

Synergy detection: A practical guide to statistical assessment of potential drug combinations

Peer-reviewed author version

Makariadou, Elli; Wang , Xuechen; Hein, Nicholas; DERESA, Negera Wakgari; Mutambanengwe, Kathy; Verbist, Bie & THAS, Olivier (2024) Synergy detection: A practical guide to statistical assessment of potential drug combinations. In: PHARMACEUTICAL STATISTICS,.

DOI: 10.1002/pst.2383

Handle: <http://hdl.handle.net/1942/42872>

1 Synergy detection: a practical guide to statistical assessment
2 of potential drug combinations.

3 Elli Makariadou¹, Xuechen Wang^{2*}, Nicholas Hein², Negera W. Deresa³, Kathy
4 Mutambanengwe⁴, Bie Verbist¹, and Olivier Thas^{3,5,6}

5 ¹Janssen Pharmaceutical Companies of Johnson and Johnson, Beerse, Belgium

6 ²Janssen Pharmaceutical Companies of Johnson and Johnson, Spring House, USA

7 ³Data Science Institute, I-Biostat, Hasselt University, Hasselt, Belgium

8 ⁴Open Analytics, Antwerpen, Belgium

9 ⁵National Institute for Applied Statistics Research Australia, University of
10 Wollongong, Wollongong, New South Wales, Australia

11 ⁶Department of Applied Mathematics, Computer Science and Statistics, Ghent
12 University, Ghent, Belgium

13 *Address correspondence to: xwang449@its.jnj.com

14 **Abstract**

15 Combination treatments have been of increasing importance in drug development across ther-
16 apeutic areas to improve treatment response, minimize the development of resistance, and/or
17 minimize adverse events. Pre-clinical in-vitro combination experiments aim to explore the po-
18 tential of such drug combinations during drug discovery by comparing the observed effect of
19 the combination with the expected treatment effect under the assumption of no interaction
20 (i.e, null model). This tutorial will address important design aspects of such experiments to
21 allow proper statistical evaluation. Additionally, it will highlight the Biochemically Intuitive
22 Generalized Loewe methodology (BIGL R package available on CRAN) to statistically detect
23 deviations from the expectation under different null models. A clear advantage of the method-
24 ology is the quantification of the effect sizes, together with confidence interval while controlling
25 the directional false coverage rate. Finally, a case study will showcase the workflow in analyzing
26 combination experiments.

27 **1 Introduction**

28 Combination therapy, a treatment modality that combines two or more therapeutic agents, is of
29 growing importance in drug development across multiple therapeutic areas. Co-administration of

30 compounds may be necessary to account for disease complexity and increase efficacy while poten-
 31 tially reducing drug resistance, and/or minimizing adverse events. Consequently, combinations of
 32 compounds are routinely screened in pre-clinical in-vitro experiments to identify the most effective
 33 drug combinations.

34 The establishment of a methodology to quantify the presence of synergistic or antagonistic effects
 35 is of critical importance. Such an assessment typically relies on the dose-response curves of individual
 36 compounds, called monotherapies. Synergy or antagonism is detected when the observed response
 37 of a drug combination is different from the expected treatment response under the assumption
 38 of no interaction (i.e., the null model) such that the direction of deviation determines synergy or
 39 antagonism. The expected treatment responses are derived solely from the monotherapies. Several
 40 null models, including Highest Single Agent (HSA) [1], the Bliss Independence Model [2] and the
 41 Loewe additivity model [3] have been proposed in the literature without an agreement on the most
 42 suitable choice [4]. The models differ on the assumptions of expectation under no interaction and
 43 thus, can differ on the conclusions about the detection or degree of synergy/antagonism. For the
 44 remainder of this tutorial paper, without loss of generality, we will focus on synergistic effects.

45 Many software packages, relying on the above concept, often referred to as deviance assessment,
 46 have been published. Alternatively, one can perform some parametric modelling of a synergy index
 47 [5, 6] which is out of scope for this tutorial paper. Table 1 shows an overview of deviance assessment
 48 software packages frequently used in the pharmaceutical industry. To the best of our knowledge the
 49 details of the software packages are correct; however, we were unable to test these software packages
 50 and we extracted the details from publicly available documentation.

Table 1: Overview of deviance assessment software packages.

Software package	Accessibility	Monotherapy	Null Models	Statistics reported	
				Overall	Individual contributions
SynergyFinder	Free	4PL, LOESS or LM	HSA, Bliss, Loewe, and ZIP	Point estimate (sd)	Point estimate (sd)
SynergyFinder Plus	Free	4PL, LOESS or LM	HSA, Bliss, Loewe, and ZIP	Point estimate p-value	Point estimate CI (normal bootstrap)
MacSynergy™	Free	No	Bliss	Point estimate	Point estimate (sd)
Genedata Screener®	Commercial	4PL, 3PL, 2PL	HSA, Bliss, Loewe	Point estimate	Point estimate
CombeneFit	Free	3PL (b ¹ =1)	HSA, Bliss, Loewe	Point estimate	Point estimate p-value
Chalice™	Commercial	3PL	HSA, Bliss, Loewe and Boost	Point estimate (sd)	Point estimate (sd)
BIGL	Free	4PL, 3PL, 2PL	HSA, Bliss, Generalized Loewe ²	Point estimate CI (wild bootstrap) p-value	Point estimate CI (wild bootstrap) p-value

¹b represents the lower asymptote.

²The Generalized Loewe null model allows for partial response of the monotherapy data.

4PL, 3PL, 2PL are abbreviations for four-, three- and two-parameter logistic regression, respectively. LM is an abbreviation for linear model. CI is an abbreviation for confidence interval.

51 Most software packages start with modelling the monotherapy data using a dose-response rela-
 52 tionship (e.g., 4-parameter logistic regression). The modelled monotherapy data are then used
 53 to predict the expected treatment effects under a specified null model. MacSynergy™ [7] uses the

54 observed responses to predict the expected treatment responses. Next, the divergence of the ob-
55 served responses from the expected treatment responses are summarized in synergy scores as excess
56 responses. For any statistical analysis, it is important to evaluate the point estimate (i.e., excess
57 response) relative to the variability of the data expressed as either a confidence interval (CI) and/or
58 p-value. Genedata Screener[®], Chalice[™], MacSynergy[™], and SynergyFinder [8–10] are reporting the
59 excess responses together with observed standard deviations for the latter three. The excess re-
60 sponses in the individual combination points are further summarized as an overall excess response
61 and is often evaluated using a threshold. Routinely, this threshold is not subjected to any hypothesis
62 testing. Failure to evaluate the excess response(s) relative to either CIs or p-values increases the risk
63 of reporting false synergies (i.e., false positives); the rate of these false positives (error rate) should
64 be controlled. SynergyFinder Plus [11], which is an updated version of the original SynergyFinder,
65 recently added a bootstrapping approach to retrieve CIs around the excess responses to reduce the
66 risk of false positive calls, a procedure similar to the Biochemically Intuitive Generalized Loewe
67 (BIGL) implementation [12]. Combenefit [13], alternatively, performs a one sample T-test on the
68 excess responses, whereas BIGL uses an F-test.

