combinations as Tables 2 and 4 in the "biomarker-stratified" trial design

354    with biomarker prevalence $\phi = 0.1$, randomization factors $\gamma_+ = \gamma_- = 0.5$, $n_+$ and $n_-$

355    determined by the prevalence of the biomarker $\phi$. As the same sample sizes are eventually

356    available in the four groups, the estimated frequencies are very similar to those observed in

357    the "randomize-all" trial design. Interestingly, the total number of scenarios violating

358    Bradley's liberal criterion of robustness in the "biomarker-stratified" design with sample

359    sizes determined by the prevalence of the biomarker (Table 5) is much higher than in the

360    "randomize-all" design (Table 3). Both regression models violate the criterion in about 9% of

361    the scenarios with $\beta_{TB} = b_{TB} = 0$ (logistic 317 times, linear 309 times). Again, the logistic

362    model tends to be liberal, violating the upper criterion bound, whereas the linear model

363    tends to be conservative, violating the lower criterion bound.

364

365    **Fig 3. Illustration of scenarios with notable power differences between regression models.**

366    Number of patients with a positive outcome. Based on a sample size of 100 in every

367    constellation in (A) scenario 48 characterized by a strong negative predictive effect of the

368    biomarker, a positive treatment effect and a strong negative interaction and in (B) scenario

369    20 characterized by no treatment effects and a negative predictive effect of the biomarker

370    with an additional interaction effect.

371

372    **Table 5. Number of scenarios in which type I error frequencies deviate from Bradley's**

373    **criterion [18] in the "biomarker-stratified" design.**

|  |  | n = 100 | n = 200 | n = 500 | n = 1000 | Σ |
|---|---|---|---|---|---|---|
|  | > 0.075 | 171 (20%) | 109 (13%) | 17 (2%) | 11 (1%) | 308 (9%) |
| logistic | < 0.025 | 2 (0%) | 7 (1%) | 0 (0%) | 0 (0%) | 9 (0%) |
|  | Σ | 173 (20%) | 116 (13%) | 17 (2%) | 11 (1%) | 317 (9%) |

19

|        |          |          |          |          |          |           |
| ------ | -------- | -------- | -------- | -------- | -------- | --------- |
|        | $> 0.075$ | 13 (2%)  | 14 (2%)  | 14 (2%)  | 14 (2%)  | 55 (2%)   |
| linear | $< 0.025$ | 72 (8%)  | 63 (7%)  | 61 (7%)  | 58 (7%)  | 254 (7%)  |
|        | $\Sigma$ | 85 (10%) | 77 (9%)  | 75 (9%)  | 72 (8%)  | 309 (9%)  |

374    Based on the likelihood ratio-based restricted deviance test in the "biomarker-stratified"

375    trial design with $n_+$ and $n_-$ determined by $\phi$. All 3456 scenarios with $\beta_{TB} = b_{TB} = 0$ are

376    considered.

377

378    Finally, Tables 6, 7 and 8 list the corresponding type I error frequency, scenarios in which the

379    type I error frequencies deviate from Bradley's criterion, and estimated power for the same

380    effect size combinations with randomization factors $\gamma_+ = \gamma_- = 0.5$ and fixed proportions of

381    biomarker positive and biomarker negative patients ($n_+ = n_- = {}^n/_2$). It is therefore

382    assumed that out of a larger patients' group with biomarker information, only a specified

383    number is selected and included in the trial, so that there is an equal number of biomarker

384    positive and negative cases. In this situation, the estimated type I error is very close to the

385    expected 0.05 in all scenarios with no interaction effect (Table 6), even in scenario 4.

386    Remarkably, in this trial design, the lowest numbers of scenarios violating Bradley's criterion

387    of robustness is observed (Table 7). The logistic model violates the criterion 36 times and the

388    linear model 81 times, both about 1% of all scenarios with $\beta_{TB} = b_{TB} = 0$ and $n_+, n_-$ fixed

389    at $\frac{n}{2}$. Unexpectedly, in this setting the linear model also tends to be liberal.

390

391    **Table 6. Estimated type I error frequency at the nominal two-sided 0.05 test-level in the**

392    **"biomarker-stratified" design with fixed proportion of biomarker positive and negative**

393    **patients.**

| | | | | | | | n = 200 | | n = 500 | |
|---|---|---|---|---|---|---|---|---|---|---|
| Scen | $\beta_T$ | $\beta_B$ | $\beta_{TB}$ | $b_T$ | $b_B$ | $b_{TB}$ | logistic | linear | logistic | linear |
| 1 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.054 | 0.054 | 0.054 | 0.054 |
| 2 | 0.0000 | -0.1000 | 0.0000 | 0.0000 | -0.4055 | 0.0000 | 0.061 | 0.059 | 0.045 | 0.045 |
| 3 | 0.0000 | -0.2000 | 0.0000 | 0.0000 | -0.8473 | 0.0000 | 0.049 | 0.046 | 0.044 | 0.045 |
| 4 | 0.0000 | -0.4000 | 0.0000 | 0.0000 | -2.1972 | 0.0000 | 0.055 | 0.047 | 0.051 | 0.050 |
| 5 | 0.4000 | -0.4000 | 0.0000 | 2.1972 | -2.1972 | 0.0000 | 0.060 | 0.055 | 0.056 | 0.049 |
| 6 | 0.0000 | -0.1667 | 0.0000 | 0.0000 | -0.6931 | 0.0000 | 0.046 | 0.047 | 0.044 | 0.041 |
| 7 | 0.1667 | -0.1667 | 0.0000 | 0.6931 | -0.6931 | 0.0000 | 0.053 | 0.045 | 0.044 | 0.045 |

394   Frequency estimates are based on the likelihood ratio-based restricted deviance test in the

395   "biomarker-stratified" trial design with biomarker prevalence $\phi = 0.1$, randomization factors

396   $\gamma_+ = \gamma_- = 0.5$ and $n_+ = n_- = \frac{n}{2}$. $\beta_0 = 0.5$ and $b_0 = 0$. Scen = Number of scenario with

397   respective effect size combination $\beta_T$, $\beta_B$, $\beta_{TB}$ or $b_T$, $b_B$, $b_{TB}$. Logistic and linear refer to the type

