

1 **Logic models to predict continuous outputs** 2 **based on binary inputs with an application to** 3 **personalized cancer therapy - SUPPLEMENT**

4 **Authors**

5 **Theo A. Knijnenburg**, Institute for Systems Biology, Seattle, US

6 **Gunnar W. Klau**, Centrum Wiskunde & Informatica, Amsterdam, The Netherlands

7 **Francesco Iorio**, European Molecular Biology Laboratory - European Bioinformatics
8 Institute, UK

9 **Mathew J. Garnett**, Wellcome Trust Sanger Institute, UK

10 **Ultan McDermott**, Wellcome Trust Sanger Institute, UK

11 **Ilya Shmulevich**, Institute for Systems Biology, Seattle, US

12 **Lodewyk F.A. Wessels**, Netherlands Cancer Institute, Amsterdam, and The Faculty of
13 EEMCS, Delft University of Technology, Delft, The Netherlands

14 **Abstract**

15 Mining large datasets using machine learning approaches often leads to models that
16 are hard to interpret and not amenable to the generation of hypotheses that can be
17 experimentally tested. We present ‘Logic Optimization for Binary Input to Continuous
18 Output’ (LOBICO), a computational approach that infers small and easily interpretable logic
19 models of binary input features that explain a continuous output variable. Applying LOBICO
20 to a large cancer cell line panel, we find that logic combinations of multiple mutations are
21 more predictive of drug response than single gene predictors. Importantly, we show that the
22 use of the continuous information leads to robust and more accurate logic models. LOBICO
23 implements the ability to uncover logic models around predefined operating points in terms
24 of sensitivity and specificity. As such, it represents an important step towards practical
25 application of interpretable logic models.

26 **Supplementary Information**

27 ***Supplementary Note 1 – Validation on the CTRP dataset***

28 We aimed to validate the logic models inferred by LOBICO on our cell line panel by
29 applying these logic models to the drug response data of another cell line panel: the Cancer
30 Therapeutic Response Portal version 2 (CTRP) ¹.

31 CTRP data was obtained from the supplementary information files of the
32 corresponding main publication, available online on the Cancer Discovery journal web-site
33 at: [http://cancerdiscovery.aacrjournals.org/content/early/2015/10/14/2159-8290.CD-15-](http://cancerdiscovery.aacrjournals.org/content/early/2015/10/14/2159-8290.CD-15-0235/suppl/DC1)
34 [0235/suppl/DC1](http://cancerdiscovery.aacrjournals.org/content/early/2015/10/14/2159-8290.CD-15-0235/suppl/DC1) (Supplemental Tables S1 – S7, file: 145780_2_supp_3058746_nrhtdz.xlsx).
35 From these files, identifiers of screened cell lines and compounds were extracted and
36 mapped to the cell lines and compound identifiers of our study (from now called GDSC for
37 Genomics of Drug Sensitivity in Cancer). In total, there were 47 overlapping drugs between
38 GDSC and CTRP. For CTRP, the drug response indicator is the AUC, i.e. the area under the
39 drug/cell-line dose response curve. IC50 values were not available for the CTRP study.

40 Across GDSC and CTRP there are 344 cell lines available in both panels, and 370 cell
41 lines only available in GDSC. We explored two scenarios that both involved a training cohort
42 and a validation cohort: 1) **Train:GDSC344-Validate:CTRP344** LOBICO models were trained
43 on GDSC data of the 344 cell lines (GDSC344) and validated on CTRP data of the same set of
44 344 cell lines (CTRP344), and 2) **Train:GDSC370-Validate:CTRP344** LOBICO models were
45 trained on GDSC data of the 370 cell lines (GDSC370) and validated on CTRP data of the
46 independent set of 344 cell lines (CTRP344). See **Supplementary Figure 2** for an overview of
47 these datasets.

48 For each of the two scenarios, we ran LOBICO across the 47 compounds on GDSC
49 using the same settings as for the original analysis. For each drug we selected the best model
50 according to cross-validation (CV) and applied that model to the cell lines dividing them into
51 a group that is predicted to be sensitive and a group that is predicted to be resistant. Then,
52 we performed a t-test comparing these two groups both for the GDSC IC50s as well as for
53 the CTRP AUCs. We also performed a t-test for the best single predictor model. The t-tests
54 were only performed when the groups consisted of at least 5 cell lines. This led to 37 and 39
55 drugs for scenarios 1 (**Train:GDSC344-Validate:CTRP344**) and scenario 2 (**Train:GDSC370-**
56 **Validate:CTRP344**), respectively, that we could test within this framework.

57 For scenario 1 (**Train:GDSC344-Validate:CTRP344**) the best performing models on
58 GDSC validated on CTRP with high statistical significance (**Supplementary Figure 3**). Overall,
59 the p-values for the t-tests on the GDSC IC50s and CTRP AUCs for the 37 models were
60 substantially correlated (**Supplementary Figure 4**, Pearson correlation: 0.94 , $p = 2.17 \times 10^{-18}$,
61 Spearman correlation: 0.31 , $p = 0.06$). Note that lower Spearman correlation indicates that
62 the correlation is mostly driven by the strong models in GDSC and also showed strong
63 validation in CTRP.

64 Selection of significant t-test p-values using a strict p-value threshold of $0.01/37$ (a
65 Bonferroni corrected p-value of 0.01) showed significant overlap between GDSC and CTRP
66 (Fisher exact test, $p = 1.3 \times 10^{-3}$).

67 Using a somewhat loose threshold of $1/37$ (a family-wise error rate of 1) we
68 identified that 18 of the 37 (49%) drugs led to statistically significant models in the GDSC
69 training cohort, i.e. for these models the cell lines predicted to be sensitive and resistant
70 showed differential drug response using the t-test. 5 of these 18 models (28%) also showed
71 statistical significance in the CTRP344 validation cohort. Importantly, we found that in many
72 cases multi-predictor models outperformed single predictor models (65% for the training
73 cohort and 51% for the validation cohort). See **Table 2** for an overview.

74 For scenario 2 (**Train:GDSC370-Validate:CTRP344**) we found that the two best
75 performing models on GDSC validated with high statistical significance in CTRP
76 (**Supplementary Figure 5**). Yet, many other good models on GDSC did not lead to a strong
77 prediction of drug response in CTRP; this was the case both for the multi-predictor models
78 and single predictor models inferred from GDSC370. Overall, the p-values for the t-tests on
79 the GDSC IC50s and CTRP AUCs for the 39 models were correlated (**Supplementary Figure 6**,
80 Pearson correlation: 0.84 , $p = 3.3 \times 10^{-18}$). This correlation is mostly driven by the top 2 as
81 evidenced by the much lower Spearman correlation of 0.17 , $p = 0.31$.

82 Related to this, selection of significant t-test p-values using the strict p-value
83 threshold of $0.01/39$ (a Bonferroni corrected p-value of 0.01) showed only a moderately
84 significant overlap between GDSC and CTRP (Fisher exact test, $p = 0.06$).

85 Using the more loose threshold of $1/46$ for the training cohort and $1/39$ for the
86 validation cohort (a family-wise error rate of 1) we identified that 25 of the 46 (54%) drugs
87 led to statistically significant models in the GDSC training cohort, i.e. for these models the
88 cell lines predicted to be sensitive and resistant showed differential drug response using the

89 t-test. 5 of these 25 models (20%) also showed statistical significance in the CTRP344
90 validation cohort. Again, we found that in many cases multi-predictor models outperformed
91 single predictor models (74% for the training cohort and 31% for the validation cohort). Of
92 the 5 validated models, 3 were multi-predictor models. Again, see **Table 2** for an overview.

93 A comparison of scenario 1 and 2 shows that for the independent validation cohort
94 (scenario 2), fewer of the multi-predictor models inferred on the training cohort are more
95 predictive than the best single predictor model (31% in scenario 2 vs. 51% in scenario 1, and
96 60% for statistically significant models in scenario 2 vs. 80% in scenario 1). This is an
97 indication that, for a considerable number of some drugs, the multi-predictor LOBICO
98 models inferred on one set of cell lines do not generalize to another set of cell lines. Yet,
99 overall these results are encouraging, also given that the sets of 370 and 344 cell lines are
100 substantially different in terms of tissue types and mutation landscape (**Supplementary**
101 **Figure 2**).

102 ***Supplementary Note 2 – Robustness across CV folds***

103 We investigated the robustness of the logic models across the ten CV training folds
104 for each of the 142 drugs. The logic models for a drug were inferred using the model
105 complexity (defined by K and M) selected by CV for that drug in the standard setting, i.e.
106 with the sample-specific weights and $t=0.05$. The use of the continuous output resulted in a
107 smaller variation in the FI scores across the CV folds (**Supplementary Figure 8a**). Particularly,
108 the FI scores for the logic models across the CV folds had an average Pearson correlation
109 coefficient larger than 0.75 for 113 drugs (80%), and 40 drugs (28%) had a correlation larger
110 than 0.95. In contrast, for the logic models based on binarized data, there were 89 drugs
111 (63%) had a correlation larger than 0.75 and 35 (25%) had a correlation larger than 0.95.

112 In comparison to changing the binarization threshold (**Figure 3a**), we observed that
113 logic models inferred from the randomly sampled subsets, i.e. the CV folds, showed more
114 variability in the FI scores, and thus a smaller correlation amongst the CV folds. We
115 hypothesized that the inclusion or exclusion of samples, especially those far away from the
116 binarization threshold, can have a large effect on the optimization function (Equation 1), and
117 therefore a large effect on the inferred optimal logic model and the resulting FI scores. To
118 test this hypothesis, we compared the similarity between CV folds with the similarity of the

119 FI scores derived from the logic models trained on these CV folds. Specifically, for each of the
120 142 drugs separately, we computed:

- 121 1. for each pair of the CV training folds, say a and b , the similarity between the
122 CV folds a and b in the following manner:
 - 123 a. We took \mathbf{w} , the $N \times 1$ continuous vector with weights for each of the
124 N samples. \mathbf{w} is the absolute difference between the IC50s and the
125 binarization threshold, normalized per class (Equation 14).
 - 126 b. We created \mathbf{w}^a and \mathbf{w}^b , where \mathbf{w}^a is identical to \mathbf{w} , except that all
127 samples that are not part of the training set of a are replaced by 0,
128 and similarly for \mathbf{w}^b .
 - 129 c. As a metric of the similarity between CV folds a and b , we computed
130 the Pearson correlation coefficient between vectors \mathbf{w}^a and \mathbf{w}^b .
- 131 2. for each pair of the CV training folds, as a metric of the similarity of the FI
132 scores between the two members a and b , the Pearson correlation
133 coefficient between the FI score vectors derived from the logic models trained
134 on CV folds a and b .

135 With the 142 drugs and 10-fold CV strategy, this resulted in $142 \times \binom{10}{2} = 6390$

136 pairwise correlation scores for the similarity in the weight vectors and 6390 pairwise
137 correlation scores for the similarity in the FI scores. We observed a clear relationship
138 between these correlation scores (**Supplementary Figure 8b**). Particularly, pairs of CV folds
139 with a small correlation between the weight vectors often had a small correlation between
140 the FI scores. We observed that in about 5% of the cases the correlation between the weight
141 vectors was quite low, i.e. the correlation coefficient was smaller than 0.7. These are cases,
142 where the two CV folds include (and exclude) different samples with extreme IC50s, i.e.
143 those far away from the binarization threshold. These ‘important’ samples have large
144 weights in the weight vector, and when set to 0 in one of the folds, but not the other, lead to
145 the low correlation scores between the folds. It is thus not surprising that the logic models
146 inferred on these distinct CV folds lead to different FI scores.

147 This analysis confirmed our hypothesis that the larger variation in FI scores observed
148 across CV folds is due to the inclusion or exclusion of samples with large weights, i.e. those
149 far away from the binarization threshold.

150 ***Supplementary Note 3 – Subsampling analysis***

151 Specifically, we randomly sampled from all 714 cell lines 90% to 1% of the cell lines
152 in 13 steps, i.e. 90%, 80%, 70%, 60%, 50%, 40%, 30%, 20%, 15%, 10%, 5%, 2%, 1%, and
153 repeated this 10 times. Then, for each case we ran LOBICO across all 142 drugs using the
154 original settings. We analyzed the CV errors and the FI scores across the repeats and
155 compared the results to the setting when we use all (100%) of the cell lines. Based on the CV
156 errors, FI scores and the permutation test (Methods section), we employed the following 5
157 criteria to identify ‘robust’ and ‘predictive’ models. All criteria must be met in order to call a
158 model robust and predictive.

- 159 1. The number of sensitive cell lines must be larger than 10 – This is prerequisite for
160 running LOBICO using 10-fold CV and is also the minimum number of cell lines in
161 the positive (sensitive class). In our cell line panel, for most drugs, the bulk of cell
162 lines are not affected, and only a small percentage (5-15% typically) end up in the
163 class of sensitive cell lines. When subsampling to 50% (~300 cell lines) still 135
164 (95%) of the models can be run. When sampling 20% and 10% of the cell lines
165 (~125 and 65 cell lines resp.) LOBICO models can only be run for 77 (54%) and 23
166 (16%) of the models. See **Supplementary Figure 9a**.
- 167 2. The CV error must be smaller than 0.4 – The statistical cutoff of $FDR < 1\%$ and
168 $p < 0.01$ that we used to identify statistically significant logic models using the
169 permutation test coincided with a CV error of approximately 0.4 (see **Figure 2**).
170 Therefore, we used this cutoff to identify predictive models. Without
171 subsampling, i.e. using all cell lines, 72 (51%) of the models are predictive. This
172 percentage remains more or less constant when subsampling. See
173 **Supplementary Figure 9b**.
- 174 3. The error across the complete dataset must be smaller than 0.4 – The optimal
175 logic model for each of the subsampling rates and repeats was applied to all cell
176 lines after which the error across the complete cohort was computed. We
177 required this error to be smaller than 0.4 in which case the logic model inferred

178 using a subset of the samples generalizes over the complete dataset. We
179 observed that the although CV-error remains constant, the error on the complete
180 dataset increases with a smaller subsampling frequency. This is an indication that
181 when using a smaller set of cell lines the inferred logic models do a good job of
182 explaining the drug response for those cell lines, but these models do not
183 generalize across a larger panel. For example, when sampling 10% of the cell lines
184 (~65 cell lines) the percentage of predictive models drops to ~30%. See
185 **Supplementary Figure 9c.**

186 4. The Pearson correlation of FI scores amongst the CV folds must be higher than 0.7
187 – We observed that the Pearson correlation coefficients of the similarity of FI
188 scores across the 10 CV folds on the complete dataset were higher than 0.7 for
189 most drugs (**Supplementary Figure 8**). Therefore, we used this cutoff to identify
190 robust models in the subsampling analysis. Without subsampling, i.e. using all cell
191 lines, 124 (87%) of the models are robust according to that definition. This
192 percentage remains more or less constant when subsampling. See
193 **Supplementary Figure 9d.**

194 5. The Pearson correlation of FI scores between the subsampled and complete
195 dataset must be higher than 0.7 – The FI scores obtained for the logic models for
196 each of the subsampling rates and repeats were correlated with the FI scores
197 from the logic models inferred across the complete cohort. We required this
198 correlation to be larger than 0.7 in which case the logic models inferred using a
199 subset of the samples are highly similar to the original logic model based on all
200 samples. We observed that correlation of these FI scores decreased quickly with a
201 smaller subsampling frequency. This is an indication that when using a smaller set
202 of cell lines the inferred logic models are different from the original model. For
203 example, when sampling 50% and 20% of the cell lines (~300 and 125 cell lines
204 resp.) 55 (41%) and 32 (28%) logic models met this criterion. See **Supplementary**
205 **Figure 9e.**

206 Overall, we observed that subsampling had a substantial influence on performance and
207 robustness as defined by our criteria. Specifically, whereas for 51% of the drugs LOBICO
208 inferred logic models that are robust and predictive when using all cell lines, this number
209 decreased to 25% of the drugs when using 50% of the cell lines (around 300 cell lines). With

210 a subsampling frequency of 10% (around 60 cell lines) only 4 (18%) of the drugs had a robust
211 and predictive model. Perhaps not surprising, these 4 drugs were amongst the ones with the
212 lowest CV-error on the complete dataset.

213 In conclusion, LOBICO can also effectively be run on smaller sets of cell lines, but in our panel
214 there was only a relatively small number of drugs for which robust and predictive models
215 were found with much smaller sets of cell lines. The subsampling analysis led to two
216 important insights: First, it is important that the classes in the dataset are not strongly
217 unbalanced; if one of the two classes is too small this leads to non-robust models or even
218 the inability to run LOBICO using CV. Second, only for the highly predictive models, i.e. those
219 with a small CV error, did we find robust and predictive models when using a small set of cell
220 lines. Thus, for smaller sets of cell lines strong effect sizes are necessary to reach
221 significance. This is something that can potentially be tested with univariate tests before
222 running LOBICO.

223 ***Supplementary Note 4 – LOBICO on a yeast cross phenotyped for*** 224 ***sporulation efficiency***

225 We re-analyzed the genetic linking map of a cross of two natural yeast strains, a
226 strain isolated from the bark of an oak tree that sporulates at 99% efficiency, and a strain
227 originating from a wine barrel that sporulates at only 3.5%². The genetic linkage map
228 consists of 225 loci genotyped in 374 segregants. For each of the 374 recombinant offspring,
229 the sporulation efficiency was measured as a percentage between 0 and 100. Gerke *et al.*²
230 used composite interval mapping based on a stepwise regression model to find loci that
231 significantly cosegregated with variation in sporulation efficiency, leading to 5 significant
232 loci, L7-9, L10-14, L13-6, L7-17 and L11-2 (Table 1 in²). Next, a second stepwise regression
233 was used to select significant predictors from the five loci and all 2 and 3-way interaction
234 terms involving these five loci. The final model included three significant 2-way interaction
235 effects and one 3-way interaction effect. All these interactions were comprised of
236 combinations of the three most significant individual loci, i.e. L7-9, L10-14 and L13-6 (Table
237 S2 in²).

238 We applied LOBICO to this dataset to evaluate which (logical) interaction effects
239 would be uncovered. The genotype information was straightforwardly transformed into
240 binary predictor variables: Alleles from the oak strain (wine strain) were set to 1 (0),

241 resulting in a truth table with $n=225$ loci and $p=374$ segregants. The sporulation
242 phenotype data was binarized by applying a threshold of 50%. Samples were weighted using
243 the distance to this threshold. No specificity and sensitivity constraints were applied. We
244 employed the eight model complexities also used for the cell line panel analysis, i.e. all
245 combinations of K and M with $K \cdot M \leq 4$.

246 The largest single effect found in Gerke *et al.*, loci L7-9, is also the best single
247 predictor uncovered by LOBICO (**Supplementary Figure 11**). The two-input AND model found
248 by LOBICO consisted of loci L7-9 and L10-14. This interaction, which is also one of the 2-way
249 interaction effects found in Gerke *et al.* has a much higher specificity and precision than the
250 single locus model, although a smaller recall. Many of the offspring with the highest
251 sporulation efficiency have both the L7-9 and the L10-14 locus from the oak strain. The best
252 model according to CV is a 2-by-2 model, which contains the same three loci as in the
253 interaction effects found in Gerke *et al.*, i.e. L7-9, L10-14 and L13-6. (Actually, the LOBICO 2-
254 by-2 model contained L13-7 instead of L13-6; they are highly correlated. The fourth feature
255 in the 2-by-2 is L7-11, which is highly correlated to L7-9.) Thus, LOBICO finds interactions
256 between the same three loci as the regression model employed by Gerke *et al.*.