69 The remainder of this tutorial paper will describe important design aspects of drug combination
70 experiments, describe proper statistical evaluation using the BIGL R package [14], and explain the
71 different assumptions of the underlying null models. The BIGL R package was chosen for its ability
72 to:

- 73 1. Incorporate several widely used null models, including HSA, Bliss, and Loewe while allowing
74 for partial response of the monotherapy data.
- 75 2. Provide flexibility on the monotherapy dose-response models.
- 76 3. Perform statistical testing with error rate control, while relaxing distributional assumptions
77 via bootstrapping.
- 78 4. Provide effect size estimates (i.e., excess responses), for each combination point and an average
79 effect size, and their confidence intervals.

80 Lastly, a case study will be presented to illustrate the use of the BIGL R package and how to
81 interpret the results and visualizations from the BIGL R package output.

82 **2 Experimental design**

83 Prior to performing any combination experiment involving two or more monotherapies, experiments
84 should be performed in which the monotherapies are profiled under the same conditions of the combi-
85 nation experiment (e.g., assay, cell line, E:T ratio, incubation time, etc). Profiling the monotherapy
86 is in the form of a dose-response relationship which describes the magnitude of response as a func-
87 tion of dose. This dose-response relationship can be described by dose-response curves and can be
88 mathematically modeled using either a 4-parameter logistic (4PL), 3PL, or 2PL regression model
89 [15].

90 The 4PL regression model is defined using four parameters that are related to the graphical
 91 properties of the curve, i.e., lower asymptote (b), upper asymptote (m), inflection point (EC_{50}), and
 92 hill slope (h). These four parameters can be used to express the magnitude of response f , given
 93 dose, d , as

$$f(d) = b + \frac{(m - b)}{1 + \left(\frac{EC_{50}}{d}\right)^{|h|}} \quad (1)$$

94 Parameters of the 4PL regression model can be fixed. For example, for a particular assay, the
 95 absence of any compound elicits no response. One may then fix the lower asymptote of the 4PL
 96 regression model to zero. Often, the 4PL regression model with fixed lower or upper asymptote is
 97 referred to a 3PL regression model. Fixing an additional parameter (e.g., hill slope) would create a
 98 2PL regression model.

99 Accurate estimation of the 4PL model parameters is paramount since prediction of the combined
 100 treatment responses, assuming no interaction, are estimated using the 4PL monotherapy model
 101 parameters. Accurate estimation of these parameters is dependent on the doses selected and the
 102 magnitude of responses they elicit. The entire sigmoidal pattern of the monotherapy curve should
 103 be covered with the selected doses. A recommendation of the ideal spread of doses is provided in
 104 Figure 1, with 2 points at each of the asymptotes and 3 points on the linear part of the curve.
 105 Depending on the hill slope of the monotherapy curve, a certain dilution series should be chosen to
 106 follow the above recommendation, which is in line with published guidelines [16].

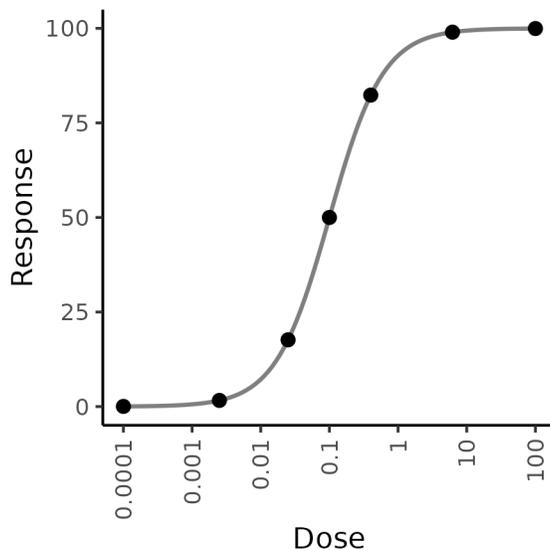


Figure 1: Ideal dilution series with 7 doses.

107 After choosing the desired doses of the monotherapies, the combination experiment needs to
 108 be set up, with two common designs being applied, the ray (dose gradients) and the checkerboard
 109 (factorial) designs. The ray design uses fixed ratios of doses while the checkerboard crosses two sets
 110 of doses. Examples are given in Figure 2a and 2b for the checkerboard and ray design, respectively.

111 For the ray design, it is typical to fix the dose titration of one of the monotherapies and multiply
112 the concentrations of the fixed doses by a scalar to get desired ratio for a ray. As such, ray designs
113 tend to have the same number of titrations for each monotherapy and require the scientist to create
114 additional dilutions of one of the monotherapies for each specified ray. Ray designs can be an efficient
115 use of resources and exploration of a set of ratios; however, prior knowledge is often required for
116 proper selection of the rays (drug ratios).

117 The checkerboard experimental design is a more comprehensive design when prior knowledge
118 is limited. The checkerboard design is a factorial design in which the doses of the monotherapies
119 are crossed with each other. As such, the checkerboard design explores many ratios; however, the
120 number of titrations for each ratio is limited (due to design) excluding the 1:1 ratio (i.e., typically the
121 diagonal of the checkerboard). Furthermore, only a single set of dose titrations for each monotherapy
122 is required; however, the checkerboard design tends to require more plate real estate than the ray
123 design. Lastly, the checkerboard is limited to pairs of monotherapies whereas the ray design can
124 more easily be scaled to triple-, quadruple-, etc., drug combinations.

125 In summary, we recommend the checkerboard design for the discovery phase of drug development
126 when prior knowledge of biological mechanism is limited, as it's easier and more convenient requiring
127 just one set of dose titrations for each monotherapy. However, when prior knowledge is available,
128 the ray design is highly valuable, offering targeted insights. Additionally, the ray design allows dose
129 response curves to be fit and compared to each ray at the expense of additional laboratory labor,
130 i.e., each ray requires a unique set of dose titrations.