398   I error frequency in the logistic and linear regression model, respectively.

399

400   **Table 7. Number of scenarios in which type I error frequencies deviate from Bradley's**

401   **criterion [18] in the "biomarker-stratified" design.**

| | | n = 100 | n = 200 | n = 500 | n = 1000 | Σ |
|---|---|---|---|---|---|---|
| logistic | > 0.075 | 27 (3%) | 9 (1%) | 0 (0%) | 0 (0%) | 36 (1%) |
| | < 0.025 | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) |
| | Σ | 27 (3%) | 9 (1%) | 0 (0%) | 0 (0%) | 36 (1%) |
| linear | > 0.075 | 18 (2%) | 15 (2%) | 6 (1%) | 18 (2%) | 57 (2%) |
| | < 0.025 | 9 (1%) | 9 (1%) | 3 (0%) | 3 (0%) | 24 (1%) |
| | Σ | 27 (3%) | 24 (3%) | 9 (1%) | 21 (2%) | 81 (2%) |

21

402    Based on the likelihood ratio-based restricted deviance test in the "biomarker-stratified"

403    trial design with $n_+ = n_- = \frac{n}{2}$. All 3456 scenarios with $\beta_{TB} = b_{TB} = 0$ are considered.

404

405    Similar as in the previous designs, if an interaction effect is present only on one scale, it is

406    hard to detect, resulting in a low power. In general, however, the pattern of the estimated

407    power is very similar to before, with an overall higher power due to balanced sample sizes.

408

409    **Table 8. Estimated power at the nominal two-sided 0.05 test-level in the "biomarker-**

410    **stratified" design with fixed proportion of biomarker positive and negative patients.**

| Scen | $\beta_T$ | $\beta_B$ | $\beta_{TB}$ | $b_T$ | $b_B$ | $b_{TB}$ | n = 200 | | n = 500 | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | logistic | linear | logistic | linear |
| 8 | 0.2000 | -0.4000 | 0.0000 | 0.8473 | -2.1972 | 0.5026 | 0.132 | 0.047 | 0.188 | 0.038 |
| 9 | 0.0000 | -0.1000 | 0.1000 | 0.0000 | -0.4055 | 0.4055 | 0.112 | 0.111 | 0.210 | 0.210 |
| 10 | 0.0000 | 0.0000 | -0.1000 | 0.0000 | 0.0000 | -0.4055 | 0.130 | 0.129 | 0.202 | 0.199 |
| 11 | 0.1000 | -0.1000 | -0.1000 | 0.4055 | -0.4055 | -0.4055 | 0.108 | 0.110 | 0.210 | 0.209 |
| 12 | 0.1000 | -0.1667 | -0.1000 | 0.4055 | -0.6931 | -0.4055 | 0.112 | 0.114 | 0.200 | 0.210 |
| 13 | 0.1667 | -0.1667 | -0.1000 | 0.6931 | -0.6931 | -0.4055 | 0.111 | 0.116 | 0.210 | 0.218 |
| 14 | 0.0000 | -0.4000 | 0.2000 | 0.0000 | -2.1972 | 1.3499 | 0.537 | 0.354 | 0.896 | 0.709 |
| 15 | 0.0000 | 0.0000 | -0.2000 | 0.0000 | 0.0000 | -0.8473 | 0.342 | 0.328 | 0.657 | 0.636 |
| 16 | 0.0000 | -0.2000 | -0.2000 | 0.0000 | -0.8473 | -1.3499 | 0.539 | 0.361 | 0.900 | 0.694 |
| 17 | 0.2000 | -0.4000 | -0.2000 | 0.8473 | -2.1972 | -0.8473 | 0.194 | 0.431 | 0.402 | 0.806 |
| 18 | 0.4000 | -0.4000 | -0.2000 | 2.1972 | -2.1972 | -0.8473 | 0.195 | 0.448 | 0.398 | 0.815 |
| 19 | 0.0000 | -0.2000 | 0.4000 | 0.0000 | -0.8473 | 1.6946 | 0.831 | 0.831 | 0.996 | 0.996 |
| 20 | 0.0000 | -0.4000 | 0.4000 | 0.0000 | -2.1972 | 2.1972 | 0.945 | 0.869 | 0.998 | 0.995 |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 21 | 0.2000 | -0.4000 | 0.4000 | 0.8473 | -2.1972 | 2.1972 | 0.924 | 0.914 | 1.000 | 0.999 |
| 22 | 0.1000 | 0.1000 | -0.0077 | 0.4055 | 0.4055 | 0.0000 | 0.037 | 0.037 | 0.055 | 0.050 |
| 23 | 0.1000 | 0.1667 | -0.0167 | 0.4055 | 0.6931 | 0.0000 | 0.045 | 0.041 | 0.049 | 0.054 |
| 24 | 0.1667 | 0.1667 | -0.0333 | 0.6931 | 0.6931 | 0.0000 | 0.041 | 0.048 | 0.041 | 0.055 |
| 25 | 0.1000 | -0.1667 | -0.0048 | 0.4055 | -0.6931 | 0.0000 | 0.053 | 0.048 | 0.038 | 0.039 |
| 26 | 0.1000 | 0.1000 | 0.0714 | 0.4055 | 0.4055 | 0.4055 | 0.094 | 0.066 | 0.213 | 0.145 |
| 27 | 0.1000 | 0.1667 | 0.0515 | 0.4055 | 0.6931 | 0.4055 | 0.092 | 0.060 | 0.181 | 0.089 |
| 28 | 0.1667 | 0.1667 | 0.0238 | 0.6931 | 0.6931 | 0.4055 | 0.088 | 0.041 | 0.175 | 0.057 |
| 29 | 0.0000 | -0.1667 | 0.0952 | 0.0000 | -0.6931 | 0.4055 | 0.107 | 0.107 | 0.202 | 0.193 |
| 30 | 0.1000 | -0.1667 | 0.0961 | 0.4055 | -0.6931 | 0.4055 | 0.110 | 0.105 | 0.204 | 0.200 |
| 31 | 0.0000 | -0.1000 | -0.0923 | 0.0000 | -0.4055 | -0.4055 | 0.125 | 0.116 | 0.187 | 0.169 |
| 32 | 0.0000 | -0.1667 | -0.0833 | 0.0000 | -0.6931 | -0.4055 | 0.118 | 0.109 | 0.182 | 0.147 |
| 33 | 0.1667 | 0.1667 | -0.1061 | 0.6931 | 0.6931 | -0.4055 | 0.101 | 0.123 | 0.186 | 0.241 |
| 34 | 0.1000 | 0.1000 | 0.1182 | 0.4055 | 0.4055 | 0.6931 | 0.205 | 0.135 | 0.442 | 0.305 |
| 35 | 0.1000 | 0.1667 | 0.0905 | 0.4055 | 0.6931 | 0.6931 | 0.179 | 0.103 | 0.401 | 0.195 |
| 36 | 0.1667 | 0.1667 | 0.0556 | 0.6931 | 0.6931 | 0.6931 | 0.154 | 0.057 | 0.366 | 0.107 |
| 37 | 0.0000 | -0.1000 | 0.1714 | 0.0000 | -0.4055 | 0.6931 | 0.235 | 0.236 | 0.520 | 0.520 |
| 38 | 0.0000 | -0.1667 | 0.1667 | 0.0000 | -0.6931 | 0.6931 | 0.235 | 0.228 | 0.491 | 0.482 |
| 39 | 0.1000 | -0.1667 | 0.1667 | 0.4055 | -0.6931 | 0.6931 | 0.216 | 0.212 | 0.505 | 0.505 |
| 40 | 0.0000 | 0.0000 | -0.1667 | 0.0000 | 0.0000 | -0.6931 | 0.249 | 0.244 | 0.498 | 0.483 |
| 41 | 0.0000 | -0.1000 | -0.1500 | 0.0000 | -0.4055 | -0.6931 | 0.235 | 0.213 | 0.476 | 0.419 |
| 42 | 0.0000 | -0.1667 | -0.1333 | 0.0000 | -0.6931 | -0.6931 | 0.212 | 0.179 | 0.427 | 0.339 |
| 43 | 0.1000 | -0.1000 | -0.1667 | 0.4055 | -0.4055 | -0.6931 | 0.231 | 0.227 | 0.474 | 0.475 |
| 44 | 0.1000 | -0.1667 | -0.1606 | 0.4055 | -0.6931 | -0.6931 | 0.216 | 0.216 | 0.467 | 0.473 |
| 45 | 0.1667 | -0.1667 | -0.1667 | 0.6931 | -0.6931 | -0.6931 | 0.228 | 0.237 | 0.451 | 0.482 |