257 It is important to point out that LOBICO uncovered these interactions using the
258 complete dataset of 225 loci, and not by first filtering on individual features as was done in
259 Gerke *et al.*. Surely, the (biological) interpretation of the logic model and the additive linear
260 model is quite different. We would argue that the logic model is more intuitive and sensible
261 than the linear model.

262 ***Supplementary Note 5 – Explanation of the Boolean Function*** 263 ***Synthesis Problem and proof that LOBICO is NP-complete***

264 The Boolean Function Synthesis Problem (BFSP) is a particular type of Boolean
265 Satisfiability Problem, where the goal is to find an algebraic sum-of-products expression for
266 an incompletely specified Boolean function $\Phi : \{0,1\}^n \rightarrow \{0,1\}$. The sum-of-products
267 expression is also called a disjunctive normal form (DNF), i.e. a disjunction of conjunctions.
268 Each Boolean function can be expression in DNF. An element of the domain of Φ is called a
269 minterm of Φ . The set of minterms for which Φ evaluates to 1 (resp. 0) is called the
270 ON-set (resp. OFF-set). An incompletely specified Boolean function is one for which

271 $|\text{ON-set}| + |\text{OFF-set}| < 2^n$. **Supplementary Figure 17** displays an incompletely specified
272 Boolean function with $n = 10$ input variables, x_1, x_2, \dots, x_{10} and an output variable y .

273 The number of rows in the Boolean truth table is given by p (
274 $p = |\text{ON-set}| + |\text{OFF-set}|$), and is 40 in this case ($40 \ll 2^{10}$). Note that in most biology
275 applications, Φ is incompletely specified. The sought after algebraic expression is a Boolean
276 DNF expression that evaluates to 1 for all minterms in the ON-set (OFF-set) and to 0 for all
277 minterms in the OFF-set. Formally, the problem is as follows: Given an ON-set and an
278 OFF-set of minterms that characterize a Boolean function Φ , find a DNF of Φ with
279 maximally K disjunctive terms having each maximally M variables. The corresponding
280 decision problem is NP-complete³.

281 The decision version of LOBICO is as follows: Given inputs \mathbf{X} , \mathbf{y} , \mathbf{w} , K, M and a
282 parameter L , does there exist a logic function $\hat{\Theta}$ expressed in $\text{DNF}(K, M)$, i.e. a DNF having
283 at most K disjunctive terms and M literals, such that the weighted sum of incorrectly
284 inferred samples as described in Equation 1 is less than or equal to ε ? Clearly, the problem is
285 in NP. It is easily shown that this problem is also NP-complete by the following polynomial-
286 time reduction from BFSP: Given an instance of BFSP we construct an instance for LOBICO by
287 deriving \mathbf{X} from the minterms, \mathbf{y} from the ON-set and OFF-set and by setting \mathbf{w} to $\mathbf{1}$, i.e.
288 $w_n = 1 \quad \forall n$. Now, BFSP can be satisfied if and only if LOBICO has a solution with an error of
289 $\varepsilon = 0$. Since the BFSP decision problem is NP-complete, the LOBICO decision problem is also
290 NP-complete.

291 ***Supplementary Note 6 – Comparison with logic regression***

292 Logic regression (LR)^{4,5} is a generalized regression methodology that can be applied
293 to data with binary predictors, although continuous predictors are also allowed. The goal of
294 LR is to find linearly weighted logic combinations of the original predictors that explain a
295 continuous response variable or class label. We configured the implementation of LR, i.e. the
296 R-package ‘logreg’, such that it infers logic models with a predefined model complexity.
297 Specifically, logreg has a scoring function for classification using sample-specific weights,
298 which we used to give it the same objective function as LOBICO (Equation 1). Also, logreg

299 can be configured to output a single logic model (a tree) with predefined logical operators
300 and size. (See below for experimental details.)

301 We ran LR for each of the 142 drugs in the cancer cell line panel using the model
302 complexity (defined by K and M) selected by CV for the associated drug when using LOBICO.
303 LR was run on the same computers (Intel(R) Xeon(R) CPU, E5645, 2.40GHz, 6 cores) as
304 LOBICO and was given the same amount of CPU time (**Supplementary Figure 18a**). Then, we
305 evaluated the logic formulas inferred by LR. Specifically, we looked at the Jaccard similarity
306 of the selected predictors in the inferred LOBICO and LR models. (In computing the Jaccard
307 similarity negated terms, e.g. \neg TP53, are treated as separate predictors from their positive
308 equivalents.) For models with $K=1$ and/or $M=1$, a Jaccard similarity of 1 indicates that the
309 exact same logic formula was found. For $K=2$ and $M=2$ (the 2x2 models) this is not
310 necessarily the case, but we were not able to restrict logreg to output a DNF with $K=2$ and
311 $M=2$ anyway. For example, for the drug 'MG-132' LOBICO inferred the 2x2 ' $(\neg$ MYC & RB1) |
312 $(\neg$ PIK3CA & \neg TP53)', whereas LR inferred ' $((\neg$ TP53) or (RB1 or NOTCH1)) and $(\neg$ PIK3CA)'.
313 The LR model is clearly not a DNF with $K=2$ and $M=2$.

314 Overall, LR found the same (optimal) logic formulas as LOBICO (**Supplementary**
315 **Figure 18b**). The main exception is the 2x2 model ($K=2$, $M=2$), but this is because of the
316 reason mentioned above. For the 4-input OR models ($K=4$, $M=1$) we observed four cases
317 where the logic formulas differed between LOBICO and LR. Upon further inspection, we
318 found that in these cases the formula inferred by LR had the same (optimal) error as the
319 LOBICO solution. These LR solutions were present in LOBICO's solution pool, i.e. they were
320 part of the set of (sub-)optimal solutions output by LOBICO. (See experimental details
321 below.)

322 In conclusion, when logreg parameters are properly set, LR can find the optimal
323 solution when given the same amount of time that was necessary for LOBICO to find the
324 optimal solution on the cancer cell line dataset. Potentially, LR finds this solution faster than
325 LOBICO on this dataset. It is however important to point out that LR cannot guarantee that
326 the obtained solution is optimal, and it is known that ILP solvers spend a long time proving
327 that the found solution is indeed optimal. In future work, we will investigate the
328 performance of LR and LOBICO on other (larger) datasets, and assess how the two methods
329 can be used in parallel to find optimal solutions faster. For example, we will investigate
330 whether LR can be used to identify initial starting models for LOBICO.

331 Importantly, LR cannot incorporate statistical performance constraints, such as
332 sensitivity and specificity (Equations 10 - 13), which we assert is the preferred and, in
333 practice, most relevant scenario for LOBICO inferences. Additionally, in contrast to LR,
334 LOBICO can output the pool of (sub-)optimal solutions, which we used to measure feature
335 importance.

336 **Experimental details:** LR was run using the R-packing logreg (version 1.5.8). We
337 used simulated annealing as this search algorithm gave the best results. We followed the
338 logreg's documentation to set the upper and lower temperature of the annealing chain
339 based on experiments with the cancer cell line panel. The number of iterations was set, such
340 that the total CPU time spent on solving the problem was comparable to the CPU time that
341 LOBICO needed to find the optimal solution (**Supplementary Figure 18a**). The simulated
342 annealing parameters were set as follows (R-code):

```
343           myanneal <- logreg.anneal.control(start = 1, end = -5, iter =  
344 T*200000, update = T*20000)
```

345 To infer a LR model with the same model complexity as LOBICO, we made sure that
346 for 2-, 3- and 4-input AND models only AND operators were allowed. Similarly, for 2-, 3- and
347 4-input OR models only OR operators were allowed. For 2x2 models we allowed both AND
348 and OR operators. The parameters of the logic 'tree' shape were set as follows (R-code):

```
349           if (K>M) mytreecontrol <- logreg.tree.control(opers=3) else  
350 mytreecontrol <- logreg.tree.control(opers=2)  
351           if (K==2&M==2) mytreecontrol <- logreg.tree.control(opers=1)
```

352 LR was run to output one logic tree (ntrees=1), where the maximum number of
353 leaves was set to $K \times M$ (nleaves= $K \times M$). In the R-code below Y is \mathbf{y} , X is \mathbf{X} and W is \mathbf{w} as
354 used in the **Methods Section** and Equation 1. LR was run as follows (R-code):

```
355           q<-logreg(resp=Y, bin=X, wgt=W, type=1, select=1, ntrees=1,  
356 nleaves=K*M,anneal.control = myanneal, tree.control = mytreecontrol)
```

357 ***Supplementary Note 7 – Comparison with sparse linear regression*** 358 ***and Random Forests***

359 We compared the LOBICO models obtained on the cancer cell line panel with Elastic
360 Net ⁶, a sparse linear regression model, and with Random Forests regression ⁷, a non-linear
361 regression model. Specifically, for the 25 drugs with the lowest CV error in the original
362 analysis, we compared the model-specific FI scores (of the model complexity selected by CV)

363 with the regression weights inferred by Elastic Net (EN) and the importance scores inferred
364 by Random Forests (RF).

365 We observed a large concordance between LOBICO's FI scores and the EN regression
366 weights (**Supplementary Figure 12a**). In the EN models, the large majority (63%) of all
367 regression weights across the 60 features and 25 drugs were 0. Importantly, all of the
368 important features according to LOBICO (FI>0.05) had a non-zero regression weight in EN.
369 Moreover, the smallest EN regression weight for which the corresponding LOBICO FI was
370 larger than 0.05, was 0.2752 (blue line in **Supplementary Figure 12a**), which was in the tail of
371 the EN weights.

372 Similarly for RF, we observed a high degree of correlation between LOBICO's FI scores
373 and the RF importance scores (**Supplementary Figure 12b**). The important features
374 according to LOBICO (FI>0.05) also had a high RF importance score. The smallest RF
375 importance score for which the corresponding LOBICO FI was larger than 0.05, was 0.014,
376 which marked the 82% percentile of the RF importance scores.

377 We note that it is not possible to do a direct comparison in terms of predictive
378 performance, because LOBICO and EN/RF use different error measures. That is, LOBICO's
379 error is the weighted sum of incorrectly inferred samples, where the weight is the distance
380 from a cell line's IC50 to the discretization boundary - an L1 norm for misclassified samples.
381 This error is different from the L2 norm (least squares) or L1 norm across all samples used in
382 (RF) regression models. At the same time, LOBICO's goal is not to be better than other
383 methods in terms of prediction performance, but instead to create interpretable models
384 with good performance.

385 **Experimental details:** For EN, we used the MATLAB 'lasso' function with an alpha
386 (mix between L2 and L1 penalty) of 0.5, 10-fold CV and sample weights \mathbf{w} , the $N \times 1$
387 continuous vector with weights for each of the N samples (Equation 14). For RF, we
388 employed the Random Forests implementation for MATLAB v0.02 downloaded from
389 <http://code.google.com/p/randomforest-matlab/>. The RF regression models were run with
390 1000 trees each and default settings for the other parameters were used. The reported
391 importance scores represent the mean decrease in accuracy. To accommodate the different
392 sample weights, we created (for each drug) a dataset of 10,000 samples, which were
393 randomly drawn with replacement from the original dataset, where the probability of being
394 drawn was proportional to the sample weights in \mathbf{w} .

395 **Citations**

396 1 Seashore-Ludlow, B. *et al.* Harnessing connectivity in a large-scale small-molecule sensitivity
397 dataset. *Cancer discovery* **5**, 1210-1223 (2015).

398 2 Gerke, J., Lorenz, K. & Cohen, B. Genetic interactions between transcription factors cause
399 natural variation in yeast. *Science* **323**, 498 (2009).

400 3 Gimpel, J. F. A method of producing a Boolean function having an arbitrarily prescribed prime
401 implicant table. *Electronic Computers, IEEE Transactions on*, 485-488 (1965).

402 4 Ruczinski, I., Kooperberg, C. & LeBlanc, M. Logic regression. *Journal of Computational and*
403 *Graphical Statistics* **12**, 475-511 (2003).

404 5 Kooperberg, C. & Ruczinski, I. Identifying interacting SNPs using Monte Carlo logic regression.
405 *Genetic epidemiology* **28**, 157-170 (2005).

406 6 Zou, H. & Hastie, T. Regularization and variable selection via the elastic net. *Journal of the*
407 *Royal Statistical Society: Series B (Statistical Methodology)* **67**, 301-320 (2005).

408 7 Breiman, L. Random forests. *Machine learning* **45**, 5-32 (2001).

409

410

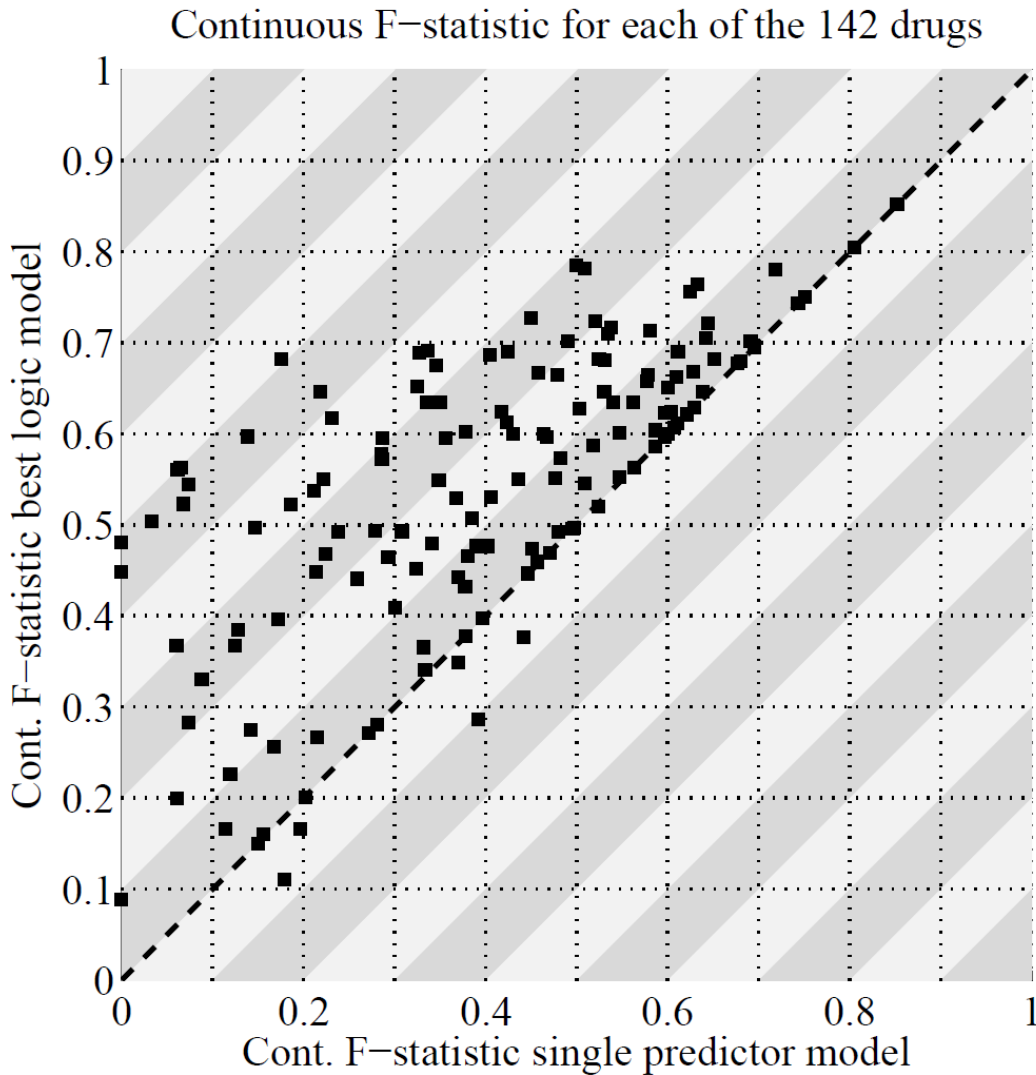
411

412

413

414 **Supplementary Figures**

415



417
418

419 **Supplementary Figure 1 | Multi-predictor models outperform single predictor models**

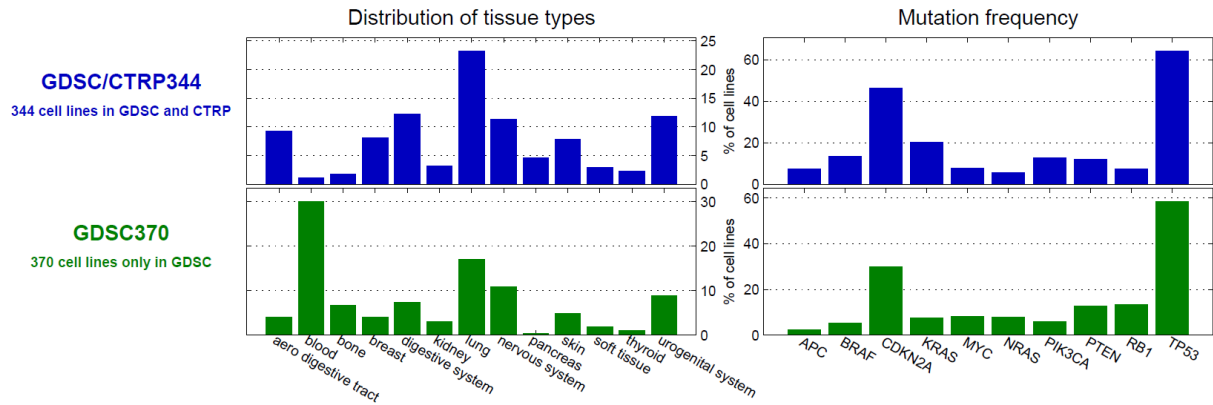
420 Scatter plot with the continuous F-statistic for single predictor models (x-axis) and the best
421 (lowest CV error) multi-predictor model (y-axis). Each point represents one of the 142 drugs.
422 The continuous F-statistic is defined as the harmonic mean of the continuous recall and
423 continuous precision. By analogy to Equations 12 and 13, the continuous recall and

424 continuous precision are defined as $R_c = \frac{\sum_{n=1}^N (w_n \cdot y'_n)}{\sum_{n=1}^N (w_n \cdot y_n)}$ and $P_c = \frac{\sum_{n=1}^N (w_n \cdot y'_n)}{\sum_{n=1}^N (w_n \cdot y'_n)}$. The

425 continuous F-statistic uses the sample-specific weights that LOBICO uses in its optimization
426 and is therefore a better performance measure than the standard F-statistic. Similarly, to the
427 CV error depicted in **Figure 2**, the continuous F-statistic was computed on the inferred class
428 labels of the samples in the test sets.