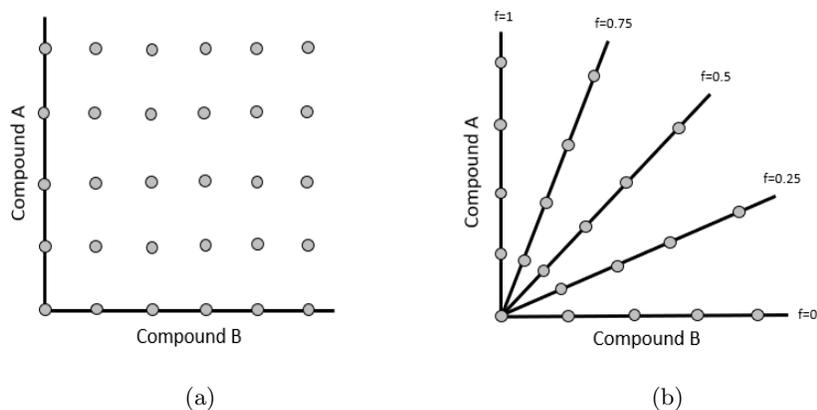


Figure 2: Drug combination designs. The dots indicate the tested doses of compounds A and B (a) Checkerboard design. (b) Ray design.

131 Please note that as with any in-vitro experiment, sufficient control wells should be included to
132 access the quality of the assay. This topic is beyond the scope of this tutorial, and we recommend
133 using established guidelines within your organization.

134 Regardless of the chosen design, the BIGL methodology provides a harmonized framework to
135 analyze drug combination experiments through statistical evaluation of the differences between the
136 expected and observed responses. To estimate variability, replicates of both the monotherapy and

137 the combination data are required. We currently recommend a minimum of 4 replicates to have an
138 accurate estimation of the observed variability and sufficient power to detect excess responses (sup-
139 portive information). Furthermore, if the experimental condition requires evaluation of donor effects
140 (e.g., biologics), we recommend against using donors as replicates due to the increased variability.
141 Instead, we recommend replicates within donor and analyzing the donors separately.

142 **3 Methods**

143 Since the BIGL methodology can be categorized as a deviance assessment methodology, a general
144 workflow to detect synergy can be followed:

- 145 1. Fit the dose-response curves for the monotherapies.
- 146 2. Predict the responses for the combinations using the chosen null model.
- 147 3. Estimate the effect sizes, together with CIs.

148 *Fitting monotherapies*

149 The BIGL R package offers flexibility in how the monotherapy dose-response relationships are
150 modelled. The default implementation models the monotherapies using 4PL regression models with
151 a shared asymptote. Depending on whether the monotherapy dose-response curves are increasing
152 or decreasing, either a common lower asymptote or upper asymptote, respectively, is assumed for
153 both drugs. The BIGL R package allows the lower asymptote, upper asymptote, and/or hill slope
154 parameter to be set to a fixed value reducing the 4PL to either a 3PL or 2PL regression model. Setting
155 a parameter to a fixed value must be done with caution and based on biological understanding.
156 Additionally, the BIGL R packages allows linear constraints on the 4PL parameters to facilitate the
157 monotherapy fitting. For instance, if the dose-response curve is decreasing, the minimum estimation
158 could be constrained to be above 0 for an improved biological interpretation.

159 *Predicting expected combination response*

160 The BIGL R package integrates three popular null models, each assumes its own underlying
161 mechanism, which characterizes the no interaction or the expected outcome under a combination
162 of drugs. The HSA quantifies the degree of synergy as the excess over the maximum monotherapy
163 response. Bliss independence assumes the drugs acted independently and synergy is evaluated as the
164 excess of the multiplicative effects of the single drugs. Loewe instead, is a dose-effect based model,
165 calculating the additive (i.e., no interaction) effects as if the single drugs where exchangeable (see
166 Table 2) resulting in an additive effect when the drug is combined with itself. It's worth noting that
167 the true mechanism is often unknown and the described null models are often a simplification of the
168 underlying biology. Therefore, the choice of a particular null model should be informed by domain
169 expertise or results should be compared across null models. Implicitly, the null models assume the
170 same maximal responses of the two monotherapies used in the combination. This assumption is
171 often violated in practice. The BIGL R package relaxes this assumption through the use of the

172 Generalized Loewe null model. Details of this approach are out of the scope of this tutorial but we
 173 refer to Van der Borghet 2017 and Thas 2022 for more information and additional alternatives. [12,
 174 17].

Table 2: Overview of the null models.

Null model	Assumption	Formula
HSA	Expected effect is highest effect of monotherapies	$f_{12}(d_1, d_2) = \max(f_1(d_1), f_2(d_2))$
Bliss independence	Drugs' effects do not interfere with one another	$f_{12}(d_1, d_2) = f_1(d_1) + f_2(d_2) - f_1(d_1)f_2(d_2)$
Loewe	Two compounds have same mode of action	$\frac{d_1}{f_1^{-1}(f_{12})} + \frac{d_2}{f_2^{-1}(f_{12})} = 1$

175 In Table 2, $f_{12}(d_1, d_2)$ or f_{12} represents the response at the combination of dose d_1 for drug 1
 176 and dose d_2 for drug 2. $f_i(d_i)$ is the response at dose d_i for drug i . $f_i^{-1}(x)$ is the inverse function of
 177 $f_i(d_i)$. It represents the dose of drug i that will produce a response of x .

178 *Estimating effect sizes*

179 The next step is the estimation of the effect sizes and their standard errors. The monotherapy
 180 data will be referred to as the on-axis points and the combination data will be referred to as the
 181 off-axis points. The effect size at off-axis point $i = 1, \dots, N$ is defined as the average deviation from
 182 the expected treatment response under a null model,

$$E_i = \frac{1}{n_i} \sum_{j=1}^{n_i} (R_{ij} - \hat{R}_i) = \bar{R}_i - \hat{R}_i,$$

183 where R_{ij} represents the observed response at off-axis point i for replicate j , n_i represents the
 184 number of replicates at point i , \bar{R}_i is the sample mean of the observed responses at point i , \hat{R}_i
 185 is the estimated treatment response under the chosen null model at off-axis point i , and N is the total
 186 number of off-axis points. An overall effect size of synergy which is also named single effect measure
 187 is defined as $\bar{E} = \frac{1}{N} \sum_{i=1}^N E_i$. With $\mathbf{E}^t = (E_1, \dots, E_N)$ the vector of all effect size estimates, the
 188 variance-covariance matrix of \mathbf{E} can be written as [12, 17]

$$\Sigma = \text{Var}(\mathbf{E}) = \sigma_0^2 \mathbf{C}_p + \sigma_1^2 \mathbf{D} \quad (2)$$

189 where σ_0^2 is the residual variance of the on-axis responses, σ_1^2 is the residual variance of the off-axis
 190 responses, \mathbf{D} is a diagonal matrix with elements $1/n_i$ ($i = 1, \dots, N$), and \mathbf{C}_p is the correlation
 191 matrix of the expected treatment responses $(\hat{R}_1, \dots, \hat{R}_N)$. Specifically, σ_0^2 is estimated as the mean
 192 squared error of the 4PL regression model fitting to the monotherapy data. And σ_1^2 is estimated as
 193 the average variance of each off-axis point $\frac{1}{N} \sum_i \sum_j A_{ij}^2$ such that $A_{ij} = R_{ij} - \hat{R}_i$. If the assumption
 194 of constant variance at the off-axis points does not hold, the variance σ_1^2 in (2) can be replaced by
 195 a diagonal matrix with model-based variances on the diagonal positions [12]. The correlation matrix

196 \mathbf{C}_p is estimated by means of a bootstrap procedure. We outline the bootstrap procedure here in
 197 some detail, because it is also needed in the next section.