23

| 46 | 0.1000 | 0.1000 | -0.1706 | 0.4055 | 0.4055 | -0.6931 | 0.249 | 0.247 | 0.488 | 0.488 |
| 47 | 0.0000 | 0.0000 | -0.4000 | 0.0000 | 0.0000 | -2.1972 | 0.939 | 0.871 | 1.000 | 0.998 |
| 48 | 0.4000 | -0.4000 | -0.4000 | 2.1972 | -2.1972 | -2.1972 | 0.718 | 0.972 | 0.979 | 1.000 |

411 Power estimates are based on the likelihood ratio-based restricted deviance test in the

412 "biomarker-stratified" trial design with biomarker prevalence $\phi = 0.1$, randomization factors

413 $\gamma_+ = \gamma_- = 0.5$ and $n_+ = n_- = \frac{n}{2}$. $\beta_0 = 0.5$ and $b_0 = 0$. Scen = Number of scenario with

414 respective effect size combination $\beta_T, \beta_B, \beta_{TB}$ or $b_T, b_B, b_{TB}$. Logistic and linear refer to the

415 power in the logistic and linear regression model, respectively.

416

417 For an overview, Table 9 shows a comparison of the estimated power across the considered

418 scenarios. Here, the number of scenarios is given in which the power in the linear and

419 logistic regression model is comparable (less than 3% difference), in which one of the models

420 is slightly better (difference between 3% and 10%), and in which one of the models is better

421 (difference greater than 10%). These numbers are given for all considered scenarios and only

422 for scenarios without extreme effect constellations. For the vast majority of scenarios, the

423 difference in estimated power of the linear and logistic model is irrelevant, i.e., the

424 difference is less than 3%, and differences are smaller with larger sample sizes. If relevant

425 power differences are observed, this is usually in favor of the logistic model. Interestingly,

426 this pattern remains the same when scenarios with extreme effect combinations are not

427 considered.

428

429 **Table 9. Power comparison for restricted deviance test.**

| | | Randomize-All | | Biomarker-Stratified | | Biomarker-Stratified* | |
|---|---|---|---|---|---|---|---|
| | | n=200 | n=500 | n=200 | n=500 | n=200 | n=500 |
| all scenarios (599) | logistic >> linear | 24 | 6 | 23 | 4 | 2 | 2 |
| | | (4.0%) | (1.0%) | (3.8%) | (0.7%) | (0.3%) | (0.3%) |
| | logistic > linear | 232 | 78 | 184 | 77 | 34 | 13 |
| | | (38.7%) | (13.0%) | (30.7%) | (12.9%) | (5.7%) | (2.2%) |
| | logistic = linear | 332 | 499 | 379 | 503 | 550 | 576 |
| | | (55.4%) | (83.3%) | (63.3%) | (84.0%) | (91.8%) | (96.2%) |
| | logistic < linear | 11 | 16 | 13 | 15 | 12 | 8 |
| | | (1.8%) | (2.7%) | (2.2%) | (2.5%) | (2.0%) | (1.3%) |
| | logistic << linear | 0 | 0 | 0 | 0 | 1 | 0 |
| | | (0%) | (0%) | (0%) | (0%) | (0.2%) | (0%) |
| excluding most extreme scenarios (535) | logistic >> linear | 24 | 6 | 23 | 4 | 2 | 2 |
| | | (4.5%) | (1.1%) | (4.3%) | (0.7%) | (0.4%) | (0.4%) |
| | logistic > linear | 212 | 75 | 164 | 75 | 32 | 13 |
| | | (39.3%) | (13.9%) | (30.4%) | (13.9%) | (5.9%) | (2.4%) |
| | logistic = linear | 297 | 450 | 343 | 453 | 498 | 516 |
| | | (55.1%) | (83.5%) | (63.6%) | (84.0%) | (92.4%) | (95.7%) |
| | logistic < linear | 6 | 8 | 9 | 7 | 6 | 8 |
| | | (1.1%) | (1.5%) | (1.7%) | (1.3%) | (1.1%) | (1.5%) |
| | logistic << linear | 0 | 0 | 0 | 0 | 1 | 0 |
| | | (0%) | (0%) | (0%) | (0%) | (0.2%) | (0%) |
| excluding extreme scenarios | logistic >> linear | 24 | 6 | 23 | 4 | 2 | 2 |
| | | (4.7%) | (1.2%) | (4.5%) | (0.8%) | (0.4%) | (0.4%) |