429
430

431



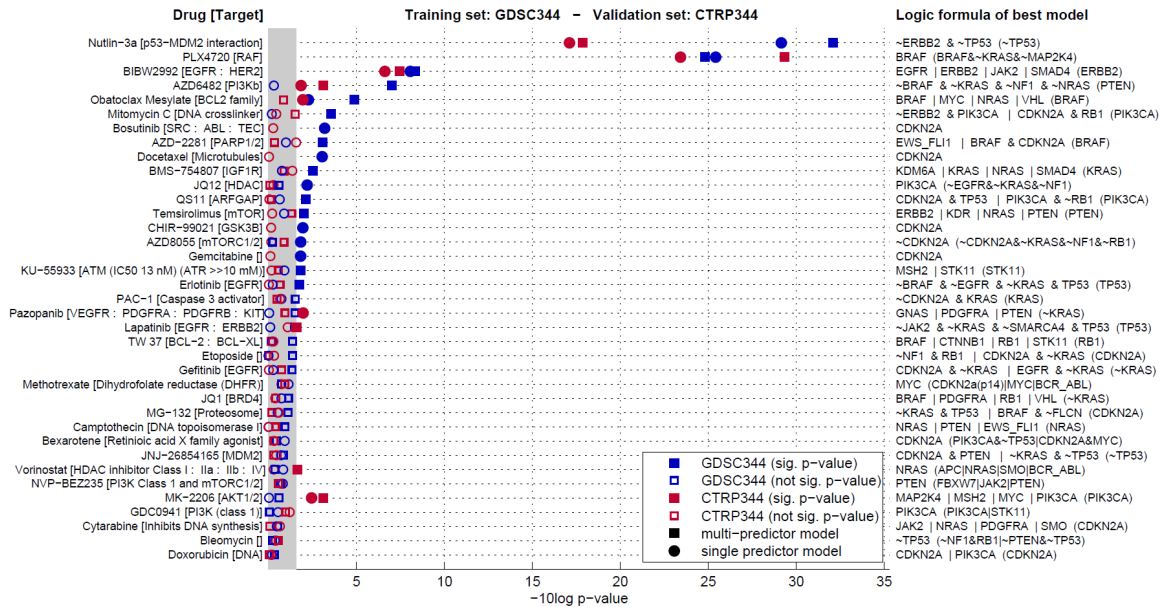
432

433

434 **Supplementary Figure 2 | Overview of the training and validation cohorts of the GDSC and**
435 **CTRP**

436 Bar graphs showing the distribution of tissue types (**left**) and mutation frequency of the ten
437 most frequently mutated genes (**right**) for the 344 cell lines available in both the GDSC and
438 CTRP (GDSC/CTRP344) (**top**) and for the 370 cell lines only available in GDSC (GDSC370)
439 (**bottom**).

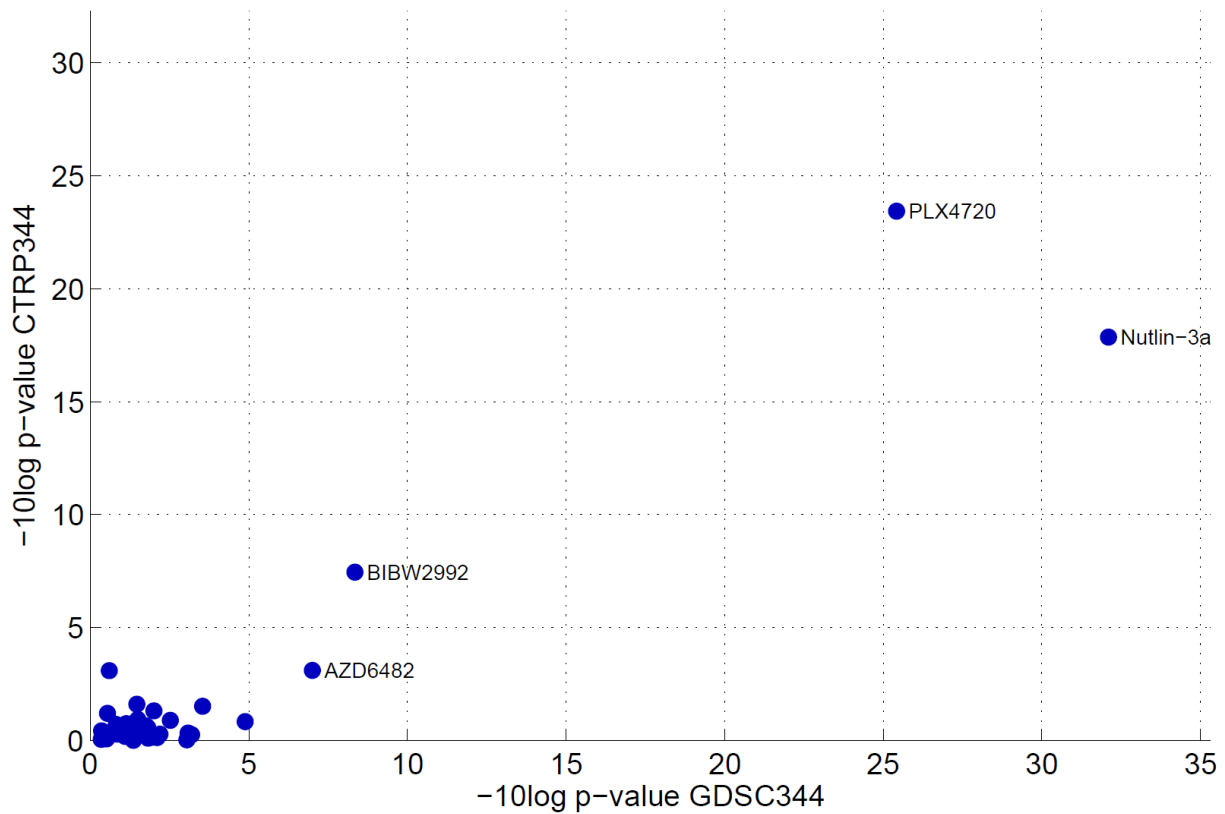
440



441
442

443 **Supplementary Figure 3 | Ordered t-test p-values for GDSC and CTRP for scenario 1 -**
444 **Train:GDSC344-Validate:CTRP344**

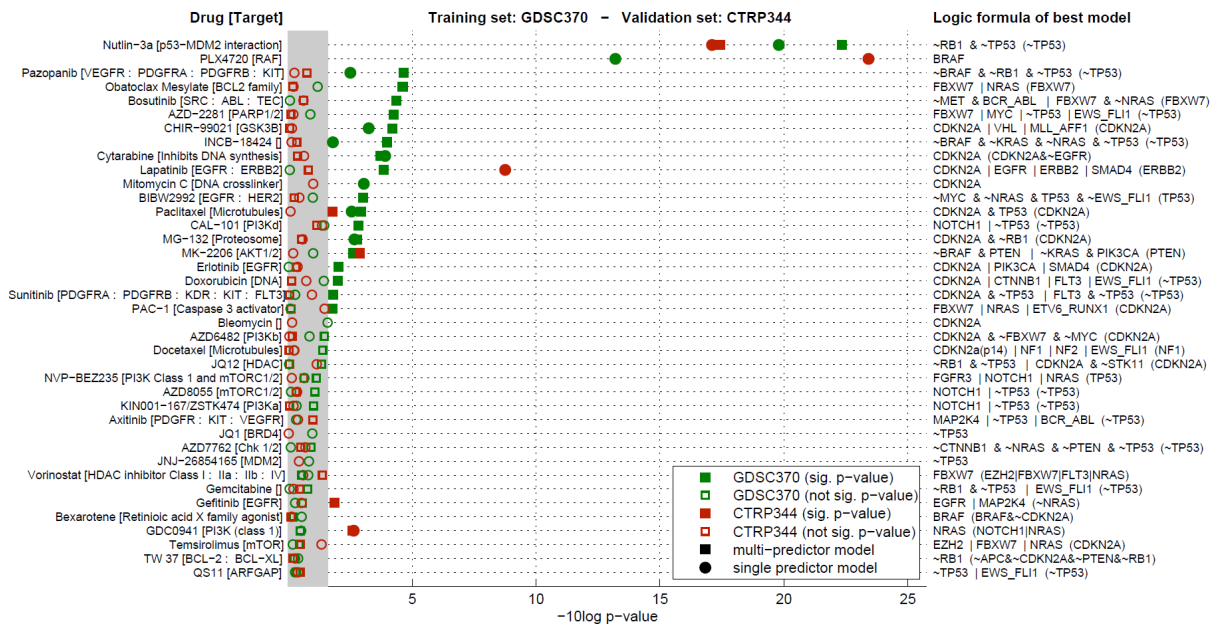
445 LOBICO models were trained on GDSC data of the 344 cell lines (GDSC344) and validated on
446 CTRP data of the same set of 344 cell lines (CTRP344). The scatter plot depicts the $-\log_{10}$ p-
447 values for t-tests that quantify the difference between cell lines predicted to be sensitive and
448 resistant according to LOBICO. In case the best model is a multi-predictor model, the $-\log_{10}$
449 p-value for the best single predictor model is also depicted, and the single-predictor formula
450 is stated in parentheses behind the formula of the multi-predictor model. In case the single-
451 predictor model led to a lower p-value, yet according to CV the multi-predictor was better,
452 the formula of the multi-predictor model is stated between parentheses. The 37 drugs are
453 sorted based on the t-test p-value derived from the GDSC344 IC50s. P-values are considered
454 significant at $p < 0.027$ (1/37).
455



456
 457
 458
 459
 460
 461
 462
 463
 464
 465
 466
 467
 468

Supplementary Figure 4 | Comparison of t-test p-values for GDSC and CTRP for scenario 1 - Train:GDSC344-Validate:CTRP344

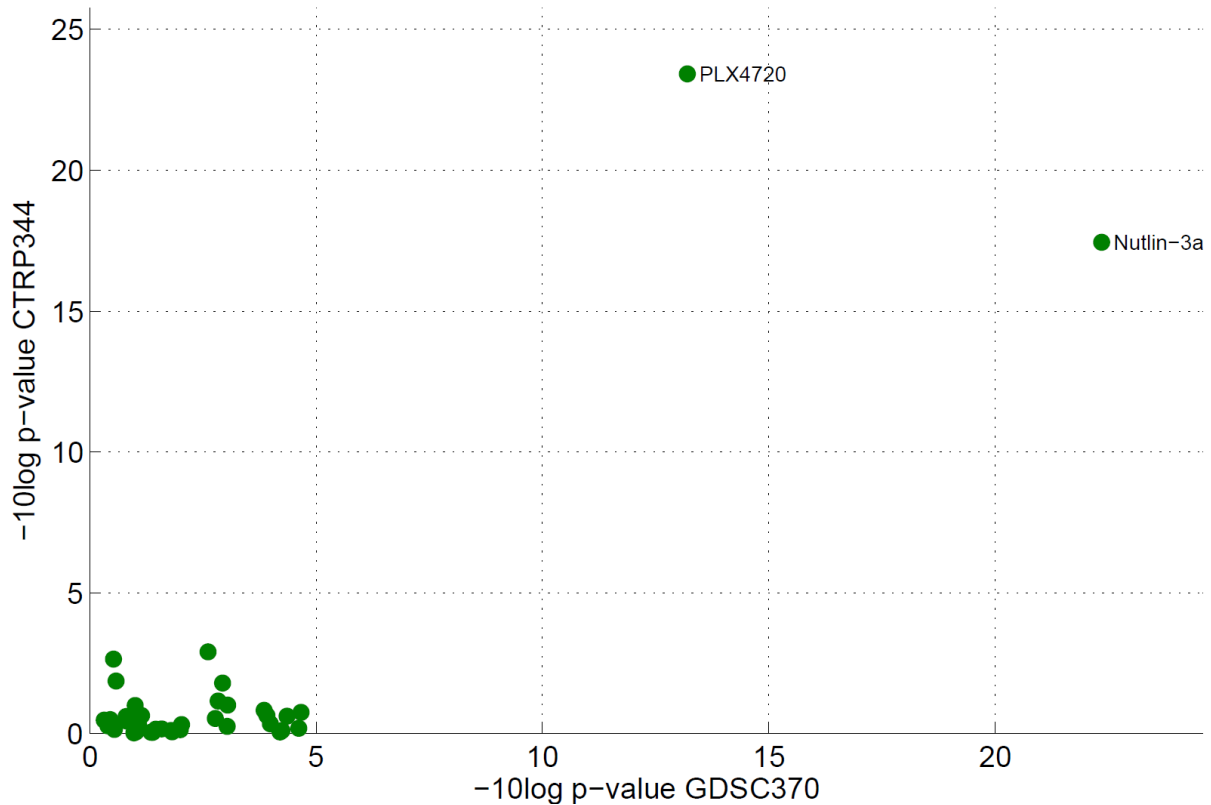
LOBICO models were trained on GDSC data of the 370 cell lines (GDSC344) and validated on CTRP data of the same set of 344 cell lines (CTRP344). The scatter plot depicts the $-\log_{10}$ p-values for t-tests that quantify the difference between cell lines predicted to be sensitive and resistant according to LOBICO. The x-axis depicts p-values for the difference between these two groups based on the IC50s within GDSC344. The y-axis depicts p-values for the difference between these two groups based on the AUCs within CTRP344. Drugs with a p-value lower than 10^{-5} are annotated.



469
470
471
472
473
474
475
476
477
478
479
480
481
482
483

Supplementary Figure 5 | Ordered t-test p-values for GDSC and CTRP for scenario 2 - Train:GDSC370-Validate:CTRP344

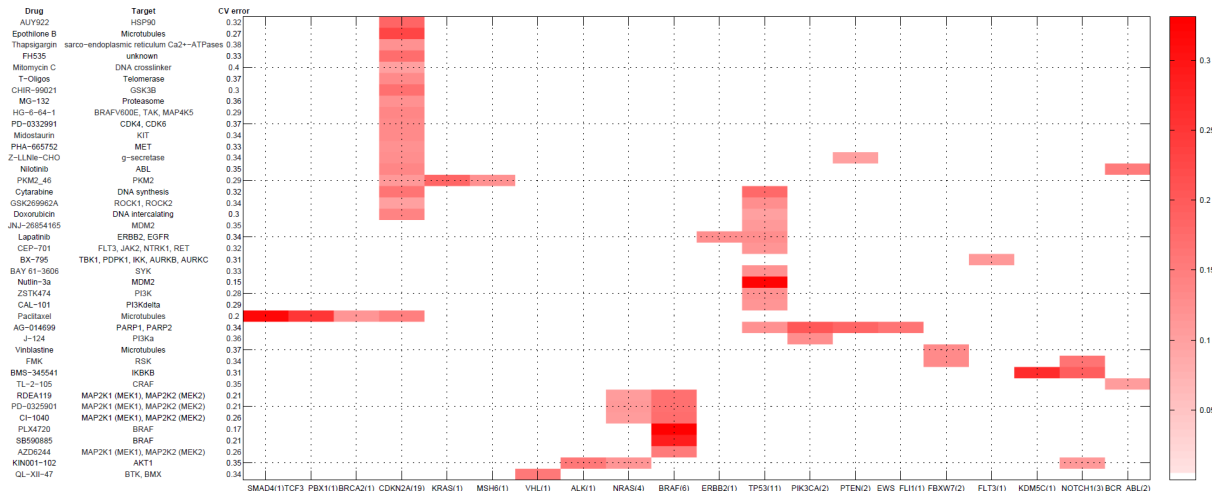
LOBICO models were trained on GDSC data of the 370 cell lines (GDSC370) and validated on CTRP data of the independent set of 344 cell lines (CTRP344). The scatter plot depicts the -log₁₀ p-values for t-tests that quantify the difference between cell lines predicted to be sensitive and resistant according to LOBICO. In case the best model is a multi-predictor model, the -log₁₀ p-value for the best single predictor model is also depicted, and the single-predictor formula is stated in parentheses behind the formula of the multi-predictor model. In case the single-predictor model led to a lower p-value, yet according to CV the multi-predictor was better, the formula of the multi-predictor model is stated in parentheses. The 39 drugs are sorted based on the t-test p-value derived from the GDSC370 IC50s. P-values are considered significant at p<0.026 (1/39).



484
485

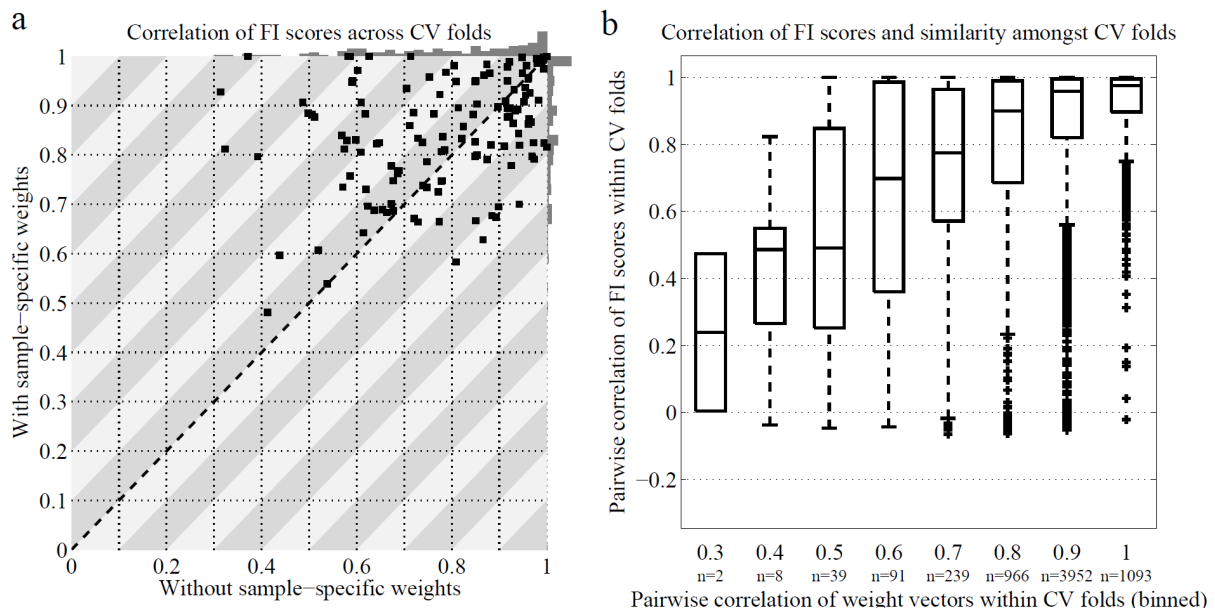
486 **Supplementary Figure 6 | Comparison of t-test p-values for GDSC and CTRP for scenario 2 -**
487 **Train:GDSC370-Validate:CTRP344**

488 LOBICO models were trained on GDSC data of the 370 cell lines (GDSC370) and validated on
489 CTRP data of the independent set of 344 cell lines (CTRP344). The x-axis depicts p-values for
490 the difference between these two groups based on the IC50s within GDSC370. The y-axis
491 depicts p-values for the difference between these two groups based on the AUCs within
492 CTRP344. Drugs with a p-value lower than 10^{-5} are annotated.
493



494
 495 **Supplementary Figure 7 | Essential gene mutation features for explaining drug response.**
 496 Of the 72 statistically significant logic models at FDR<1% and p<0.01, we selected those for
 497 which at least one gene mutation feature had a FI score of 0.1 or higher. The statistical
 498 cutoff of FDR<1% and p<0.01 that we used to identify statistically significant logic models
 499 coincided with a CV error of approximately 0.4 (see **Figure 2**). Since the FI score for a feature
 500 is the increase in error when the feature is left out of the inferred logic model, and the
 501 randomly expected CV error is 0.5, we classified gene mutation features with a FI score
 502 larger than 0.1 as ‘essential’ features for explaining the drug response. This heatmap
 503 visualizes the FI score of those features with the drugs on the rows and the features on the
 504 columns. The number in parentheses behind the gene labels indicates the number of drugs
 505 for which these genes had an FI score larger than 0.1. These results show that most gene
 506 mutations are essential for not more than one drug, but that BRAF (6), CDKN2A (19) and
 507 TP53 (11) play an important role for many of the drugs in our panel.