198 The estimation of \mathbf{C}_p is as follows:

- 199 1. Construct a bootstrap sample of on-axis observations. Specifically, resample the residuals from
 200 the 4PL monotherapy regression with replacement. Add these resampled residuals to the fitted
 201 values of the 4PL regression model creating a bootstrap sample. Using the bootstrap sample,
 202 re-fit the monotherapeutic dose-response curves.
- 203 2. Based on the re-fitted dose-response curves, new estimates of the expected treatment responses
 204 under the null model are computed for each off-axis point i : \hat{R}_i^b for bootstrap replicate b .
- 205 3. Repeat steps 1 and 2 many times (e.g., 1000 times, $B = 1000$).
- 206 4. Calculate the sample correlation matrix of the B vectors $(\hat{R}_1^b, \dots, \hat{R}_N^b)$, denoted by $\hat{\mathbf{C}}_p$.

207 The sample correlation matrix $\hat{\mathbf{C}}_p$, is an estimator of \mathbf{C}_p . The estimator of the covariance matrix
 208 $\mathbf{\Sigma}$ can now be written as $\hat{\mathbf{\Sigma}} = \hat{\sigma}_0^2 \hat{\mathbf{C}}_p + \hat{\sigma}_1^2 \mathbf{D}$. The square roots of its diagonal elements (denoted by
 209 s_i) are estimates of the standard errors of the effect sizes. The variance of the overall effect size is
 210 given by $\text{Var}(\bar{E}) = \text{Var}\left(\frac{1}{N} \sum_{i=1}^N E_i\right)$ and can be estimated as $\frac{1}{N^2} \mathbf{1}^t \hat{\mathbf{\Sigma}} \mathbf{1}$.

211 *Controlling the directional false coverage rate*

212 The original BIGL methodology [17] framed synergy detection in a classical multiple hypothesis
 213 testing paradigm, by constructing hypothesis tests based on the effect size estimates, E_i , and their
 214 standard errors, s_i , and by controlling the familywise error rate (FWER) at some nominal level.
 215 Despite the correctness of this procedure and a positive empirical evaluation[12] we have replaced the
 216 hypothesis testing method with a procedure that makes use of simultaneous confidence intervals to
 217 control the *directional false coverage rate* (dFCR). Before defining the dFCR, we give two drawbacks
 218 of the original approach: (1) controlling the FWER at 5% results in a very conservative detection
 219 method (small sensitivity); (2) the results from this testing procedure do not always agree with
 220 what would be concluded if confidence intervals were used instead. The first issue could have been
 221 resolved by controlling e.g. the false discovery rate (FDR) instead of the FWER, and a solution to
 222 the second problem could have been found in aligning the test and CI procedures. However, instead
 223 we have chosen to develop a procedure that controls the dFCR.

224 First, we formally describe a generic method for synergy/antagonism detection based on con-
 225 fidence intervals. Let $[L_i, U_i]$ denote a confidence interval for the effect size at off-axis point i . If
 226 $0 \notin [L_i, U_i]$, then we conclude that there is a synergistic or antagonistic effect. We will use the
 227 notation τ_i for the true effect size at off-axis point i (i.e., E_i is an estimate of τ_i).

228 The original definition of the FCR[18] can then be formulated as

$$\text{FCR} = \mathbb{E}\left(\frac{F}{m}\right), \quad (3)$$

229 where m is the number of off-axis points, and F is the number of intervals among these m points,
 230 that do not cover the true effect size τ_i :

$$F = \# \{i : \tau_i \notin [L_i, U_i]\}.$$

231 This FCR makes sense in *selective inference*, i.e. statistical inference after the data-driven selection of
 232 a subset of parameters. However, in our context we want all confidence intervals to be interpretable,
 233 we therefore change the definition from a conditional to a marginal interpretation. Based on very
 234 early ideas of the concept of the FDR [19], we further adapt the definition towards a directional
 235 FCR. Equation (3) still applies, but now with

$$F = \# \{i : \tau_i \notin [L_i, U_i] \text{ and } d(L_i, U_i) \neq \text{sign}(\tau_i)\} \quad (4)$$

236 where $\text{sign}(\tau)$ equals 1 if $\tau > 0$, 0 if $\tau = 0$ and -1 if $\tau < 0$, and $d(L, U)$ equals 1 if $L, U > 0$, 0 if
 237 $L < 0$ and $U > 0$ and -1 if $U, L < 0$. In other words, F counts the number of off-axis points for
 238 which the conclusion (synergistic / antagonistic / no-effect) is wrong. This is illustrated in Figure
 239 3. Thus, if the dFCR is controlled at 10%, then, on average, 90% of the CIs either contain the true
 240 effect size, or at least these CIs result in a correct (directional) conclusion.



Figure 3: Illustration of dFCR. The vertical lines represent confidence intervals and the dots are the true effect sizes τ_i . Left: intervals that contribute to F in the definition of the dFCR. Right: intervals that do not contribute to F .

241 For controlling the dFCR at a nominal level α , the lower and upper bounds L_i and U_i can be
 242 found by means of the bootstrap procedure that was described earlier, but with two additional steps:

- 243 1. At each off-axis point i , a Wild bootstrap sample of n_i responses is obtained as $R_{ij}^b = \bar{R}_i +$
 244 $\nu_{ij}^b A_{ij}$, where ν_{ij}^b is randomly sampled from a distribution with stochastic representation

$$\left(\delta_1 + \frac{V_1}{\sqrt{2}}\right) \left(\delta_2 + \frac{V_2}{\sqrt{2}}\right) - \delta_1 \delta_2$$

245 with V_1 and V_2 two independent standard normal random variables, $\delta_1 = 0.5(\sqrt{17/6} + \sqrt{1/6})$
 246 and $\delta_2 = 0.5(\sqrt{17/6} - \sqrt{1/6})$ [20].

247 2. Compute the averages of the bootstrap responses $\bar{R}_i^b = \frac{1}{n_i} \sum_{j=1}^{n_i} R_{ij}^b$, the bootstrap effect sizes
 248 $E_i^b = \bar{R}_i^b - \hat{R}_i^b$ and their standard errors s_i^b , for all off-axis points $i = 1, \dots, N$.