|  |  |  |  |  |  |  |  |
|---|---|---|---|---|---|---|---|
| (515) | logistic > linear | 206 | 74 | 157 | 74 | 32 | 13 |
|  |  | (40.0%) | (14.4%) | (30.5%) | (14.4%) | (6.2%) | (2.5%) |
|  | logistic = linear | 282 | 428 | 331 | 431 | 474 | 492 |
|  |  | (54.8%) | (83.1%) | (64.3%) | (83.7%) | (92.0%) | (95.5%) |
|  | logistic < linear | 3 | 7 | 4 | 6 | 6 | 8 |
|  |  | (0.6%) | (1.4%) | (0.8%) | (1.2%) | (1.2%) | (1.6%) |
|  | logistic << linear | 0 | 0 | 0 | 0 | 1 | 0 |
|  |  | (0%) | (0%) | (0%) | (0%) | (0.2%) | (0%) |

430    Power estimates are based on the likelihood ratio-based restricted deviance test. Biomarker

431    prevalence $\phi = 0.1$, randomization factors $\gamma = \gamma_+ = \gamma_- = 0.5$. $\beta_0 = 0.5$ and $b_0 = 0$.

432    "Biomarker Stratified*" is with $n_+ = n_- = {}^n/_2$.

433    All = All scenarios with both $b_{TB} \neq 0$ and $\beta_{TB} \neq 0$.

434    Excluding most extreme scenarios = All scenarios with both $b_{TB} \neq 0$ and $\beta_{TB} \neq 0$ and

435    excluding scenarios with 2 or 3 linear regression parameters $\geq \pm 0.4$.

436    Excluding extreme scenarios = All scenarios with both $b_{TB} \neq 0$ and $\beta_{TB} \neq 0$ and excluding

437    scenarios with 2 or 3 linear regression parameters $\geq \pm 0.3$.

438    ">>" indicates power difference $> 10\%p$. ">" indicates power difference $> 3\%p$. "="

439    indicates power difference $\leq 3\%p$.

440

441    The above results were obtained from using the likelihood-based restricted deviance test for

442    interaction. Using a Wald-like test instead produces the same results in the linear model, but

443    lower type I and type II errors in the logistic model. The number of scenarios in which the

444    type I error frequencies deviate from Bradley's criterion in the Wald-like test are shown in S3

445    to S5 Tables. In addition, we presented only a limited selection of the simulation results, but

26

446    the preceding descriptions are also valid for the other simulation settings, and a compilation

447    of all results can be found in S6 Table (note that the numbers of the effect size combinations

448    in S6 Table are not the same as in Tables 2, 4, 6, 8).

## 449    Discussion and conclusions

450    The predictiveness of a biomarker can be evaluated via the treatment-by-biomarker

451    interaction in linear or logistic regression models for a binary outcome, and we have derived

452    the relationship between the effects of the linear model and the logistic model (S1

453    Appendix). The translation between ORs from the logistic and AARs from the linear model

454    might be useful, since the ARRs can in turn be used to calculate the NNT which is helpful for

455    the clinical interpretation. In a comprehensive simulation study, we compared the power of

456    the linear and logistic regression models to detect the predictiveness of a biomarker under a

457    variety of scenarios in the randomize-all and the biomarker-stratified design. In general, we

458    found that the differences in power to detect interaction were minor. Visible differences in

459    power were detected in rather unrealistic scenarios of effect size combinations and were

460    usually in favor of the logistic model. If the number of biomarker-positive and biomarker-

461    negative patients in the biomarker-stratified design was guided by the prevalence of the

462    biomarker, we did not find notable differences compared to the randomize-all design.

463    However, if equal subgroups of biomarker-positive and biomarker-negative patients could

464    be selected in the biomarker-stratified design, the power was decidedly greater owing to the

465    balanced samples sizes.

466    Different baseline probabilities were not considered in our simulations. These could have

467    impact on the power of both regression models and the power differences as well, especially

468    if they are close to 0 and 1. However, we assume that these values only play a minor role in

469    applications.

470     For choosing between the logistic and the linear model for a clinical trial that aims at

471     showing predictiveness of a biomarker one should therefore consider the following factors:

     1. The linear regression model has statistical disadvantages. For example, the predicted

       probability might be out of the 0-1-range of possible values. Furthermore, the model

       fit is rather poor if the predicted probabilities are close to 0 or 1. In the logistic

       regression model, the error terms follow a binomial distribution, and statistical

       properties are generally good for a binary outcome [19].

     2. As expected, the type I error frequency was adequate in both models, unless the

       scenarios were extreme, where the linear model was sometimes conservative.

     3. Power was comparable, again unless the effect size combinations were highly

       unusual. If there were differences, the logistic model usually had higher power than

       the linear probability model.

     4. The effects from the linear model can be interpreted in a more straightforward way,

       which was also pointed out be Hellevik [14] in the case of main effects, and ARR and

       OR can be translated into each other.

485     Thus, the choice of the appropriate regression model should always be driven by the

486     primary aim of a study [19] and is influenced by two different currents, the statistical

487     properties and the ease of interpretation. From the statistical viewpoint one should favor

488     the most sparse model. Following this, one could estimate both models and select the one

489     with the least number of non-zero estimates. However, our simulations have shown that it

490     is hard to find effect size combinations with non-zero effects on only one scale. Thus, from a

491     practical point of view one should favor the logistic regression model, and inference based

492     on the logistic regression model estimates should be theoretically more valid than inference

493     based on linear regression model estimates. Consequently, the logistic model should be

494    used if the presence of an interaction effect is to be tested. Concerning the interpretation

495    regarding the treatment effect in different groups, the linear model seems recommendable.