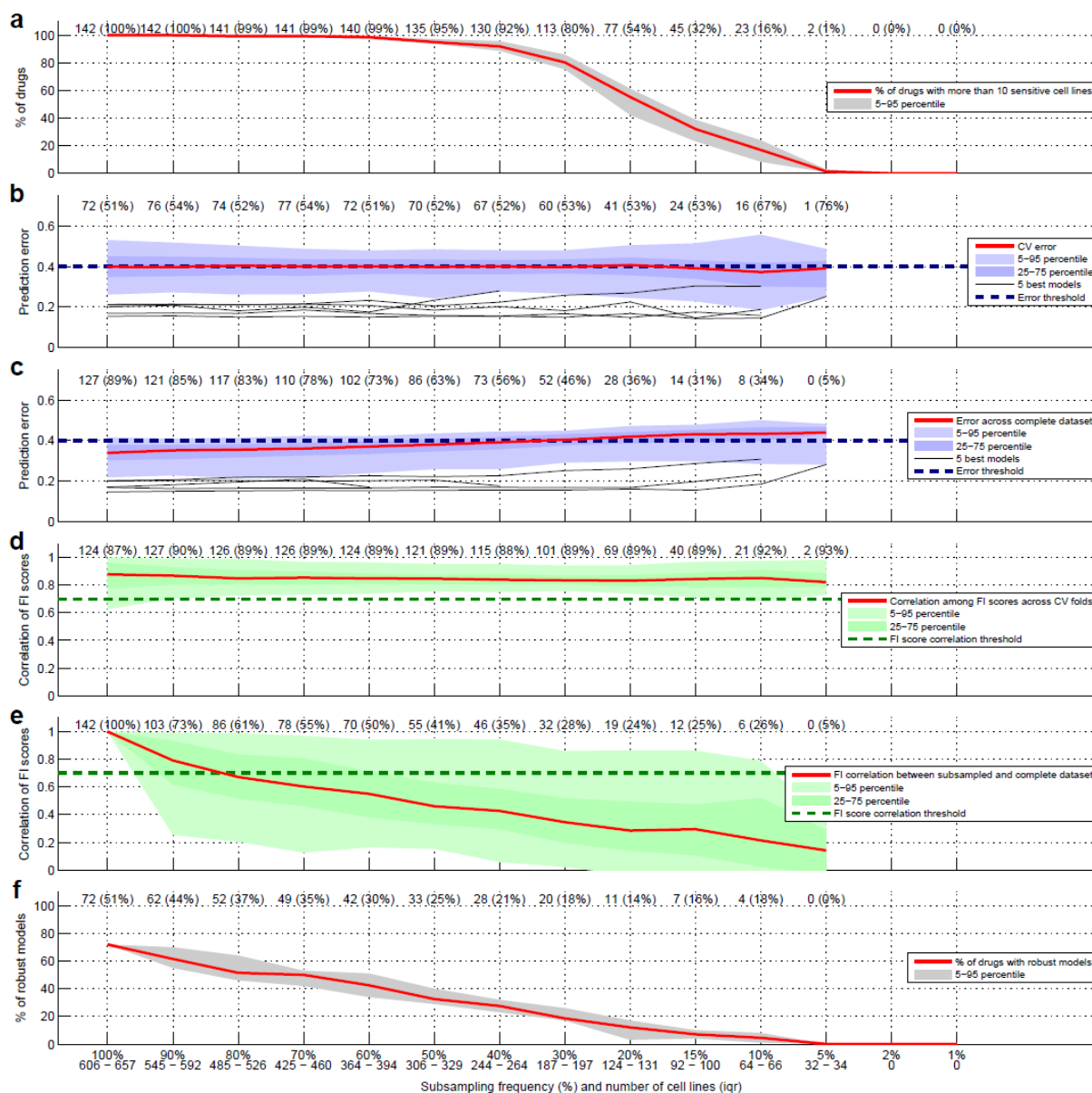
508
 509
 510
 511
 512
 513
 514



515
 516
 517
 518
 519
 520
 521
 522
 523
 524
 525
 526
 527
 528

Supplementary Figure 8 | Robustness across CV folds

a) Scatter plot with the average Pearson correlation coefficients of the similarity of FI scores across the 10 CV folds for inferred logic models without (x-axis) and with (y-axis) the sample-specific weights. Each point represents one of the 142 drugs. The correlation scores are computed using the model-complexity-specific FI scores. The grey bars on top and to the right of the scatter plot represent histograms of these correlation scores for models without and with the sample-specific weights, respectively. **b)** Boxplot comparing the pairwise correlation of weight vectors between CV folds (x-axis) with the pairwise correlations of FI scores between the same CV folds (y-axis). The pairwise correlation of weight vectors were binned by rounding the correlation to the nearest decimal. The number of correlations per box is indicated below the box.



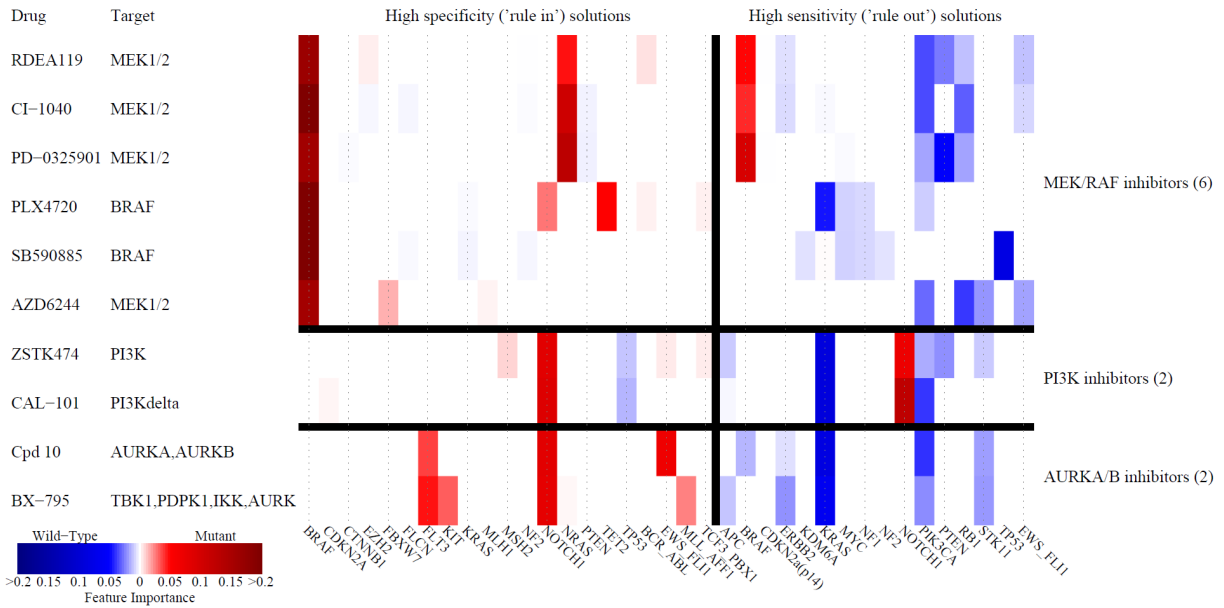
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545

Supplementary Figure 9 | Subsampling analysis to identify robust logic models with fewer cell lines

LOBICO models were run on randomly selected sets of cell lines for an array of subsampling frequencies (x-axis). The subsampling frequency and interquartile range of the number cell lines associated with the subsampling frequency is found in the bottom of the figure. From top to bottom: **a**) The percentage of drugs for which LOBICO models could be run, i.e. for which the number of sensitive cell lines was 10 or greater. Percentile values show variation across the repeats. Absolute numbers and percentages are stated in the top of the plot. **b**) The CV-error. Percentile values show variation across the repeats and across drugs. Average number and percentage of drugs with a CV-error lower than the threshold of 0.4 are stated in the top of the plot. Note that the percentages are based on the number of drugs that can be run for a particular subsampling frequency (a), not based on all 142 drugs. This is also the case for c), d), e) and f). **c**) The error across the complete dataset based on the logic models inferred from the subsampled datasets. Percentile values show variation across the repeats and across drugs. Average number and percentage of drugs with an error lower than the threshold of 0.4 are stated in the top of the plot. **d**) The Pearson correlation among the FI

546 scores across the CV folds. Percentile values show variation across the repeats and across
547 drugs. Average number and percentage of drugs with a FI score correlation larger than the
548 threshold of 0.7 are stated in the top of the plot. **e)** The Pearson correlation between the FI
549 scores of the logic models from the subsampled datasets and the complete dataset.
550 Percentile values show variation across the repeats and across drugs. Average number and
551 percentage of drugs with a FI score correlation higher than the threshold of 0.7 are stated in
552 the top of the plot. **f)** The percentage of robust and predictive models, i.e. those that meet
553 all criteria. Percentile values show variation across the repeats. Absolute numbers and
554 percentages are stated in the top of the plot.
555
556

557

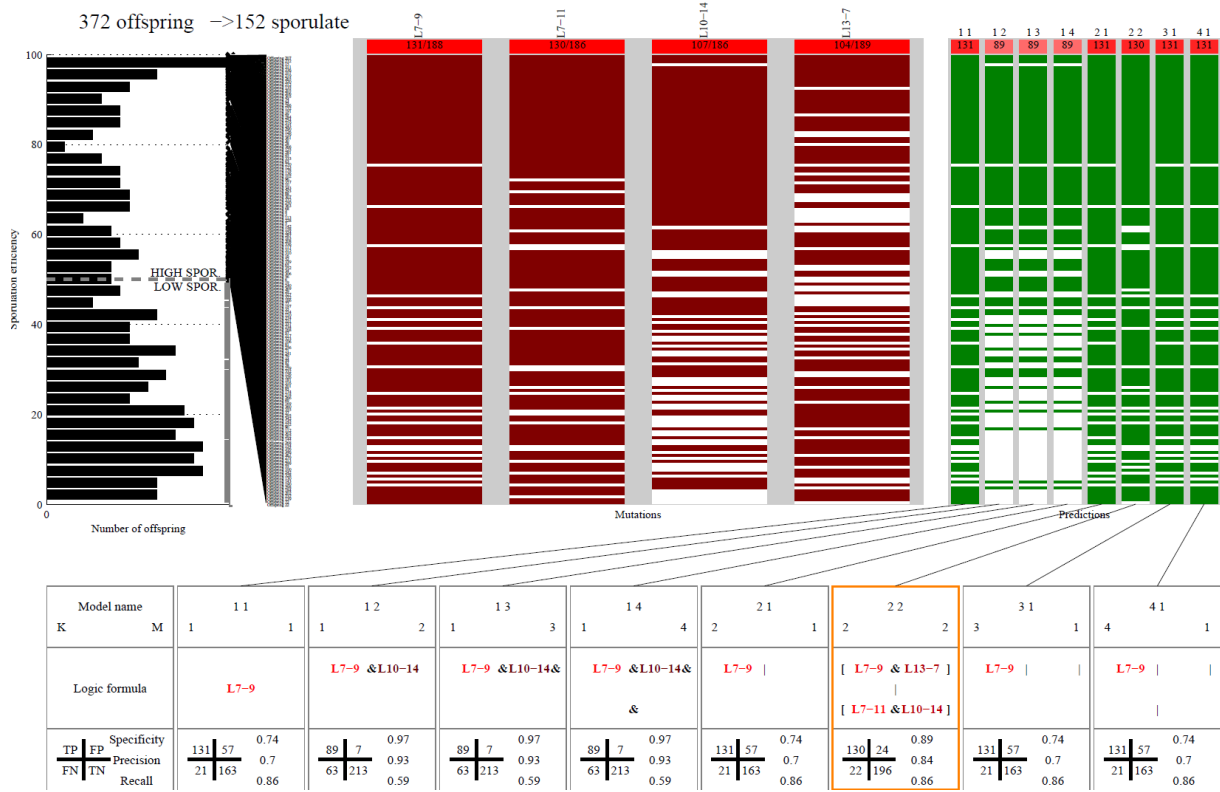


558

559

560 **Supplementary Figure 10 | Feature importance scores for 'rule in' and 'rule out' solutions**
 561 FI scores of 6 MEK/RAF, 2 PI3K and 2 AURKA/B inhibitors (rows) for high specificity ('rule in')
 562 solutions (left) and high sensitivity ('rule out') solutions (right). High specificity solutions
 563 were defined as solutions with FPR<10%. Conversely, high sensitivity solutions were defined
 564 as solutions with TPR>90%. We distinguished between positive terms, indicating mutations
 565 (red) and negated terms, indicating wild-type (blue).
 566

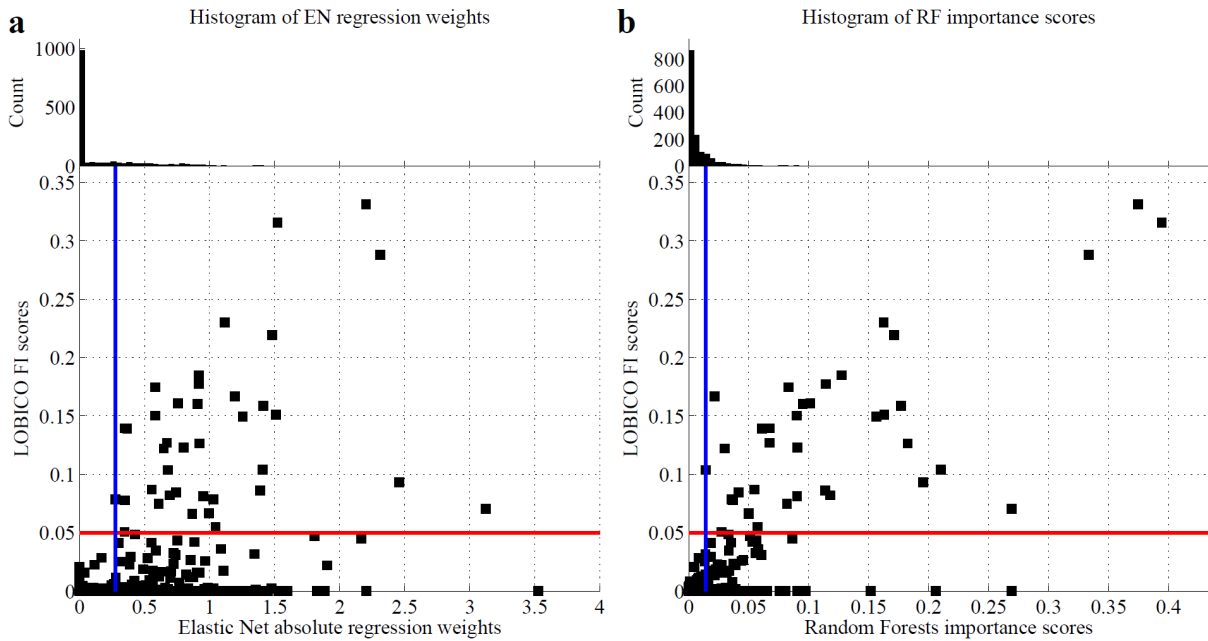
567



568
569
570
571
572
573
574
575

Supplementary Figure 11 | LOBICO results of the yeast cross phenotyped for sporulation efficiency
Standard LOBICO visualization of the uncovered logic models for the yeast cross dataset (Supplementary Note 4). See Supplementary Data 1 for an explanation of the visualization.

576



577

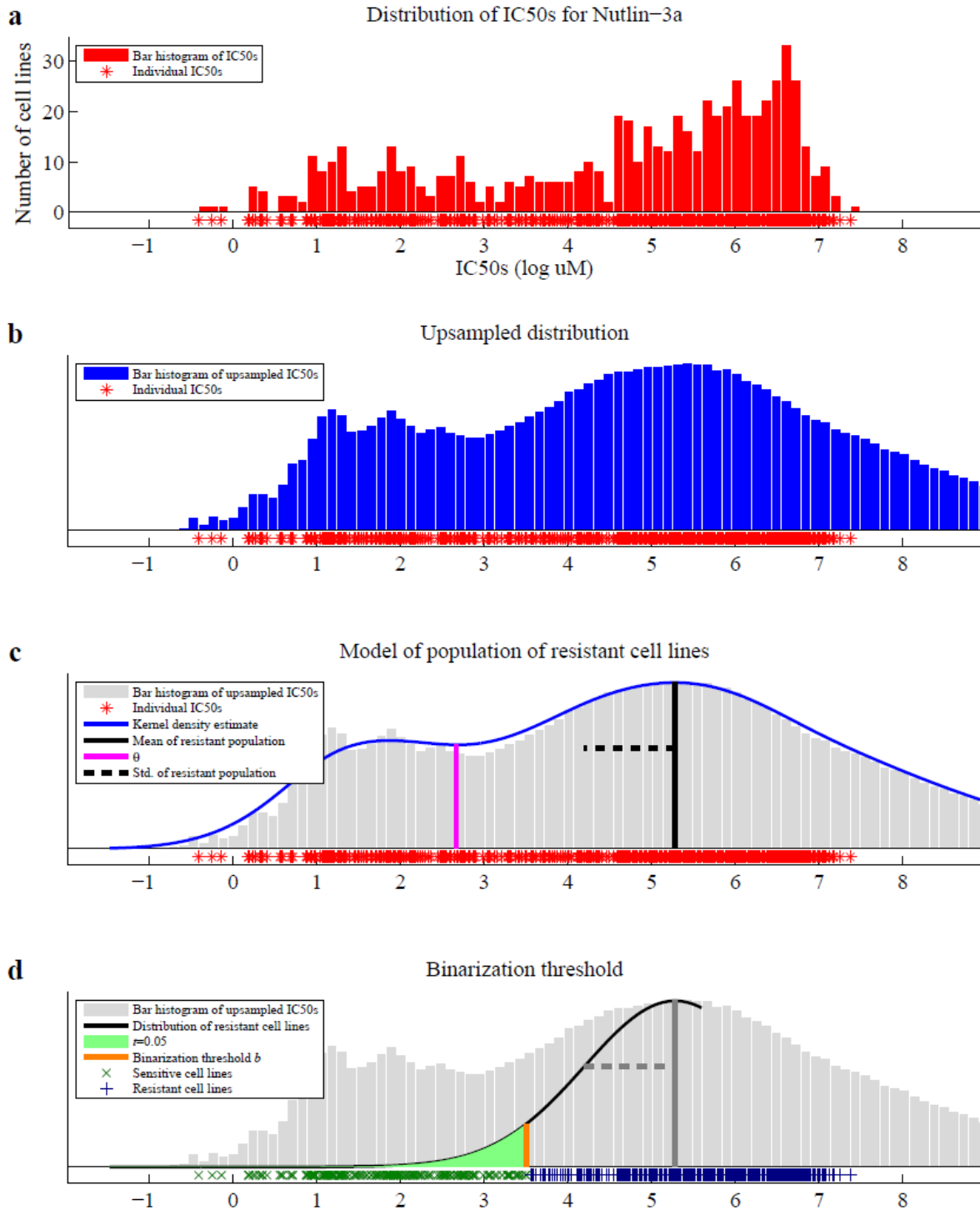
578

579 **Supplementary Figure 12 | Comparison of feature importance scores between LOBICO,**
580 **Elastic Net and Random Forests**

581 **a)** Scatter plot comparing LOBICO's FI scores with EN's absolute regression weights. These
582 scores and weights are derived from inferred models of the 25 drugs with the lowest CV
583 error in the LOBICO analysis (**Supplementary Note 6**). The red line depicts the FI score of
584 0.05; features with a FI>0.05 are considered important predictors. The blue depicts the
585 minimal EN regression weight for which the corresponding LOBICO FI was larger than 0.05.