249 After the bootstrap procedure has finished, compute the following intervals for a sequence of S
 250 constants t^s , $s = 1, \dots, S$ (e.g. 1, 1.1, 1.15, 1.2, ..., 2.95, 3): $[L_i^b(t^s), U_i^b(t^s)]$, with

$$L_i^b(t) = E_i^b - ts_i^b \quad \text{and} \quad U_i^b(t) = E_i^b + ts_i^b.$$

251 For each $s = 1, \dots, S$ and each $b = 1, \dots, B$, the number of directional false coverages F can
 252 be computed as in Equation (4); let F^{sb} denote this number. Averaging over the B bootstrap
 253 runs, gives numbers $F^s = \frac{1}{B} \sum_{b=1}^B F^{sb}$, and F^s/m may be seen as an approximation of the dFCR if
 254 threshold t^s was used for the CI calculations. Now find the smallest t^s that still results in $\text{dFCR} \leq \alpha$:
 255 $t_\alpha = \min\{t^s : F^s/m \leq \alpha\}$. This is the threshold used for the final calculation of the simultaneous
 256 confidence intervals and it will result in the control of the dFCR at the α level. In particular, the
 257 CIs are computed as $[L_i, U_i]$ with

$$L_i = E_i - t_\alpha s_i \quad \text{and} \quad U_i = E_i + t_\alpha s_i.$$

258 *Empirical evaluation of the methodology*

259 The formal method for synergy testing, which aims at controlling the dFCR, deviates from what
 260 was presented earlier in our papers [12, 17] and hence a thorough evaluation of the new methodology
 261 is needed, particularly for an assessment in terms of dFCR control and sensitivity. Since this paper
 262 is meant to be a tutorial, we have decided to move the details of the simulation study to supporting
 263 information, and report here only the main findings. Briefly, the simulation settings are adopted
 264 from the extensive simulation study of our previous work [12].

265 The results of the simulation studies demonstrate that our procedure succeeds in controlling
 266 the dFCR at the nominal level, while showing a sensitivity that is generally larger than what was
 267 obtained with our previous testing procedures (maxR and meanR). In the supporting information
 268 we give a more detailed discussion on the simulation results.

269 **4 Case study**

270 We will illustrate the described methodology and present data visualizations for synergy analysis
 271 using the BIGL R package and data from a drug combination experiment of direct-acting antivirals.
 272 Data are from the directAntivirals sample dataset included in the BIGL R package, consisting of 11
 273 drug combination experiments of direct-acting antivirals. To facilitate illustration, we will focus on
 274 the 4th experiment. We refer you to the supplementary material for R code implementation.

275 Measurements from experiments following the designs above often need to be normalized to the
 276 control wells. In this 4th experiment, the controls are the wells with 0 dose of both drug A and drug
 277 B. We chose to normalize the measurements by taking the ratio of each measurement to the average
 278 measurements of the control wells. We will refer to the normalized measurements as responses.

279 *Fitting monotherapies*

280 In the first step, the monotherapy curves for both drugs were estimated utilizing 4PL regression
 281 models (see Figure 4). Given the decreasing dose-response curves, a common upper asymptote
 282 assumption was made when fitting the model to both drugs. The 4PL regression models fit the data
 283 well (responses equally distributed above and below fitted line) and the sigmoidal patterns are well
 284 defined (i.e., two points defining upper asymptote, two points defining lower asymptote, and three
 285 points defining hill slope) with the selected doses. Notice that the lower asymptotes of Drug A and
 286 Drug B are not equal indicating a need for the Generalized Loewe null model.

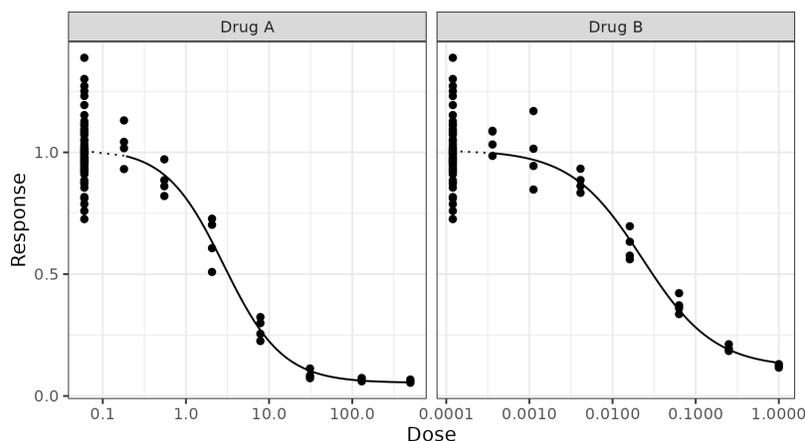


Figure 4: Monotherapy dose-response curves from directAntivirals dataset (BIGL R package), experiment 4. The y-axis is the response (normalized measurement) and the x-axis is the logarithmically transformed dose.

287 *Predicting expected combination response*

288 The expected responses for the chosen null models were calculated at all dose combinations,
 289 based on the estimated monotherapy dose-response curves. For illustration purposes, we used three
 290 different null models (the HSA, the Bliss, and the Generalized Loewe), enabling a sensitivity analysis
 291 of the different assumptions. In Figure 5, the observed and expected responses for the selected null
 292 models were visualized with stratifying the response surface by Drug A dose, creating a 2-dimensional
 293 trellis plot (the BIGL R package allows either Drug to be used as the stratifying variable). In this
 294 experiment, for a particular dose combination, points below the expected response are in the direction
 295 of synergy and points above the expected response are in the direction of antagonism. Furthermore,
 296 each null model has different underlying assumptions, as such, each null model predicts different
 297 expected treatment responses which is clearly visible in trellis 2 (top right facet).

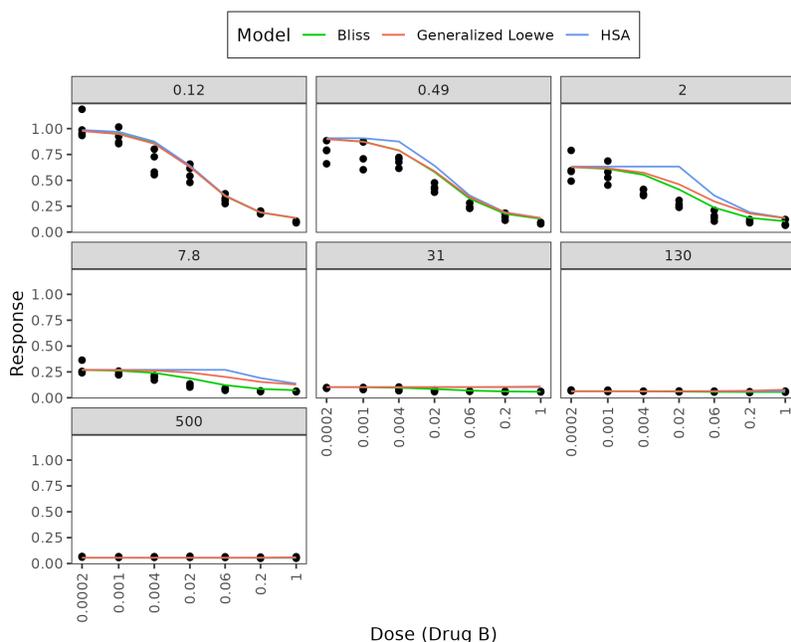


Figure 5: 2-Dimensional stratified predicted response surface plot, stratified by Drug A dose. The points are the observed responses whereas the colored lines are the expected responses derived from the different null models. The panels correspond to dose levels of drug A and the x-axis shows the dose levels of drug B.