496    With our results in mind, it therefore seems recommendable to estimate logistic regression

497    models because of their statistical properties, test for interaction effects and calculate and

498    report both ARRs and ORs from these using the formulae provided in the appendix.

499

500    **References**

501    1.    Biomarkers and surrogate endpoints: preferred definitions and conceptual framework.

502        Clin Pharmacol Ther. 2001;69(3):89-95.

503    2.    Buyse M, Michiels S, Sargent DJ, Grothey A, Matheson A, de Gramont A. Integrating

504        biomarkers in clinical trials. Expert Rev Mol Diagn. 2011;11(2):171-82.

505    3.    Riley RD, Hayden JA, Steyerberg EW, Moons KG, Abrams K, Kyzas PA, et al. Prognosis

506        Research Strategy (PROGRESS) 2: prognostic factor research. PLoS Med.

507        2013;10(2):e1001380.

508    4.    Ziegler A, Koch A, Krockenberger K, Grosshennig A. Personalized medicine using DNA

509        biomarkers: a review. Hum Genet. 2012;131(10):1627-38.

510    5.    Bjermer L, Lemiere C, Maspero J, Weiss S, Zangrilli J, Germinaro M. Reslizumab for

511        inadequately controlled asthma with elevated blood eosinophil levels: A randomized

512        phase 3 study. Chest. 2016;150(4):789-98.

513    6.    Corren J, Weinstein S, Janka L, Zangrilli J, Garin M. Phase 3 study of reslizumab in

514        patients with poorly controlled asthma: effects across a broad range of eosinophil

515        counts. Chest. 2016;150(4):799-810.

516    7.    FitzGerald JM, Bleecker ER, Nair P, Korn S, Ohta K, Lommatzsch M, et al. Benralizumab,

517        an anti-interleukin-5 receptor alpha monoclonal antibody, as add-on treatment for

518  patients with severe, uncontrolled, eosinophilic asthma (CALIMA): a randomised,

519  double-blind, placebo-controlled phase 3 trial. Lancet. 2016.

520 8. Wang SJ, O'Neill RT, Hung HM. Approaches to evaluation of treatment effect in

521  randomized clinical trials with genomic subset. Pharm Stat. 2007;6(3):227-44.

522 9. Brookes ST, Whitely E, Egger M, Smith GD, Mulheran PA, Peters TJ. Subgroup analyses in

523  randomized trials: risks of subgroup-specific analyses; power and sample size for the

524  interaction test. J Clin Epidemiol. 2004;57(3):229-36.

525 10. Pant S, Martin LK, Geyer S, Wei L, Van Loon K, Sommovilla N, et al. Baseline serum

526  albumin is a predictive biomarker for patients with advanced pancreatic cancer treated

527  with bevacizumab: a pooled analysis of 7 prospective trials of gemcitabine-based

528  therapy with or without bevacizumab. Cancer. 2014;120(12):1780-6.

529 11. Elferink A, Van Zwieten-Boot B. Analysis based on number needed to treat shows

530  differences between drugs studied. Brit Med J. 1997;314:603.

531 12. vanderWeele T, Knol M. A tutorial on interaction. Epidemiol Methods. 2014;3(1):33-72.

532 13. Bokemeyer C, Bondarenko I, Hartmann JT, de Braud F, Schuch G, Zubel A, et al. Efficacy

533  according to biomarker status of cetuximab plus FOLFOX-4 as first-line treatment for

534  metastatic colorectal cancer: the OPUS study. Ann Oncol. 2011;22(7):1535-46.

535 14. Hellevik O. Linear versus logistic regression when the dependent variable is a

536  dichotomy. Qual Quant. 2009;43:59-74.

537 15. Kraft P, Yen YC, Stram DO, Morrison J, Gauderman WJ. Exploiting gene-environment

538  interaction to detect genetic associations. Human Heredity. 2007;63(2):111-9.

539 16. R Core Team. R: A Language and environment for statistical computing. Vienna, Austria.

540  https://www.R-project.org/. R Foundation for Statistical Computing; 2016.

541  17. Lang M, Bischl B, Surmann D. batchtools: Tools for R to work on batch systems. J Open

542      Source Softw. 2017;2(10).

543  18. Bradley JV. Robustness? Br J Math Stat Psychol. 1978;31:144-52.

544  19. Ganzach Y, Saporta I, Weber Y. Interaction in linear versus logistic models: a substantive

545      illustration using the relationship between motivation, ability, and performance. Organ

546      Res Meth. 2000;3(3):237-53.

547

# Supporting information

**S1 Appendix. Relation between absolute risk reductions from linear probability models and odds ratios from logistic regression models.**