586 **b)** Similar to **a)**, except LOBICO's FI scores are compared to RF's importance scores.

587



589

590

591

Supplementary Figure 13 | Four-step-procedure to binarize IC50s for Nutlin-3a

592

a) Histogram plot for the distribution of IC50s for the drug Nutlin-3a. **b)** Histogram plot for

593

the upsampled distribution **c)** Visualization of an empirical model (obtained through density

594

estimation) of the upsampled IC50s (depicted in blue). θ was computed using rule i. (See

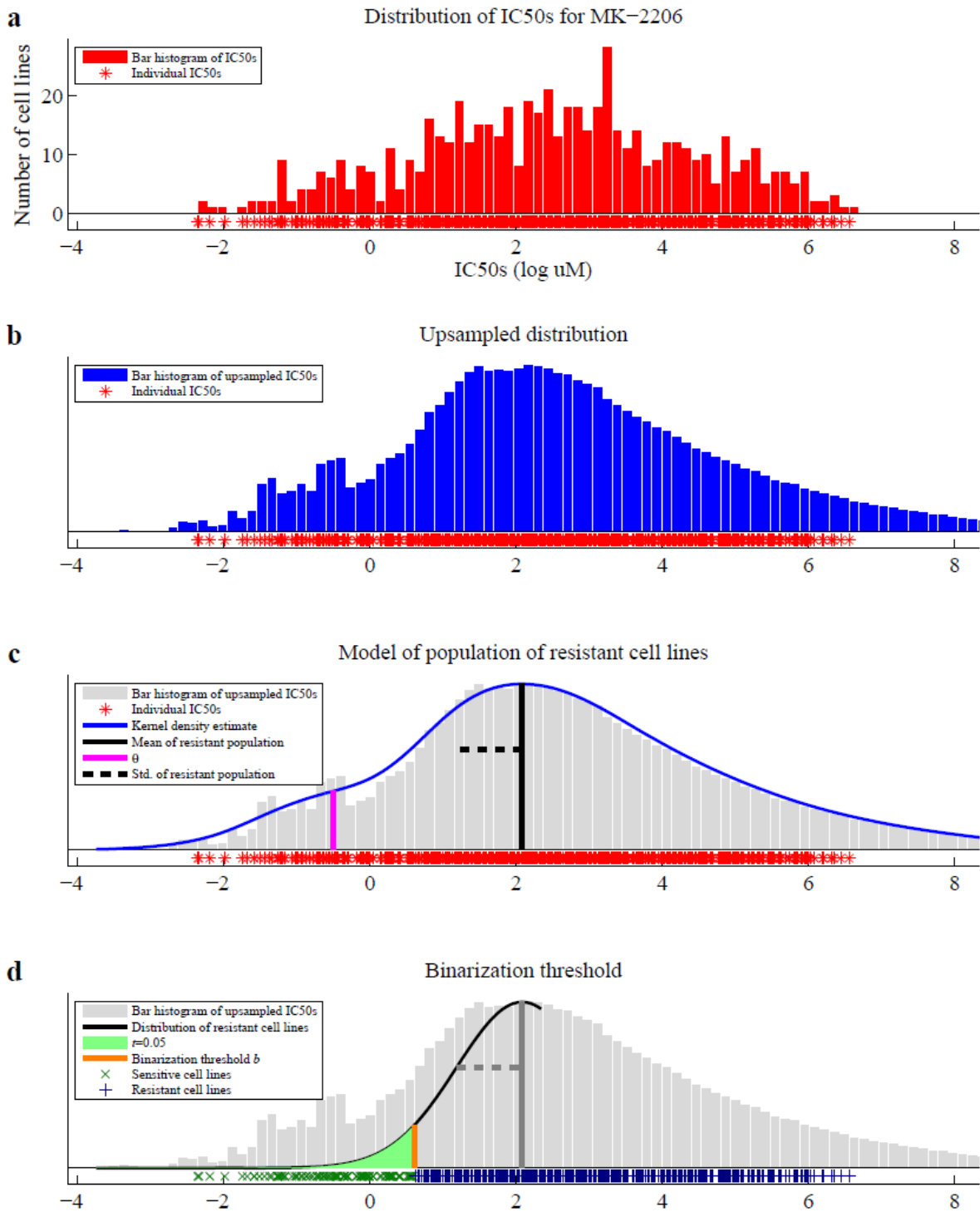
595

Methods Section for details.) **d)** Visualization of the model of resistant cell lines (depicted in

596

black), from which the binarization threshold b (depicted in orange) is derived.

597



599

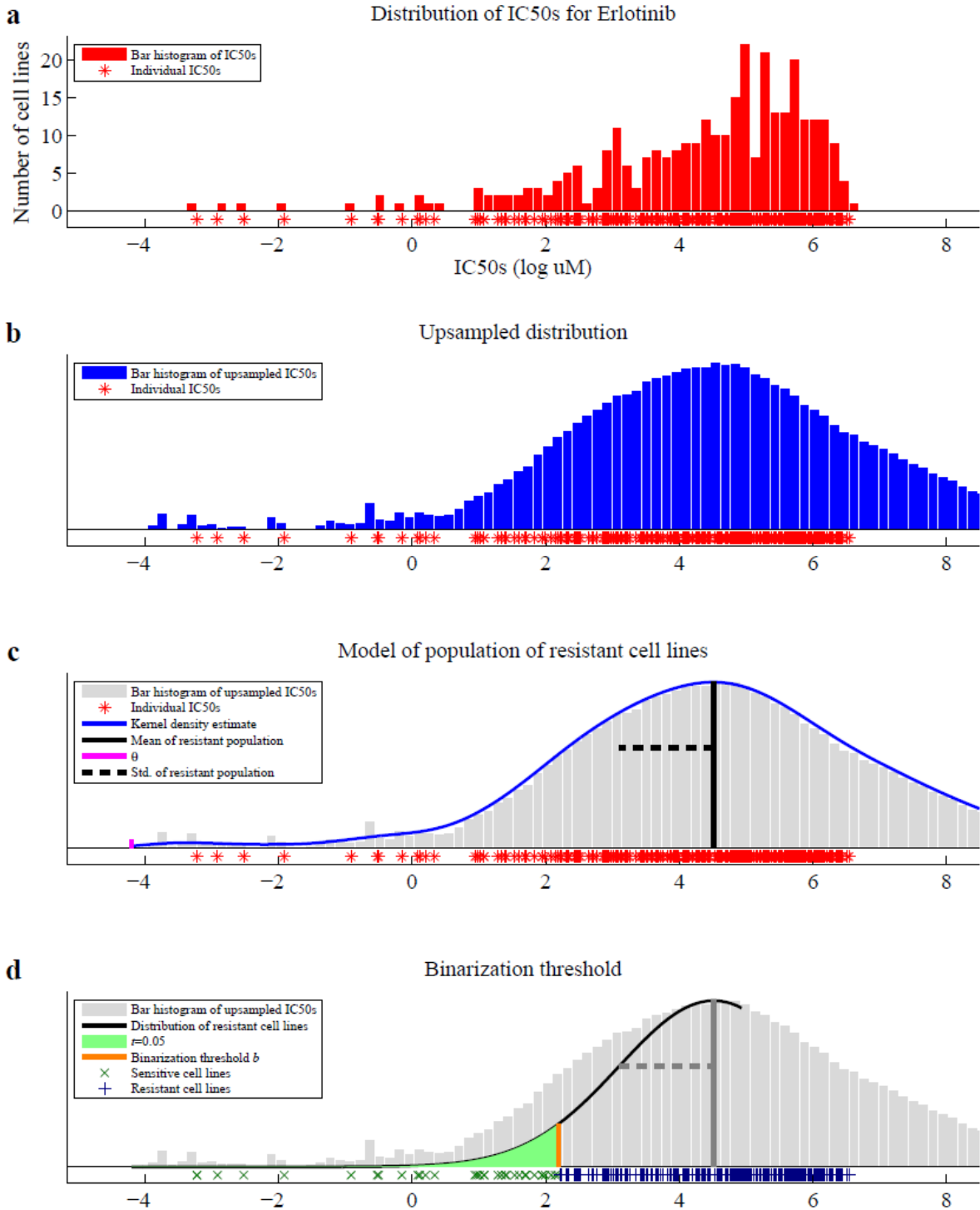
600

601 **Supplementary Figure 14 | Four-step-procedure to binarize IC50s for MK-2206**

602 Similar to **Supplementary Figure 13**, except showing the procedure for drug MK-2206, and

603 the use of rule ii to find θ .

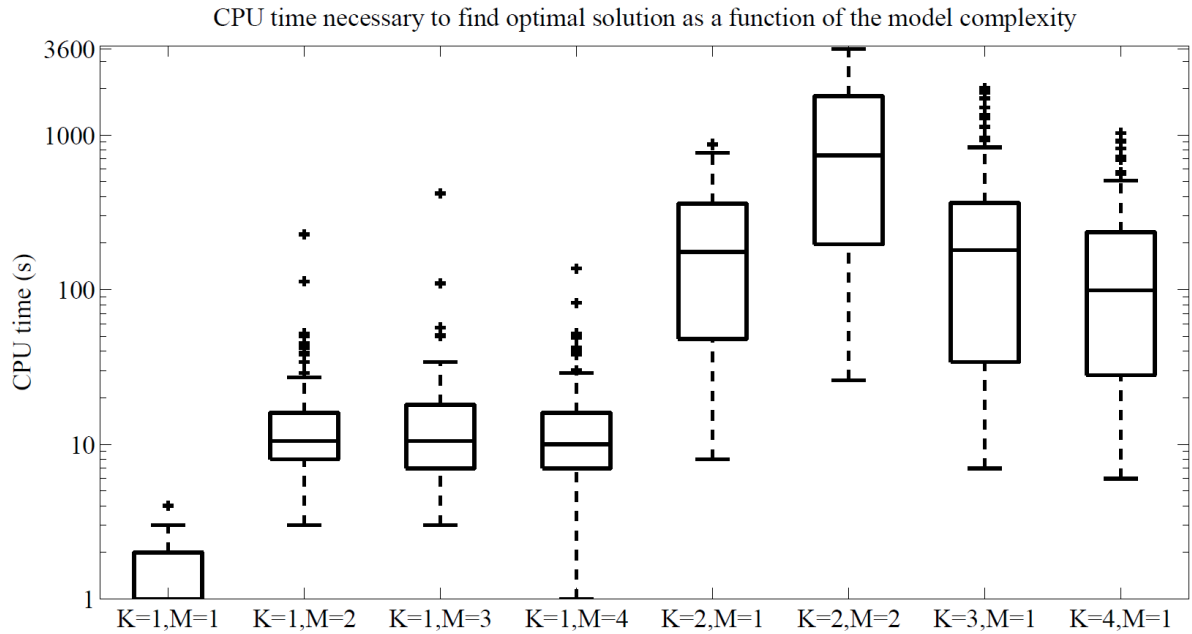
604



606
 607
 608
 609
 610
 611
 612

Supplementary Figure 15 | Four-step-procedure to binarize IC50s for Erlotinib
 Similar to **Supplementary Figure 13**, except showing the procedure for drug Erlotinib, and the use of rule iii to find θ .

613



614

615

616 **Supplementary Figure 16 | Time needed to find optimal solution**

617

618 Boxplot of CPU time (y-axis) necessary to find the optimal solution as a function of

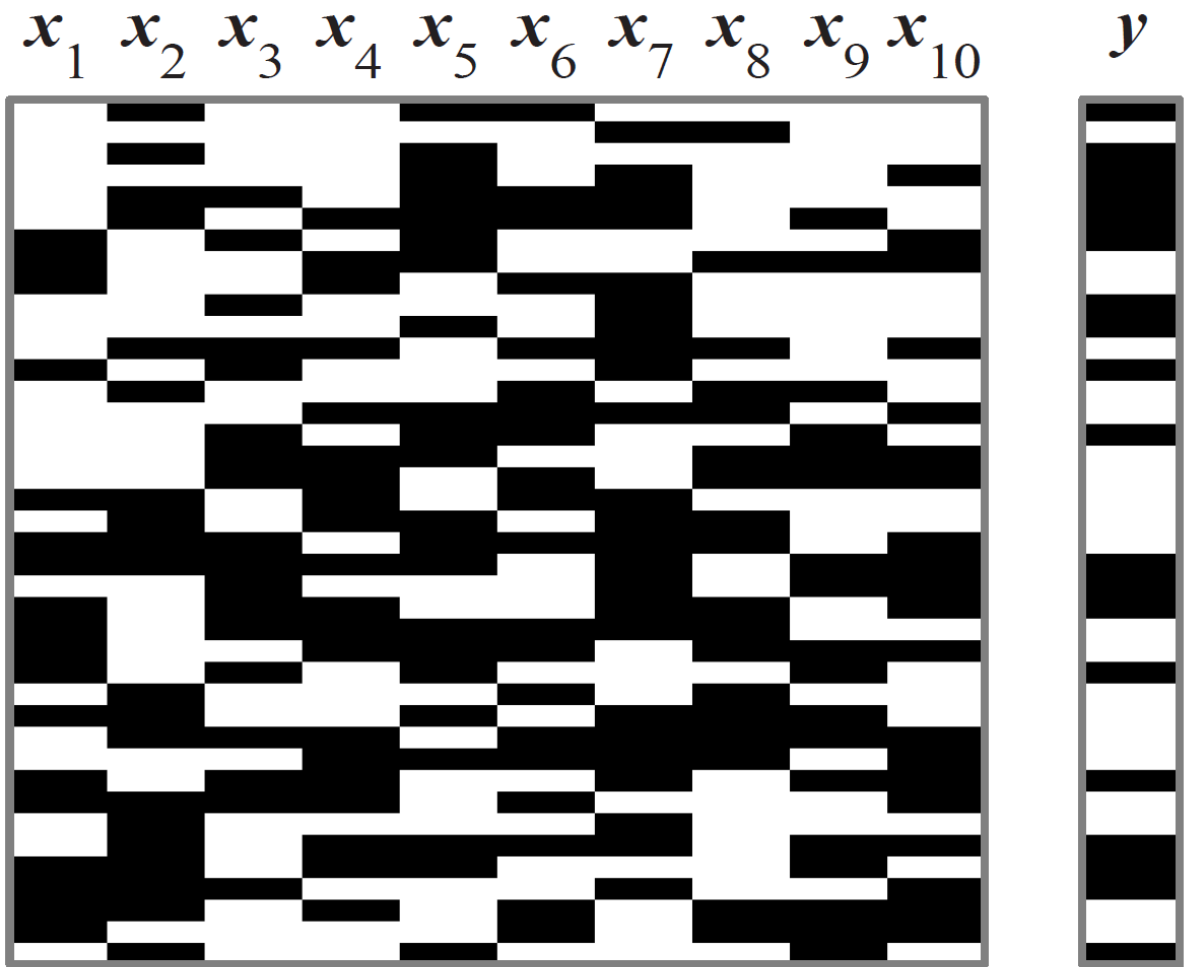
619 the model complexity (x-axis). Each box is comprised of 142 values, i.e. the time necessary

620 for CPLEX to find the optimal solution with the indicated model complexity for each of the

621 142 drugs. These experiments were performed on a computer cluster, where each ILP was

622 run on one node (Intel(R) Xeon(R) CPU, E5645, 2.40GHz, 6 cores) at a time.

623



$$y = x_5 \bar{x}_8 + x_3 \bar{x}_6 x_7$$

624

625

626 **Supplementary Figure 17 | Example Boolean truth table**

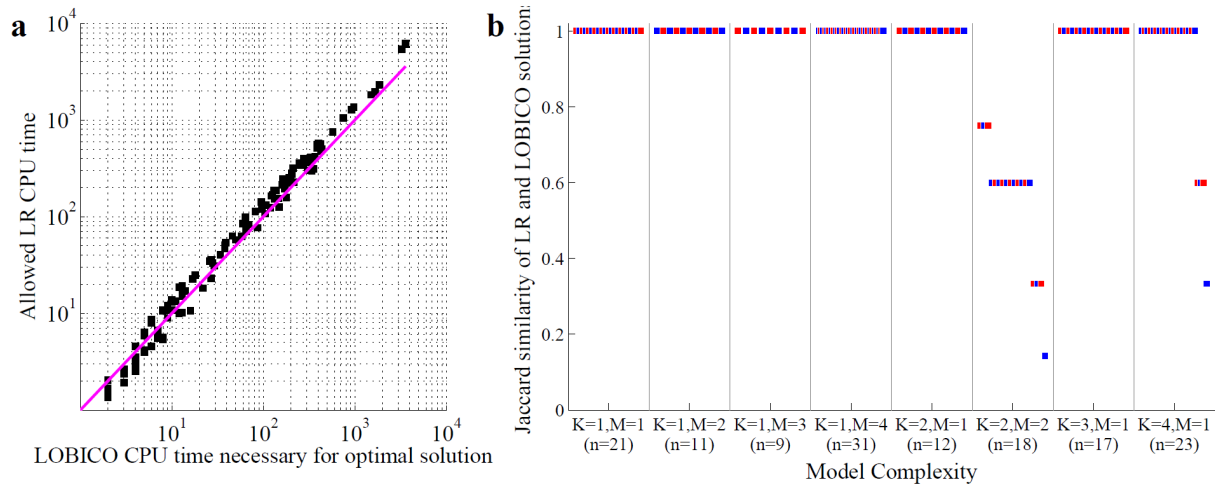
627 Boolean truth table (black=1, white=0) with 10 input variables x_1, x_2, \dots, x_{10} and

628 output variable y . The DNF expression with $K = 2$ and $M = 3$ for the truth table depicted is

629 below the table.

630

631



632

633

634 **Supplementary Figure 18| Comparison of solution time and uncovered logic models**
635 **between LOBICO and logic regression**

636 **a)** Scatter plot comparing the CPU time needed for LOBICO to find the optimal
637 solution (x-axis) and the CPU time given to LR to find the best solution (y-axis). Each of the
638 142 drugs is represented by a point. The magenta line is $y=x$. **b)** Plot of the Jaccard similarity
639 between the LR and LOBICO solutions. Each of the 142 drugs is represented by a point. The
640 points are alternately colored in blue and red for visibility. The drugs are grouped based on
641 the model complexity (x-axis) for which the LOBICO and LR models were inferred.

642

643 **Supplementary Data Explanation**

644 **Supplementary Data 1 | LOBICO visualization of the inferred logic models for all 142 drugs**
645 PDF with a visualization of the LOBICO results for each drug (pages 8-149). The first 7 pages
646 provide a visual explanation of the visualization.

647
648 **Supplementary Data 2 | Drug information**

649 Tab separated Excel table containing information about the 142 drugs. Specifically, (from left
650 to right in the table), information about the drugs (ID, name and target), binarization
651 thresholds, ground truth mapping to the gene mutation features, model performance
652 statistics of the inferred logic models, and aggregated feature importance scores for the
653 gene mutation features in the inferred logic models.