298 *Estimating effect sizes*

299 The BIGL R package provides a plot to visualize the effect sizes along with the corresponding
300 confidence intervals for the null models of interest. Simultaneously, the plot allows users to visualize
301 the synergy or antagonism calls by color. Figure 6 represents the effect sizes and corresponding
302 CIs under the Generalized Loewe null model. In particular, a square highlighted in light blue
303 represents a synergy call at the corresponding dose combination. Should antagonism been detected,
304 the corresponding combination would have been highlighted in light pink. For example, synergy was
305 detected at the combination of 2, 0.016 (i.e., dose of Drug A = 2 and dose of Drug B = 0.016). The
306 observed response at this combination was 0.198 units lower than the expected response derived from
307 the Generalized Loewe null model with a 95% confidence interval of (-0.284, -0.111). In addition,
308 synergy/antagonism calls can also be presented in a bi-dimensional contour plot (see Figure 7). The
309 size of the point in the bi-dimensional contour plot represents the magnitude of the effect size. As
310 in Figure 6, the color of the bi-dimensional contour plot indicates the direction of deviance and 0
311 was not included in the 95% CI. The BIGL R package also allows for these data to be displayed
312 in a table format (supplementary material). All confidence intervals simultaneously control the dFCR
313 at 5%.

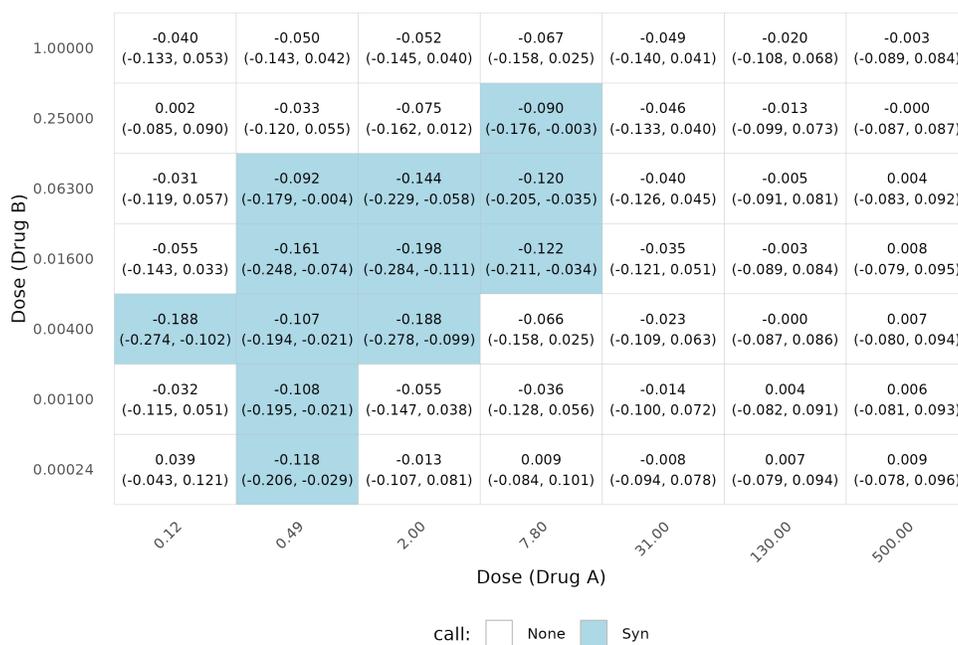


Figure 6: Effect sizes with simultaneous 95% confidence intervals (controlling the dFCR at 5%). Synergy calls are highlighted in light blue.

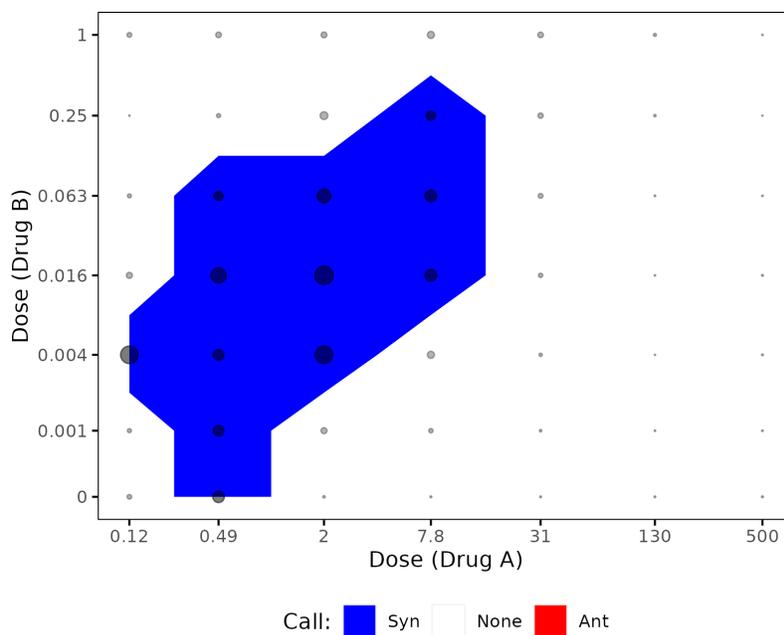


Figure 7: Contour plot. Synergy calls are highlighted in blue. The size of the points are relative to the magnitude of the effect sizes.

5 Discussion

The term synergy is used extensively to justify the potential of drug combinations. However, the meaning of it is often obscured as multiple null models exist, with most frequently used ones described in this tutorial paper. It is important, while evaluating the potential of combinations to keep the assumptions in mind and report accordingly. The HSA null model is the most liberal one, and in a technical sense not truly an evaluation of synergy as it only compares to the maximum response of either one of the monotherapies. Nevertheless, it can be an important evaluation where one would show the benefit of a combination versus a monotherapeutic effect. The other two null models, Bliss and Loewe, evaluate a particular type of additive effect. The Bliss null model is often the most conservative; however, the Bliss null model assumes the two drugs have a different mechanism of action. All these approaches are overly simplistic in that they fail to capture the underlying biological complexity. Hence, if no clear preference on the null model is present based on the assumptions, we suggest to run the models in parallel and evaluate the synergy calls with the assumptions in mind. Often, strong synergistic effects will be detected regardless of null model choice. When synergy is detected only under the HSA null model, this means that the combination only produces an increased effect when compared to a single monotherapy.