**S2 Appendix. Simulation code.** Refer to included README for further information.

**S1 Table. Estimated type I error frequency at the nominal two-sided 0.05 test-level in the "biomarker-stratified" design with biomarker prevalence 0.1.** Frequency estimates are based on the likelihood ratio-based restricted deviance test in the "biomarker-stratified" trial design with biomarker prevalence $\phi = 0.1$, randomization factors $\gamma_+ = \gamma_- = 0.5$ and $n_+$ and $n_-$ are determined by $\phi$. $\beta_0 = 0.5$ and $b_0 = 0$. Scen = Num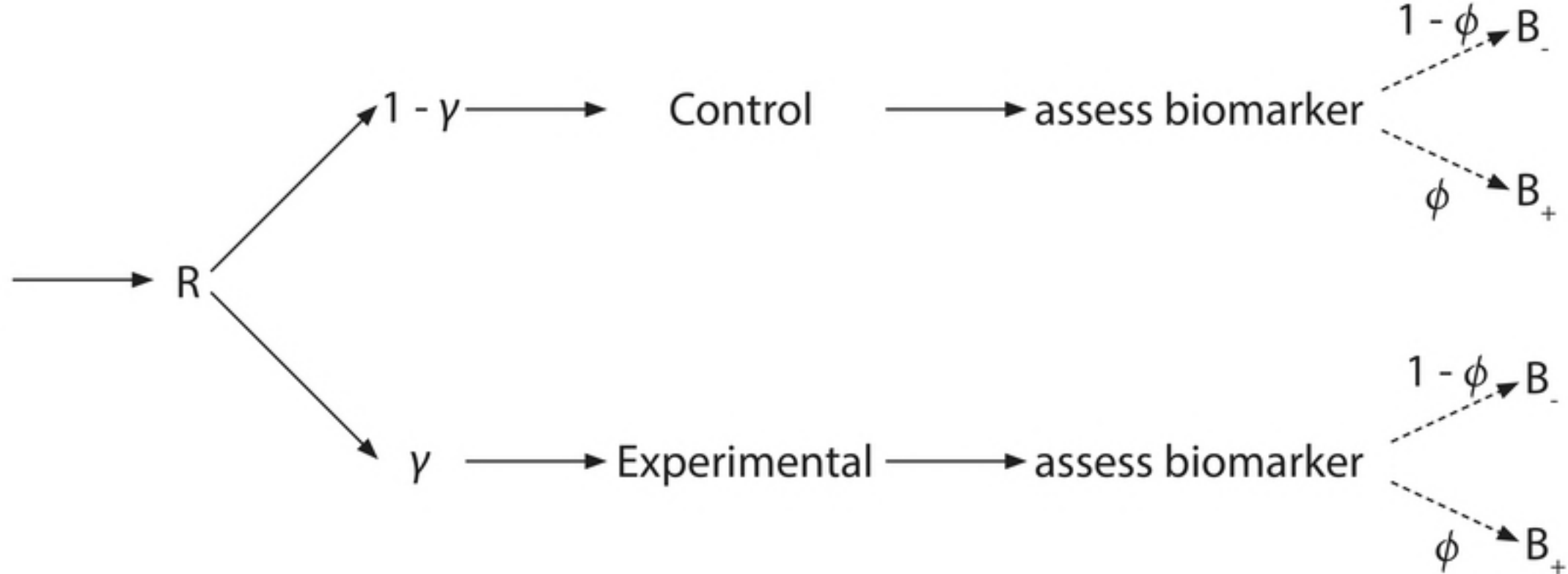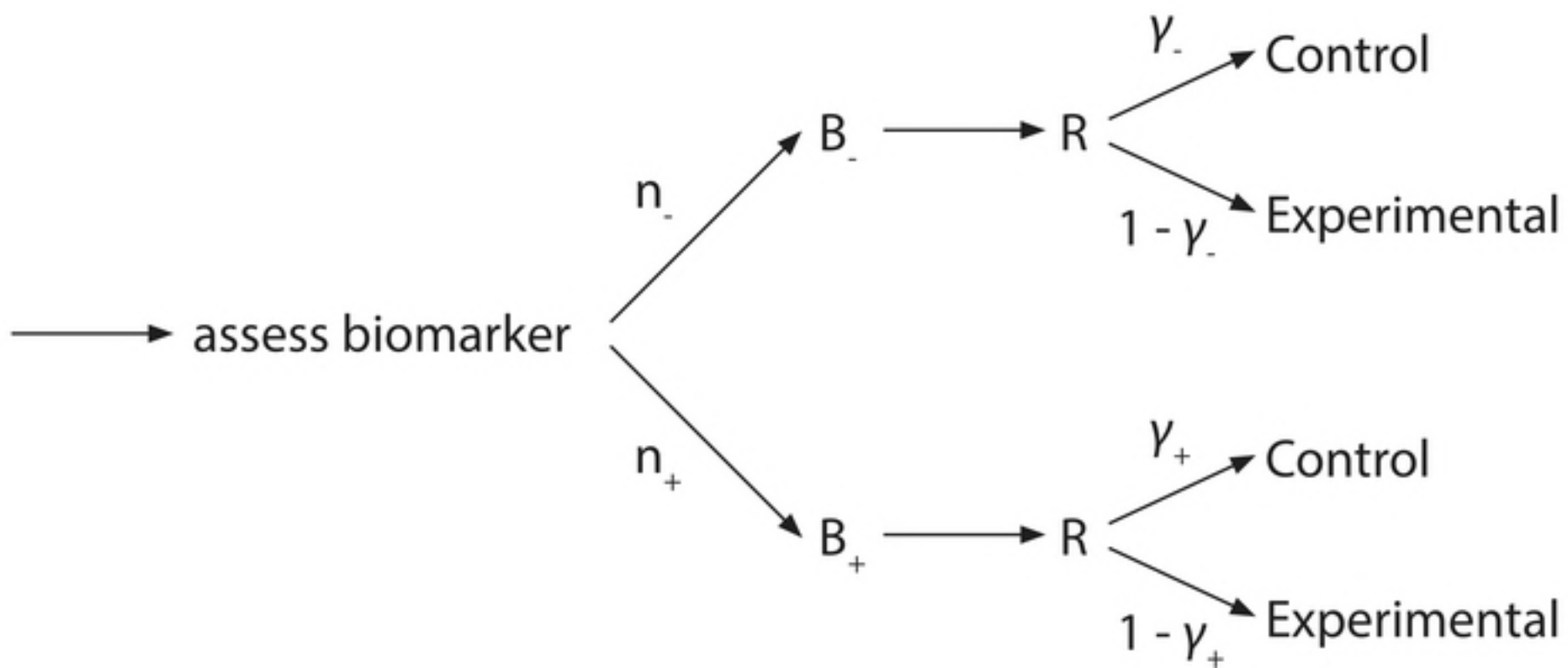ber of scenario with respective effect size combination $\beta_T$, $\beta_B$, $\beta_{TB}$ or $b_T$, $b_B$, $b_{TB}$. Logistic and linear refer to the type I error frequency in the logistic and linear regression model, respectively.