654
655 **Supplementary Data 3 | 25 ROC models visual**

656 PDF with visualizations of LOBICO solutions in the ROC space for each of 25 drugs with the
657 lowest CV error in the original analysis. Blue crosses indicate the true positive rate (TPR) and
658 false positive rate (FPR) at which the solution was found. The logic formula of the solutions is
659 printed next to the blue crosses. The color of the genes in a formula indicate their FI. Colors
660 range from black (moderately important) to bright red (highly important). For comparison,
661 the best single predictor solutions are visualized in green. The inlay depicts the histogram of
662 IC50s of the drug together with the binarization threshold.

663 The PDF consists of 50 pages; each of the 25 drugs is represented by two visualizations: one
664 using the standard (binary) definitions sensitivity (or TPR) and specificity (or 1-FPR), and one
665 using the weighted definitions of sensitivity and specificity (see Equations 12 and 13). These
666 visualizations are similar to **Figure 4a**. The visualizations are generated automatically; text
667 strings are partially overlapping in some cases.

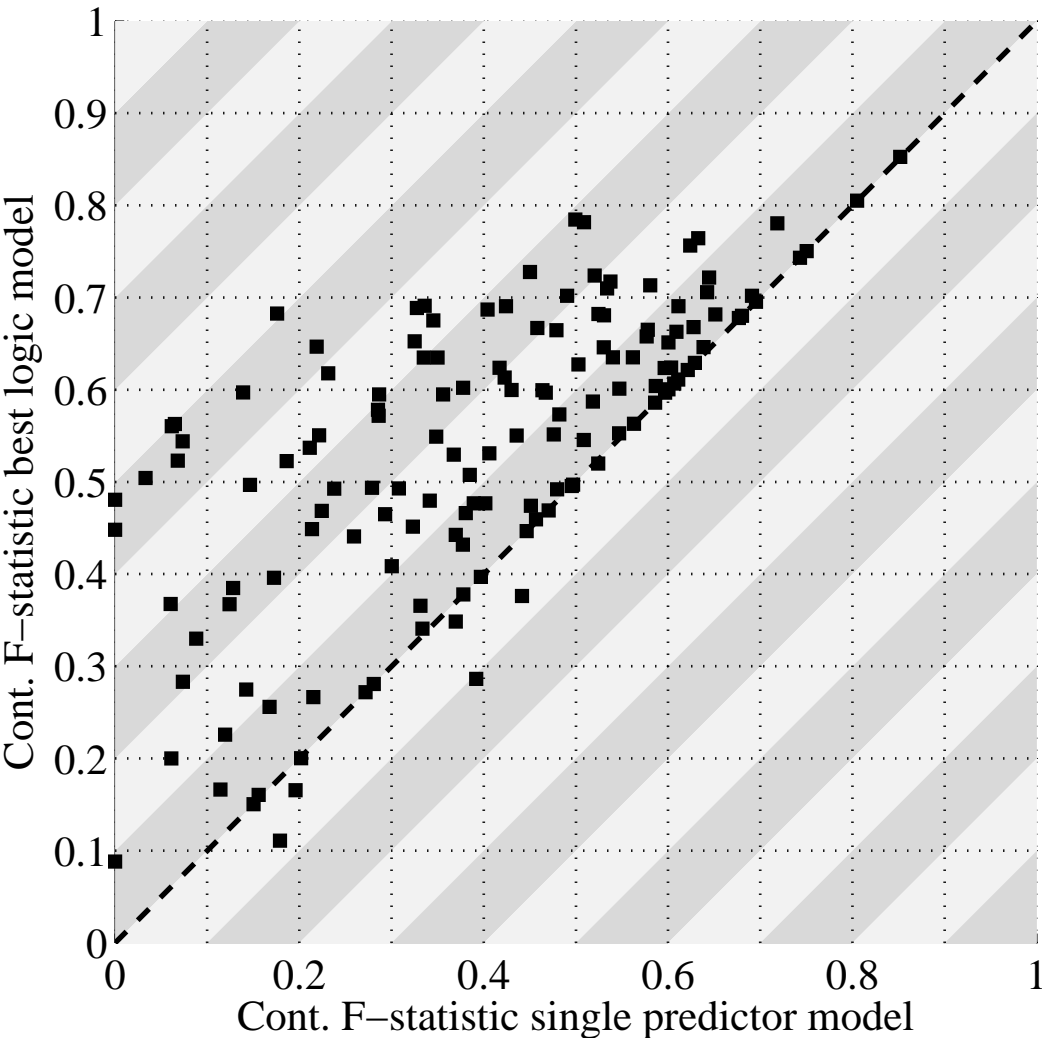
668
669 **Supplementary Data 4 | Cell line information**

670 Tab separated Excel table containing information about the 714 cell lines. Specifically, (from
671 left to right in the table), information about the cell lines (ID, name and description), the
672 binary mutation status of the 60 gene mutation features, and the IC50s for each of the 142
673 drugs.

674

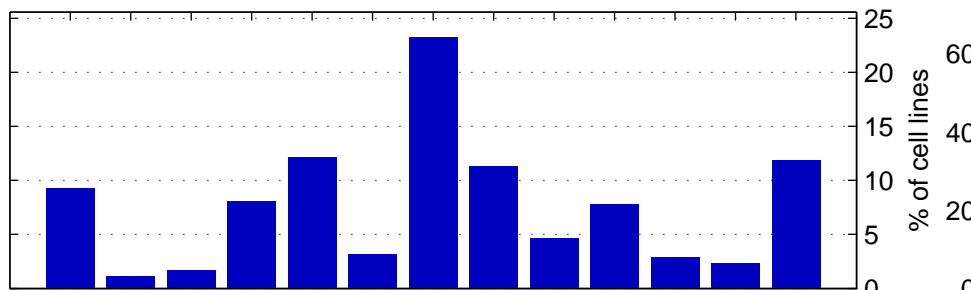
Continuous F–statistic for each of the 142 drugs

SF1



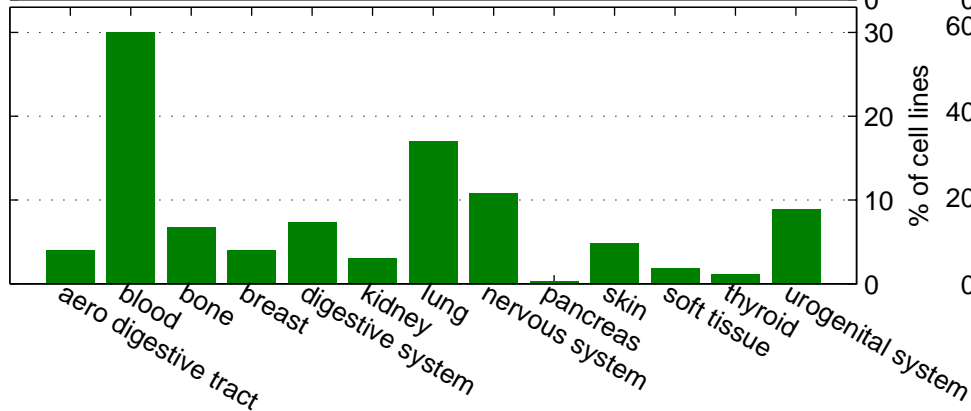
Distribution of tissue types

GDSC/CTRP344
344 cell lines in GDSC and CTRP

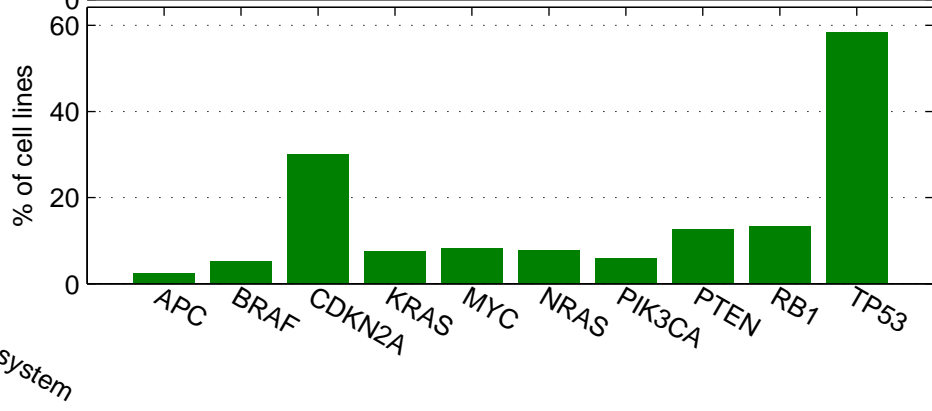
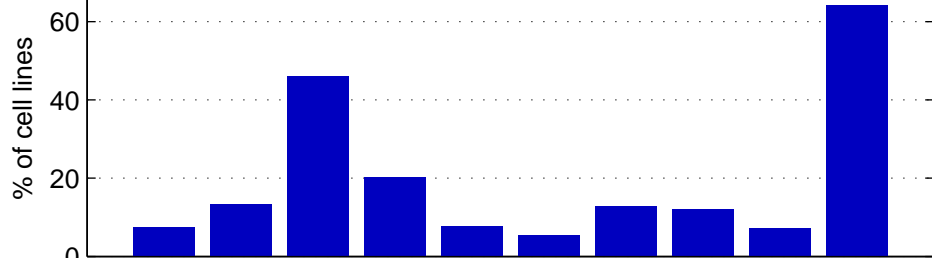


GDSC370

370 cell lines only in GDSC



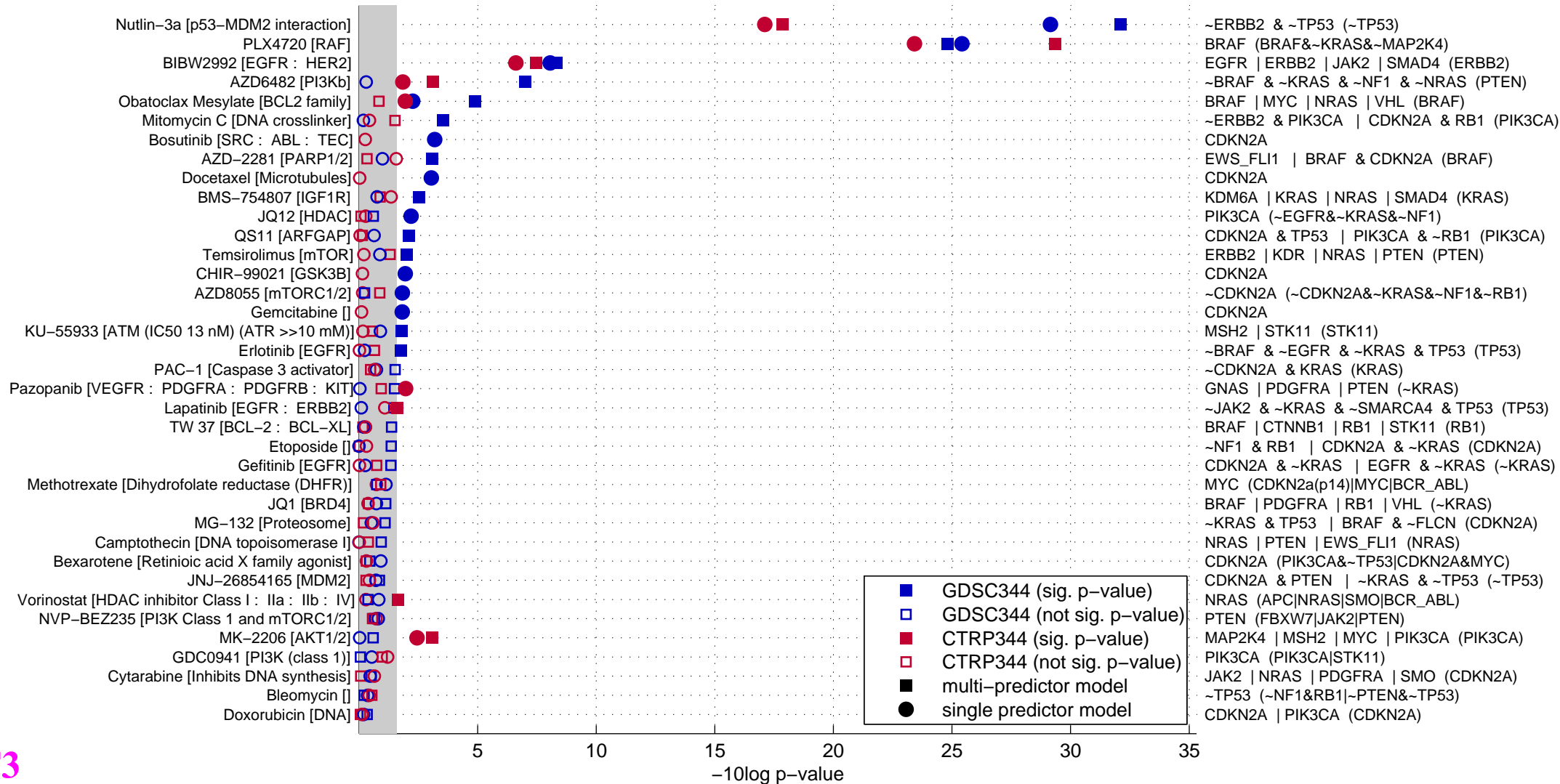
Mutation frequency

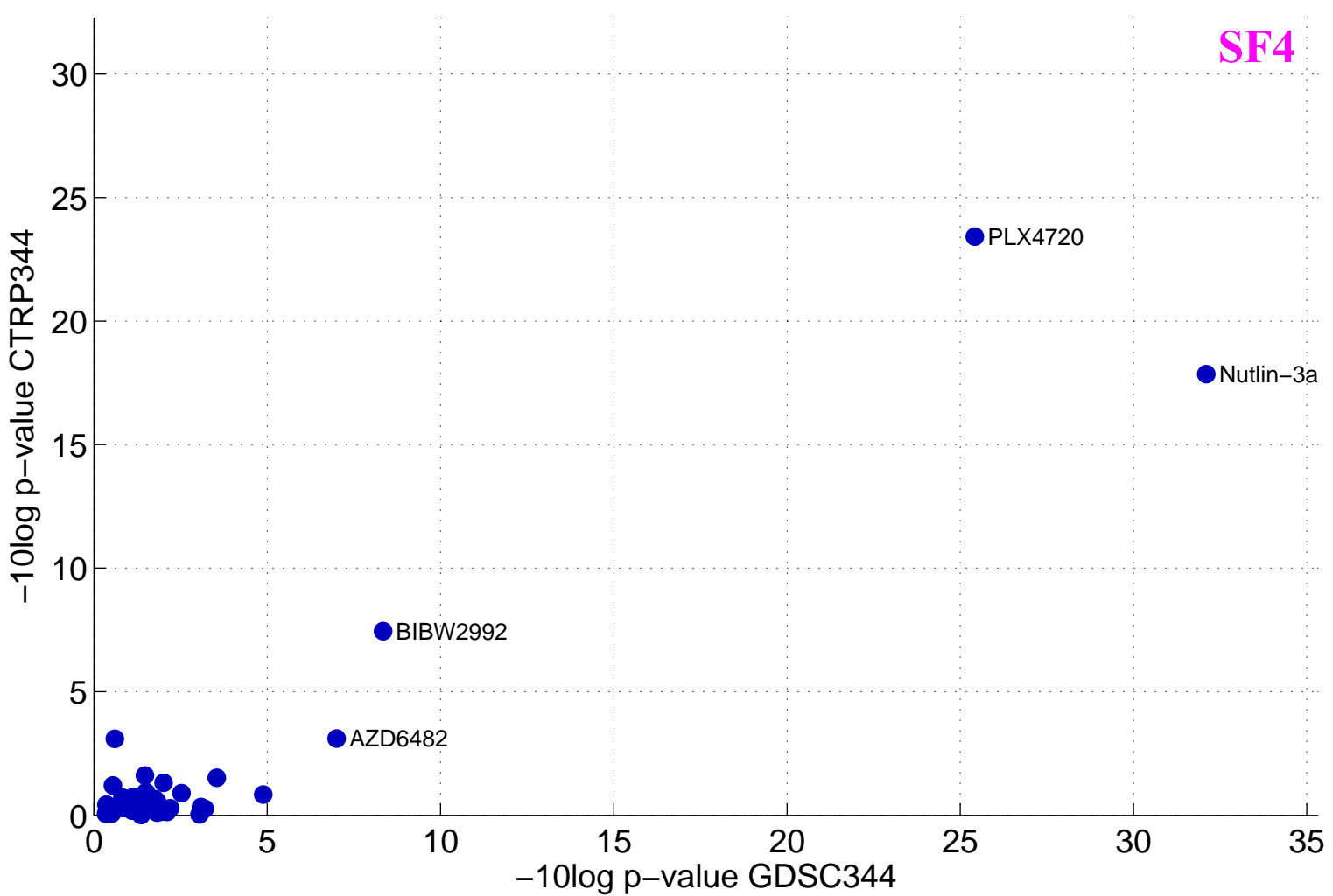


Drug [Target]