As mentioned before, it is important to evaluate the excess response versus the null models considering the variability in the data as observed deviations can be due to inherent variability of the experiment. Hence, our preference goes to methodologies evaluating the variability in the data, and not based on arbitrary thresholds (see Table 1). Additionally, it is important to express the excess response as an interpretable effect size. The latter made us shift from the maxR test in BIGL [17] to reporting the effect sizes together with CIs while controlling the directional false coverage rate (dFCR) as our default methodology. In an in-vitro screening setting, it is important to control the proportion of false calls to avoid validating irrelevant hypothesis downstream while keeping the sensitivity as high as possible. However, in these early screens, it is less crucial to have a precise estimate of the excess itself, as the expected gain of the combination will likely be evaluated downstream in more translatable experiments (e.g. in-vivo). This is why we chose to control the dFCR. In an extensive simulation study, we have demonstrated that our testing procedure controls the dFCR under a wide range of scenarios with acceptable sensitivity. To the best of our knowledge, the evaluation of error rate control, specifically dFCR, which is most important in early stages, is what discriminates BIGL from all other deviance methods described in literature; making it our preferred method.

To maintain a reasonable sensitivity, while controlling the dFCR, it is crucial to minimize the variability in the experiment which can be achieved during assay development and proper experimental design. Firstly, a capture of the full dose response of the monotherapies is required with sufficient data defining the hill slope. Often, the slopes are the areas where synergy could occur as asymptotes generally represent either no or full effect. Synergy calls, in neighbouring sampling points, are more convincing compared to singletons, as biologically, we expect some dynamic concentration ranges where compounds enhance each other. If variability can't be further reduced by identifying sources of variation and controlling key parameters in the assay protocol, it is crucial

354 to have enough replicates to ensure correct estimation of the variability in the experiment. From
355 experience, we learned that 4 replicates of the full checkerboard is often sufficient. However, research
356 is ongoing to evaluate the required number of replicates that results in adequate sensitivity to detect
357 synergy.

358 As cancer research is an important field in the exploration of combination therapies and its
359 growing attention in immunotherapies, we feel there is a need to dig deeper in the nature of the
360 replicates. The immune compartment in these type of experiments is often introduced using hu-
361 man donor material. Hence, replicates in this setting, cannot be considered as technical replicates.
362 Replicates derived from biological donors tend to enlarge the variability seen in the experiments. It
363 is crucial to disentangle these different sources of variability. Currently, we are exploring if expect-
364 ations, based on chosen null models can be calculated within a donor and the excess responses can
365 be pooled between donors for statistical testing.

366 Finally, the proposed methodology can be extended to three- or four-drug combinations by having
367 one or both, respectively, of the monotherapies be a drug combination. Under this scenario, the
368 EC_{50} s of the monotherapy are less interpretable; however, the methodology remains the same and
369 the effect sizes still represent excess responses.

370 **Acknowledgments**

371 Special thank you to Maxim Nazarov from Open Analytics who is the maintainer of the R-package
372 on CRAN. Additionally, we would like to thank Annelies Tourny and Stijn Hawinkel as they are
373 former developers and their simulation setup remained the foundation of the extensions presented
374 in this paper. Finally, we would like to thank the many scientists using our methodologies for the
375 usefull discussions which encouraged us to continuously improve the methods with clear applications
376 in mind.

377 **Conflicts of Interest**

378 Elli Makariadou, Xuechen Wang, Nicholas Hein and Bie Verbist are currently employed by Janssen
379 Pharmaceutical Companies of Johnson and Johnson. The authors declare no potential conflict of
380 interests.

381 **Data Availability**

382 We refer to the BIGL CRAN R package ([https://cran.r-project.org/web/packages/BIGL/](https://cran.r-project.org/web/packages/BIGL/index.html)
383 [index.html](https://cran.r-project.org/web/packages/BIGL/index.html)), which has an embedded directAntivirals dataset including 11 combination experi-
384 ments.

Supporting information

Simulation Study

We have conducted a simulation study using the same methodology as in our previous work [12]. More specifically, we have implemented scenarios 2 and 3, with slightly different parameter settings: the on-axis standard deviation σ_0 ranges from 0.05 to 0.2 and the numbers of replicates are 2, 3, 4 and 6. The results of the simulation study are also presented in the same way as in [12]: Appendix S1. In this html file the complete set of parameter settings is provided.

The simulation results are evaluated in terms of several criteria. Most of them are the conventional criteria: FDR, sensitivity, specificity, However, since we now aim to control the dFCR, we have included a few extra criteria.

The positive predicted value (PPV) and negative predicted value (NPV) refer to the expected proportion of true positive and true negative calls among the positive and negative calls, respectively. More formally, using the notation in Section 3, the PPV is defined as

$$\text{PPV} = \text{E} \left\{ \frac{\#\{i : 0 \notin [L_i, U_i] \text{ and } \tau_i \neq 0\}}{\#\{i : 0 \notin [L_i, U_i]\}} \right\}$$

and the NPV is given by

$$\text{NPV} = \text{E} \left\{ \frac{\#\{i : 0 \in [L_i, U_i] \text{ and } \tau_i = 0\}}{\#\{i : 0 \in [L_i, U_i]\}} \right\}.$$

We also introduced the concept of *neighbouring* in the evaluation of the method to better reflect a realistic use of the testing procedure. In practice the data analyst often looks at the dose combinations in the (d_1, d_2) plane for which zero is not contained in the CI (i.e. the method gives a synergy or antagonism call). Let us refer to such a point as a *positive point*. If such a positive point is isolated in the sense that no neighbouring points are positive, then many researchers will ignore this point and suspect it as a false positive. A neighbouring point is defined as a point that is within a certain distance in the (d_1, d_2) plane. For example, in a checkerboard design this distance can be chosen such that all closest points along the horizontal, vertical and diagonal directions are considered as neighbouring points. On the other hand, if a positive point has at least one positive neighbouring point, then scientist are willing to believe that these are true deviations from additivity and these points will be considered as *strong positive calls*. In the light of this reasoning we now define the *neighbouring PPV* as

$$\text{nPPV} = \text{E} \left\{ \frac{\#\{i : i \in \mathcal{N}^+ \text{ and } \tau_i \neq 0\}}{\#\{i : i \in \mathcal{N}^+\}} \right\},$$

where \mathcal{N}^+ is the set of positive off-axis points that have at least one positive neighbour.