**S2 Table. Estimated power at the nominal two-sided 0.05 test-level in the "biomarker-stratified" design with biomarker prevalence 0.1.** Power estimates are based on the likelihood ratio-based restricted deviance test in the "biomarker-stratified" trial design with biomarker prevalence $\phi = 0.1$, randomization factors $\gamma_+ = \gamma_- = 0.5$ and $n_+$ and $n_-$ are determined by $\phi$. $\beta_0 = 0.5$ and $b_0 = 0$. Scen = Number of scenario with respective effect size

31

564 combination $\beta_T$, $\beta_B$, $\beta_{TB}$ or $b_T$, $b_B$, $b_{TB}$. Logistic and linear refer to the power in the logistic and

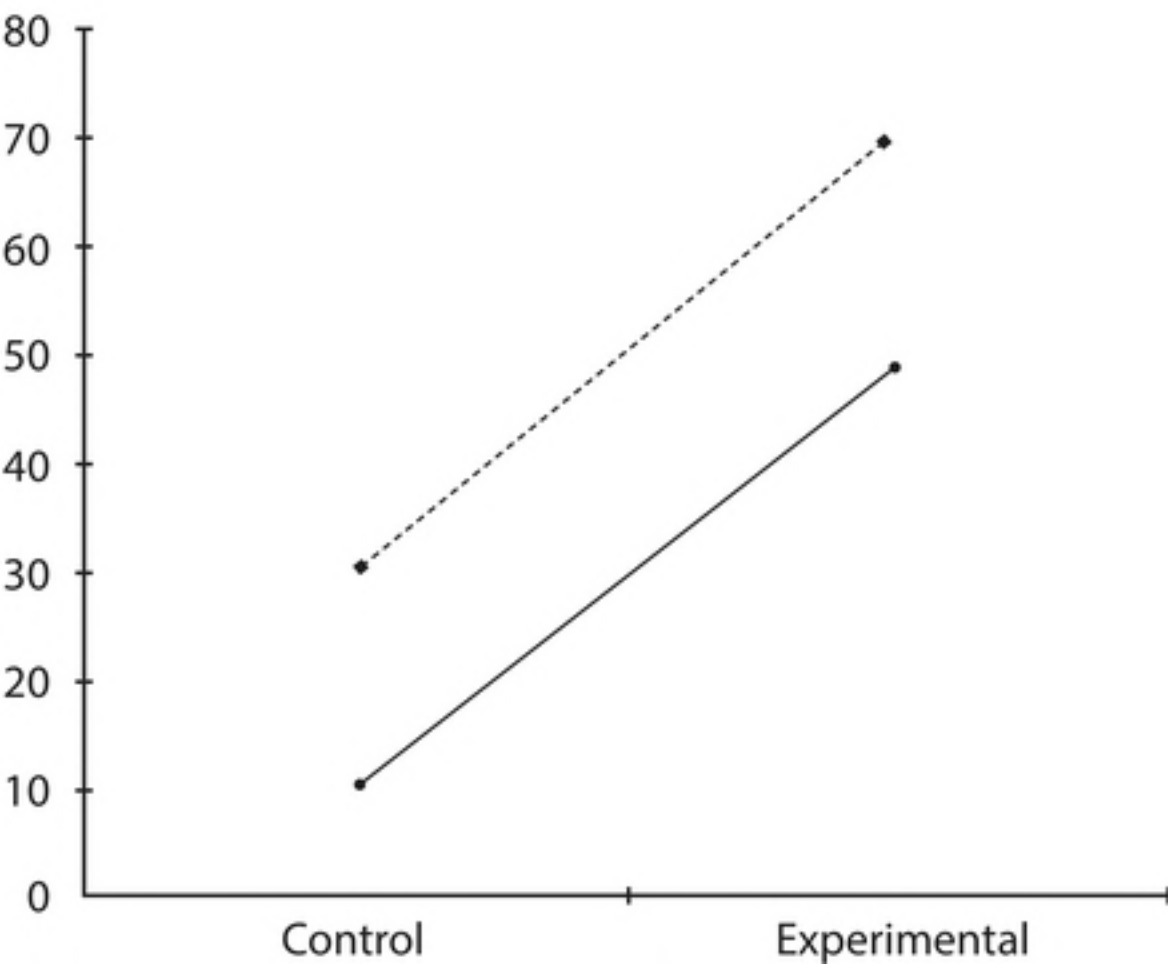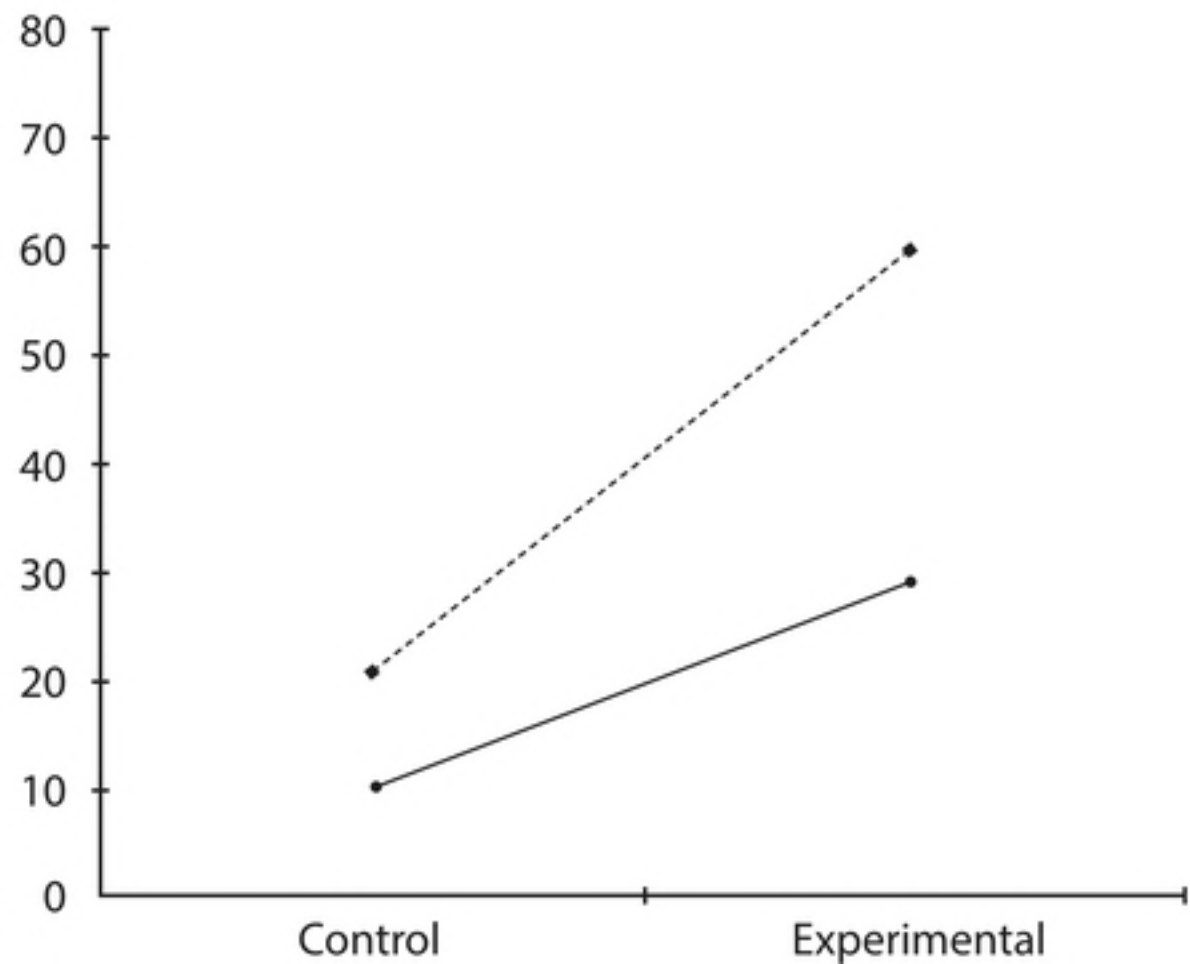565 linear regression model, respectively.