Training set: GDSC344 - Validation set: CTRP344

Logic formula of best model

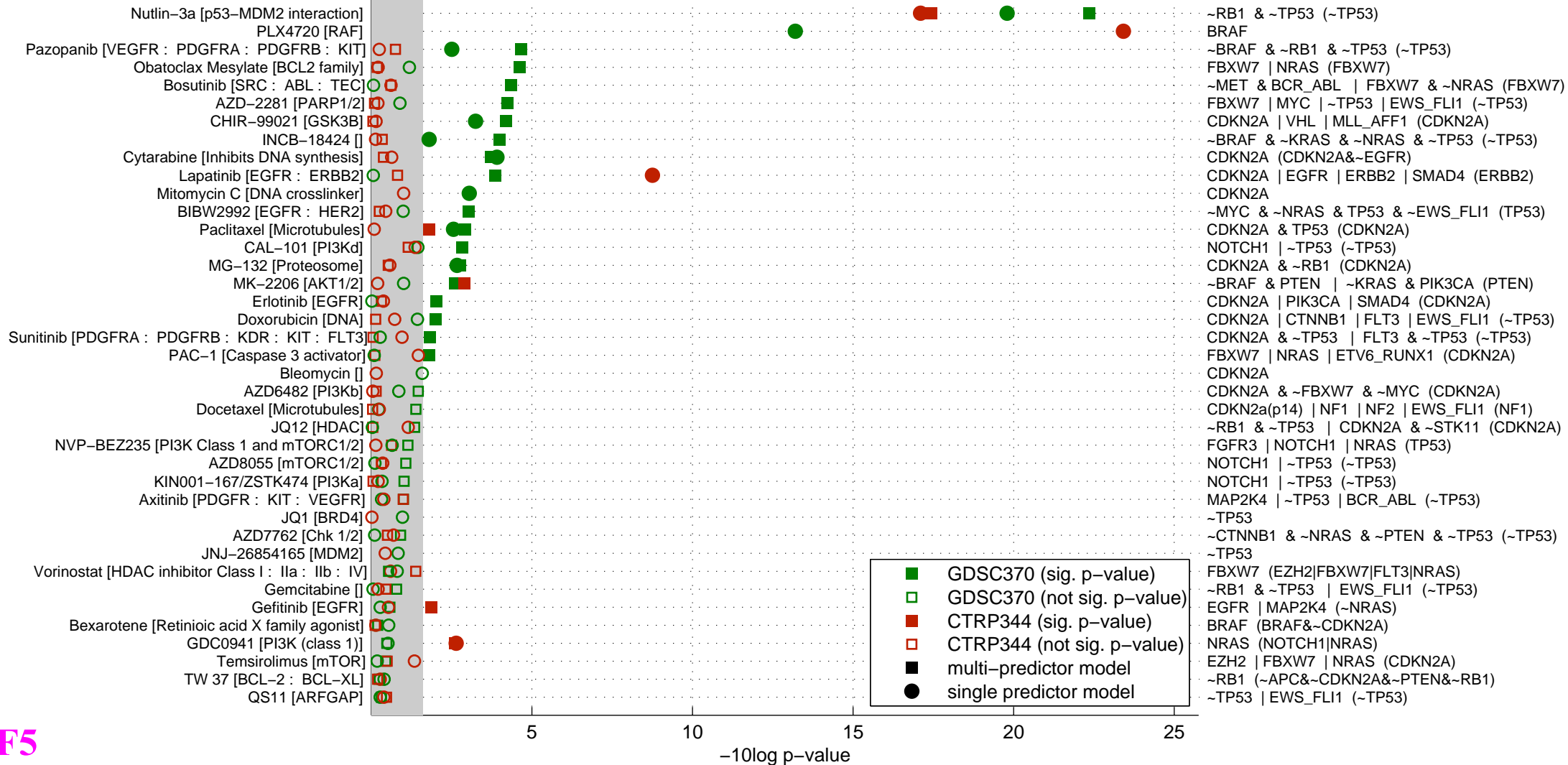




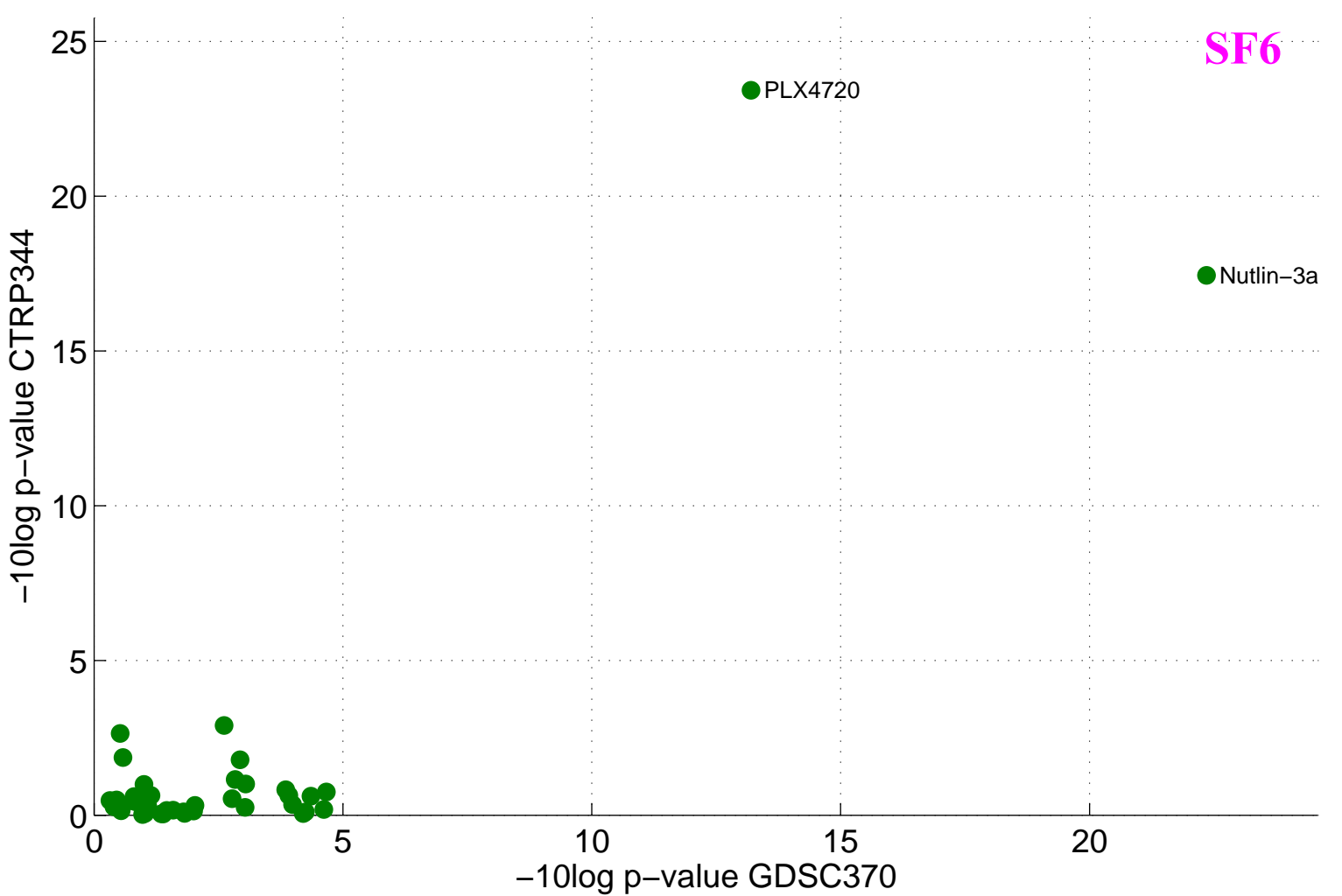
Drug [Target]

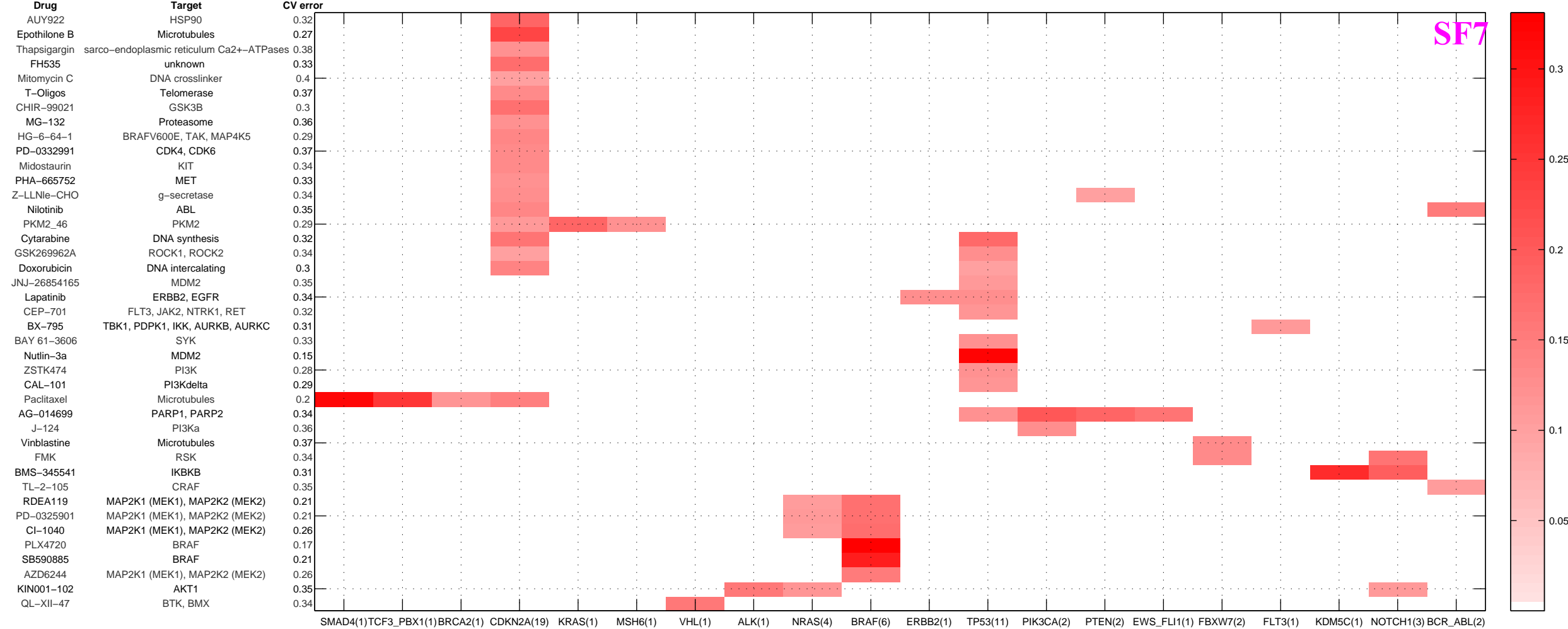
Training set: GDSC370 - Validation set: CTRP344

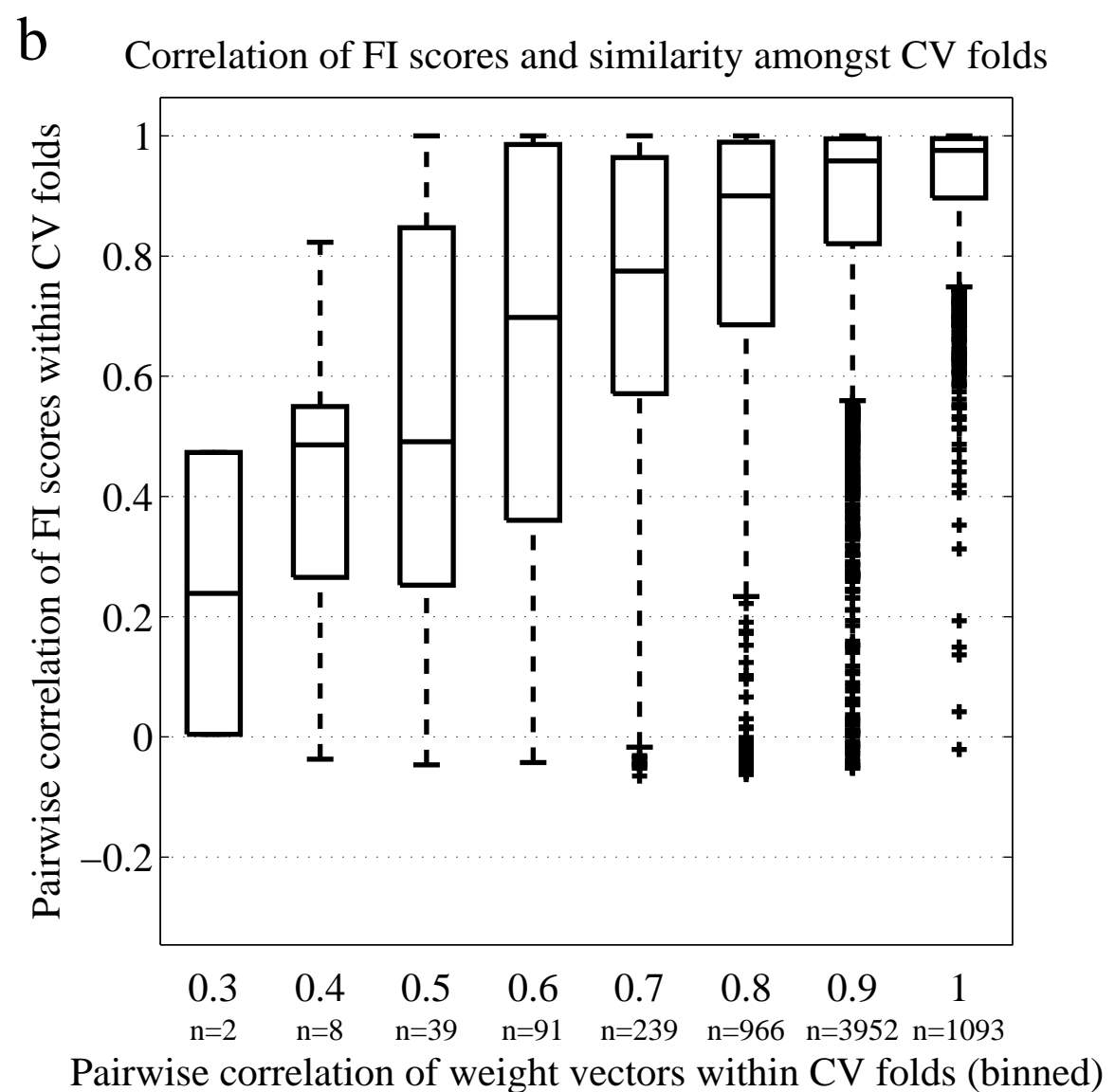
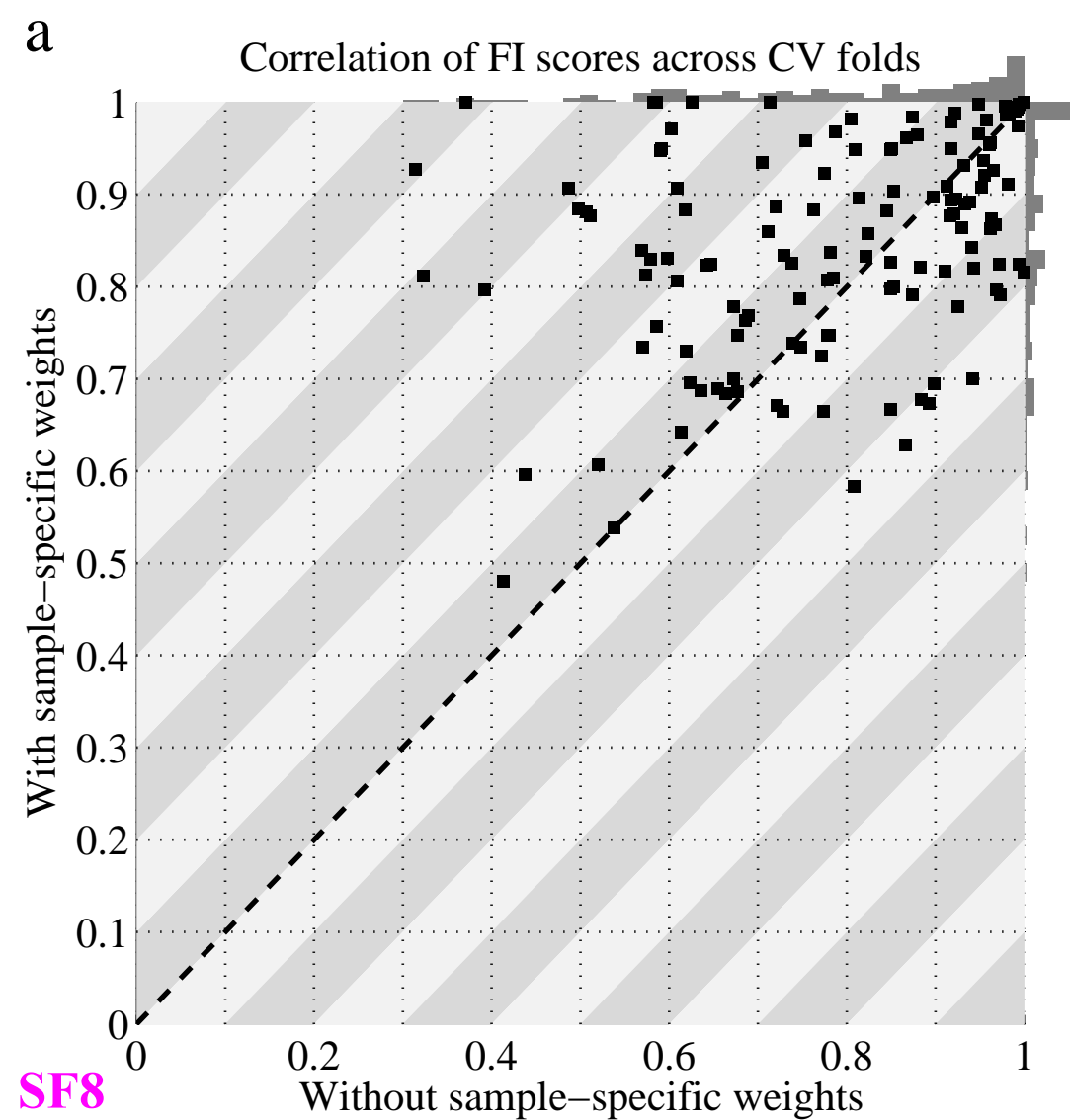
Logic formula of best model