In a similar fashion the *neighbouring false discovery rate* (nFDR) is given by

$$\text{nFDR} = \text{E} \left\{ \frac{\#\{i : i \in \mathcal{N}^+ \text{ and } \tau_i = 0\}}{\#\{i : i \in \mathcal{N}^+\}} \right\}.$$

413 Finally, we included criteria related to power. We defined *power3* as the probability to correctly
 414 detect at least three synergistic points,

$$\text{power3} = P \{ \#\{i : 0 \notin [L_i, U_i] \text{ and } \tau_i \neq 0\} \geq 3 \}$$

415 and *power_all* as the probability to detect all synergistic points,

$$\text{power_all} = P \{ \#\{i : 0 \notin [L_i, U_i] \text{ and } \tau_i \neq 0\} = \#\{\tau_i \neq 0\} \}.$$

416 R Code

417 Load the packages that are needed to generate the results. Set seed to get the same results.

```
418 library(BIGL)
419 library(ggplot2)
420 library(dplyr)
421 set.seed(1)
```

422 A function to subset data to a single experiment and, optionally, select the necessary columns
 423 only, and create the normalized measurements/responses.

```
424 subsetData <- function(data, i) {
425   subset(data, experiment == i)[, c("effect", "d1", "d2")] %>%
426   mutate(effect = effect / mean(effect[d1==0 & d2==0]))
427 }
```

428 Extract data of the 4th experiment.

```
429 data <- subsetData(directAntivirals, 4)
```

430 Step 1: Fit monotherapy dose-response models

```
431 mf <- fitMarginals(data, method = "nls", names = c("Drug A", "Drug B"))
```

432 Step 2: Predict expected combination responses using 3 different null models: the HSA, the Bliss
 433 and generalized loewe, controlling the dFCR at 5%.

```
434 rs <- list()
435 rs[["hsa"]] <- fitSurface(data, mf, null_model="hsa", statistic="both",
436   parallel=4, B.B=20, wild_bootstrap=TRUE,
437   wild_bootType="normal", control="dFCR")
438 rs[["bliss"]] <- fitSurface(data, mf, null_model="bliss", statistic="both",
439   parallel=4, B.B=20, wild_bootstrap=TRUE,
440   wild_bootType="normal", control="dFCR")
441 rs[["loewe"]] <- fitSurface(data, mf, null_model="loewe", statistic="both",
442   parallel=4, B.B=20, wild_bootstrap=TRUE,
443   wild_bootType="normal", control="dFCR")
```

444 Plot the monotherapy dose-response curves. (Figure 4)

```
445 plot(mf) + labs(x="Dose", y="Response")
```

446 Plot the 2-Dimensional stratified predicted response surface plot, stratified by Drug A dose.
447 (Figure 5)

```
448 synergy_plot_bycomp(rs, color = TRUE, plotBy = "Drug A",
449                    xlab = "Dose (Drug B)", ylab="Response")
```

450 Make the plot of effect sizes with simultaneous 95% confidence intervals. (Figure 6)

```
451 plotConfInt(rs[["loewe"]], color = "effect-size")
```

452 Make the contour plot. (Figure 7)

```
453 contour(rs[["loewe"]], colorBy = "effect-size", digits=3, main = NULL)
```

454 Print the effect sizes for all the combinations

```
455 rs[["loewe"]]$confInt$offAxis
```

456 References

- 457 1. Berenbaum MC. What is synergy? *Pharmacological reviews* 1989;41:93–141.
- 458 2. Bliss CI. The toxicity of poisons applied jointly 1. *Annals of applied biology* 1939;26:585–615.
- 459 3. Loewe S. Die quantitativen probleme der pharmakologie. *Ergebnisse der Physiologie* 1928;27:47–
460 187.
- 461 4. Vlot AH, Aniceto N, Menden MP, Ulrich-Merzenich G, and Bender A. Applying synergy met-
462 rics to combination screening data: agreements, disagreements and pitfalls. *Drug discovery*
463 *today* 2019;24:2286–98.
- 464 5. Harbron C. A flexible unified approach to the analysis of pre-clinical combination studies.
465 *Statistics in medicine* 2010;29:1746–56.
- 466 6. Straetemans R, O’Brien T, Wouters L, et al. Design and analysis of drug combination ex-
467 periments. *Biometrical Journal: Journal of Mathematical Methods in Biosciences* 2005;47:299–
468 308.
- 469 7. Prichard MN and Shipman Jr C. A three-dimensional model to analyze drug-drug interactions.
470 *Antiviral research* 1990;14:181–205.
- 471 8. Ianevski A, He L, Aittokallio T, and Tang J. SynergyFinder: a web application for analyzing
472 drug combination dose–response matrix data. *Bioinformatics* 2017;33:2413–5.
- 473 9. Ianevski A, Giri AK, and Aittokallio T. SynergyFinder 2.0: visual analytics of multi-drug
474 combination synergies. *Nucleic acids research* 2020;48:W488–W493.
- 475 10. Ianevski A, Giri AK, and Aittokallio T. SynergyFinder 3.0: an interactive analysis and con-
476 sensus interpretation of multi-drug synergies across multiple samples. *Nucleic Acids Research*
477 2022;50:W739–W743.

- 478 11. Zheng S, Wang W, Aldahdooh J, et al. SynergyFinder plus: toward better interpretation and
479 annotation of drug combination screening datasets. *Genomics, Proteomics and Bioinformatics*
480 2022;20:587–96.
- 481 12. Thas O, Tourny A, Verbist B, et al. Statistical detection of synergy: New methods and a
482 comparative study. *Pharmaceutical Statistics* 2022;21:345–60.
- 483 13. Di Veroli GY, Fornari C, Wang D, et al. Combeneft: an interactive platform for the analysis
484 and visualization of drug combinations. *Bioinformatics* 2016;32:2866–8.
- 485 14. Turner H, Tourny A, Thas O, et al. BIGL: Biochemically Intuitive Generalized Loewe Model.
486 R package version 1.7.0. 2023. URL: [https://cran.r-project.org/web/packages/BIGL/
487 index.html](https://cran.r-project.org/web/packages/BIGL/index.html).
- 488 15. Baud M. *Methods of Immunological Analysis 1: Fundamentals*. 1993.
- 489 16. Sebaugh J. Guidelines for accurate EC50/IC50 estimation. *Pharmaceutical statistics* 2011;10:128–
490 34.
- 491 17. Van der Borgh K, Tourny A, Bagdziunas R, et al. BIGL: Biochemically Intuitive General-
492 ized Loewe null model for prediction of the expected combined effect compatible with partial
493 agonism and antagonism. *Scientific Reports* 2017;7:17935.
- 494 18. Benjamini Y and Yekutieli D. False discovery rate-adjusted multiple confidence intervals for
495 selected parameters. *Journal of the American Statistical Association* 2005;100:71–81.
- 496 19. Benjamini Y, Hochberg Y, and Kling Y. False discovery rate control in pairwise comparisons.
497 Research Paper, Department of Statistics and Operations Research, Tel Aviv University 1993.
- 498 20. Liu RY. Bootstrap procedures under some non-iid models. *The annals of statistics* 1988;16:1696–
499 708.