566 **S3 Table. Number of scenarios in which type I error frequencies deviate from Bradley's**

567 **criterion [18] in the "randomize-all" design.** Based on the Wald-test in the "biomarker-

568 stratified" trial design. All 1152 scenarios with $\beta_{TB} = b_{TB} = 0$ are considered.
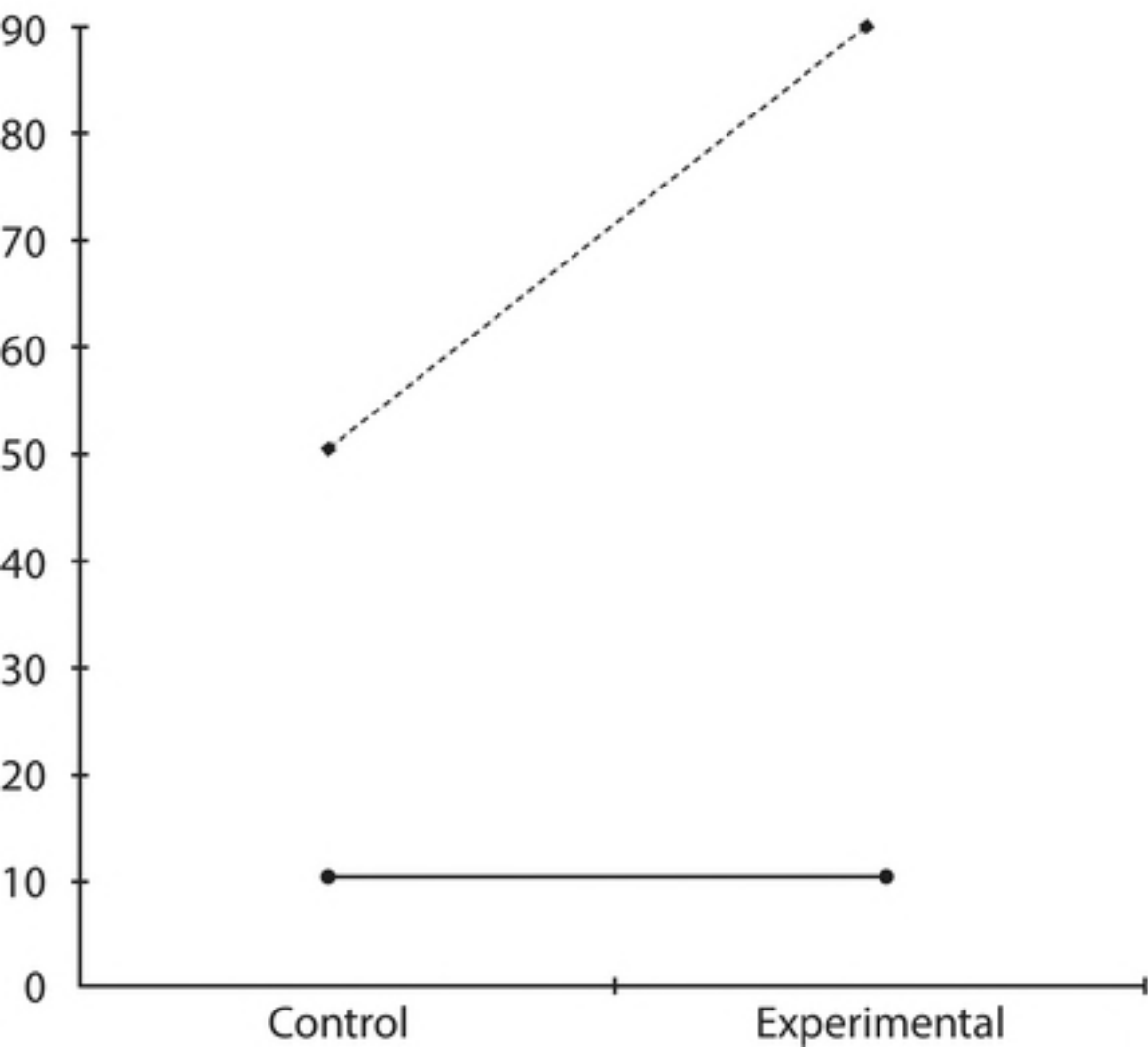
569 **S4 Table. Number of scenarios in which type I error frequencies deviate from Bradley's**

570 **criterion [18] in the "biomarker-stratified" design.** Based on the Wald-test in the

571 "biomarker-stratified" trial design with $n_+$ and $n_-$ determined by $\phi$. All 3456 scenarios with

572 $\beta_{TB} = b_{TB} = 0$ are considered.

573 **S5 Table. Number of scenarios in which type I error frequencies deviate from Bradley's**

574 **criterion [18] in the "biomarker-stratified" design.** Based on the Wald-test in the

575 "biomarker-stratified" trial design with $n_+ = n_- = \frac{n}{2}$. All 3456 scenarios with $\beta_{TB} = b_{TB} = 0$

576 are considered.

577 **S6 Table. Compilation of all simulation results.** The numbers of the effect size combinations

578 are not the same as in Tables 2, 4, 6, 8.

**A**

**B**