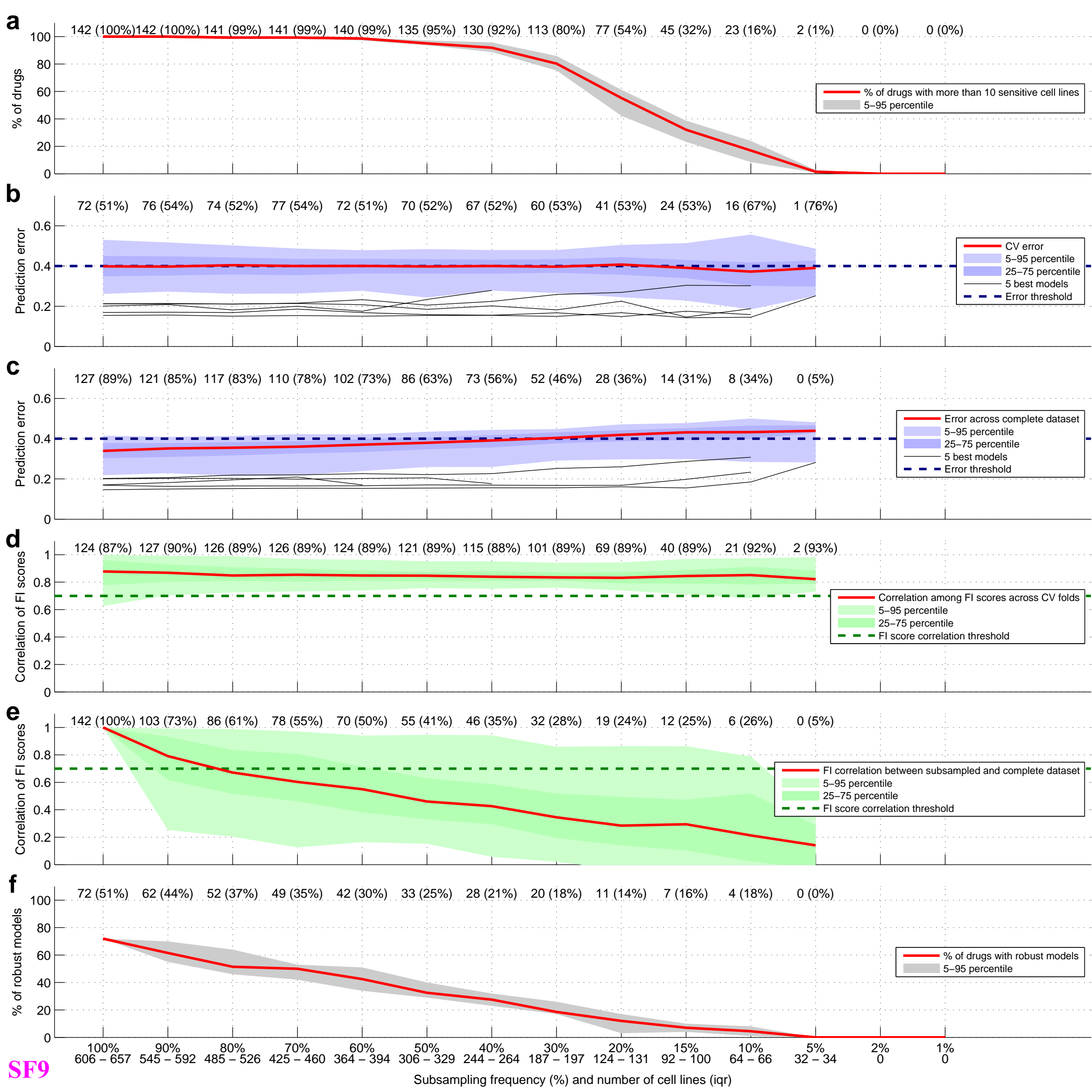


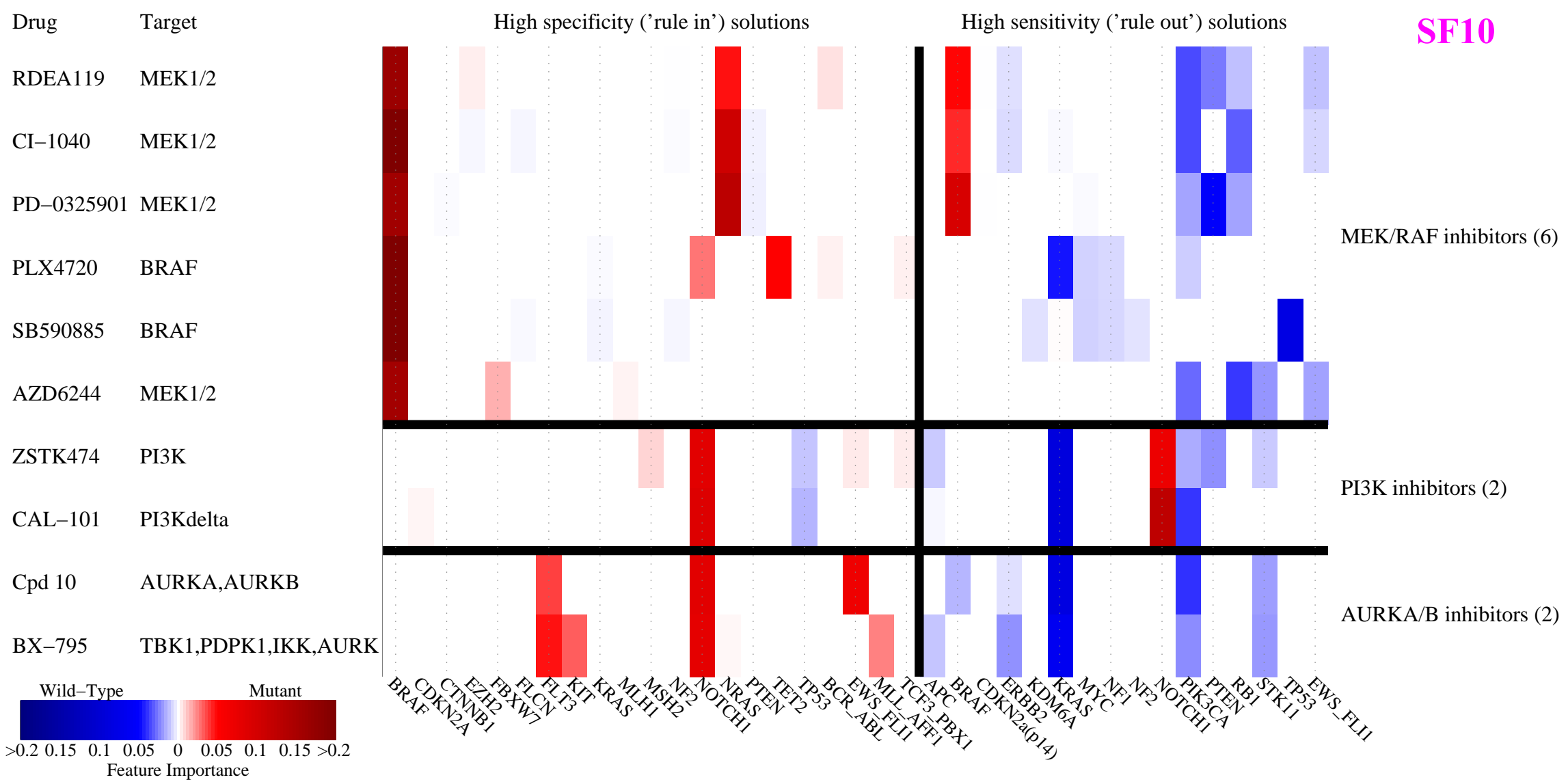
-10log p-value





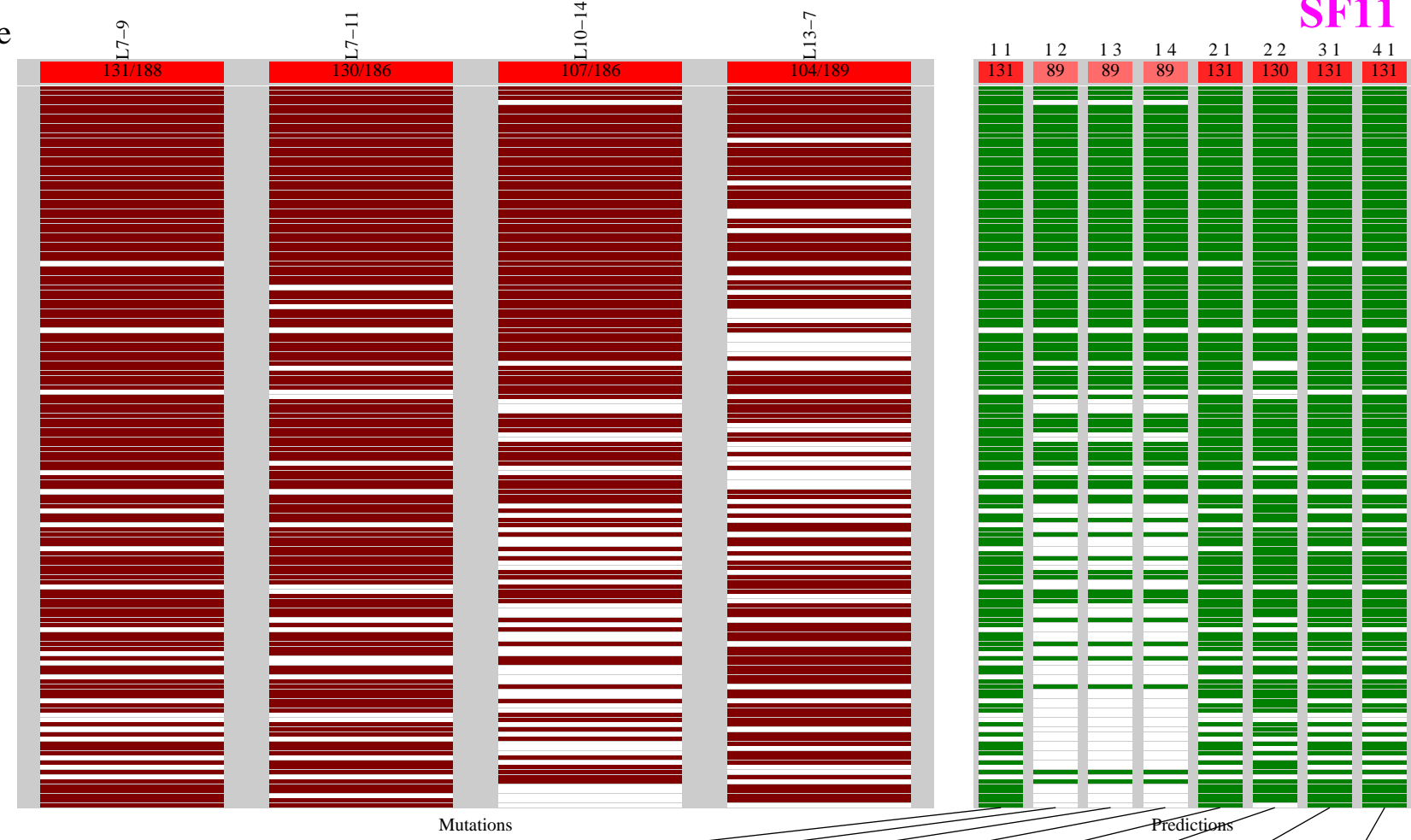
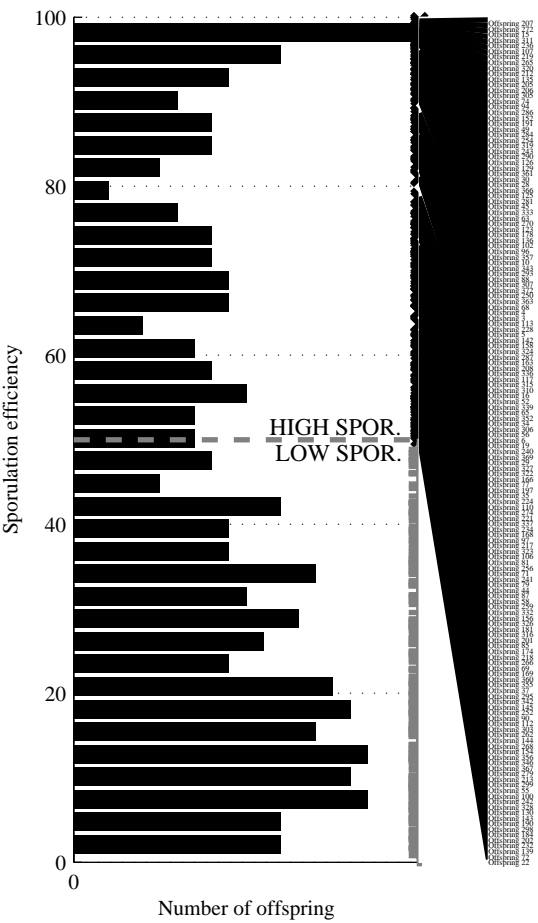




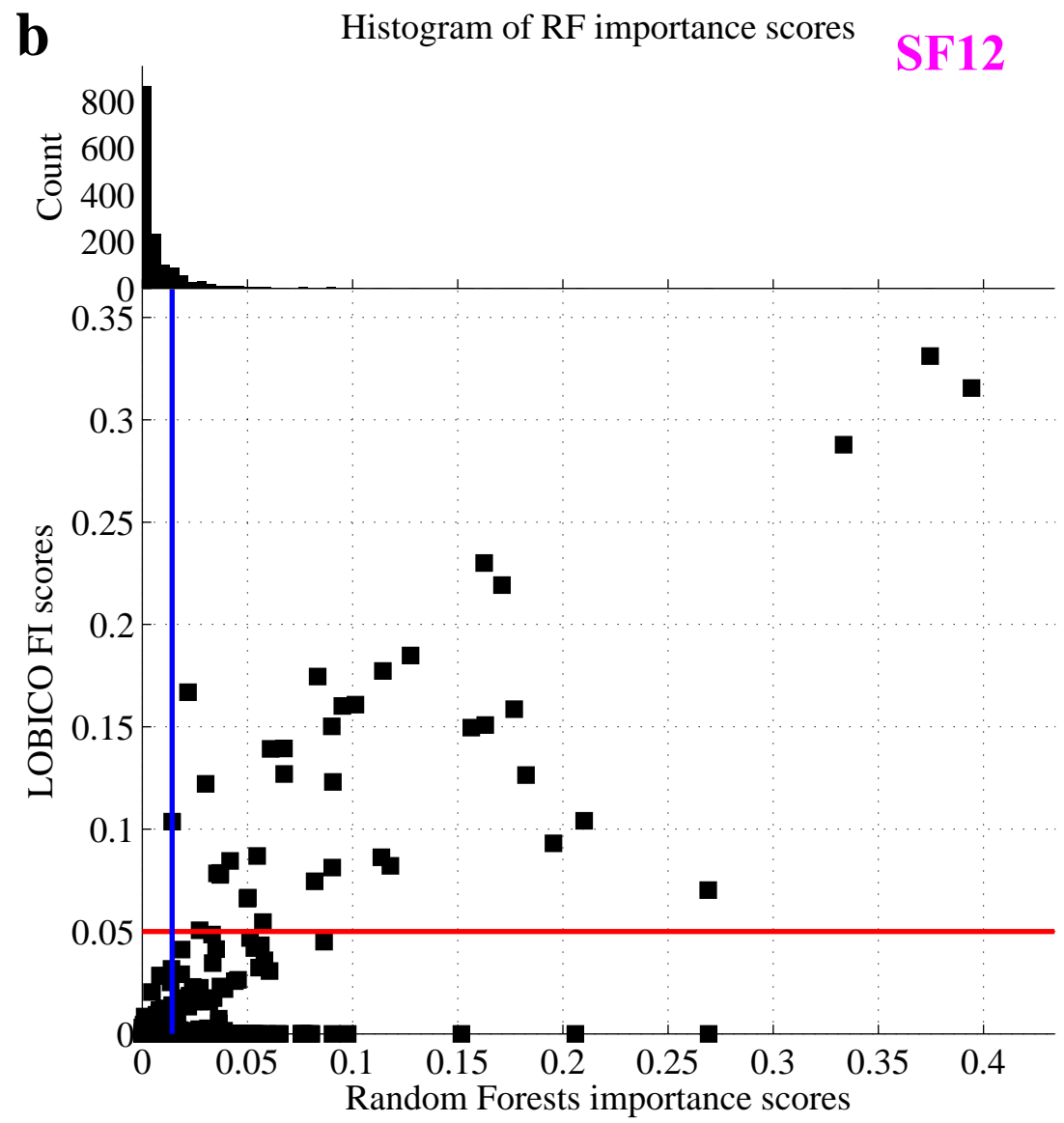
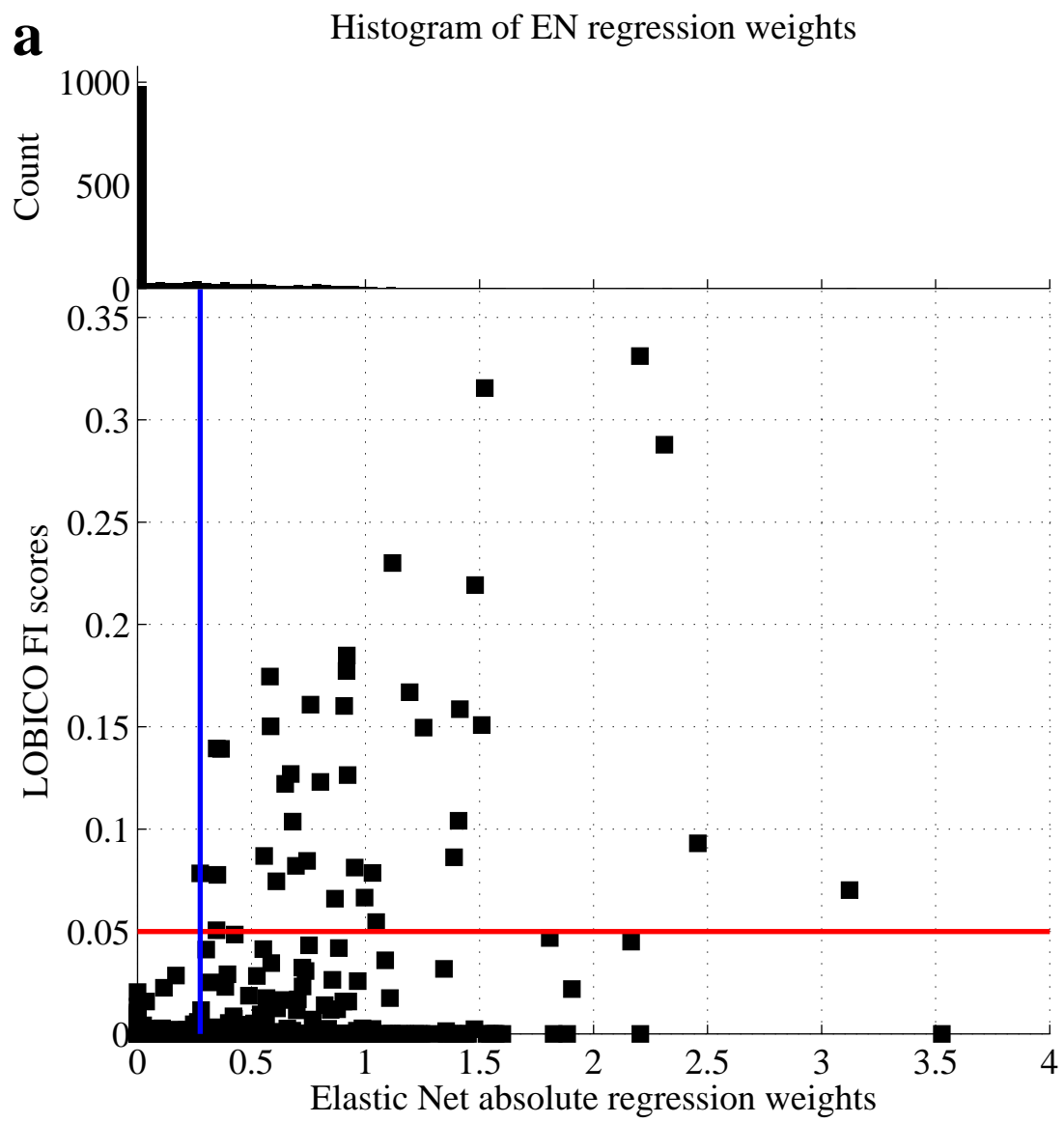


372 offspring → 152 sporulate

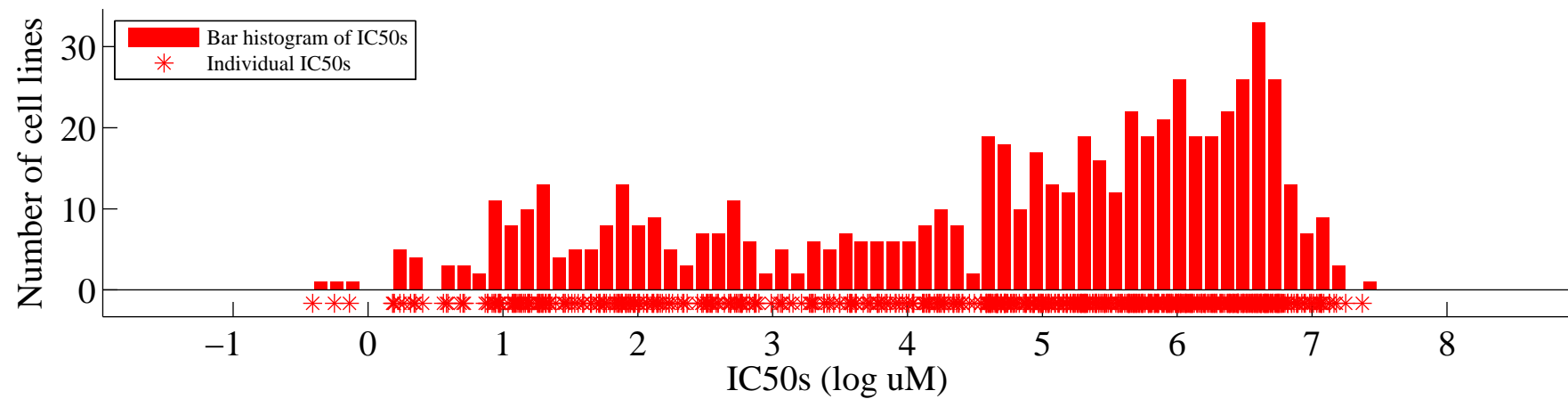
SF11



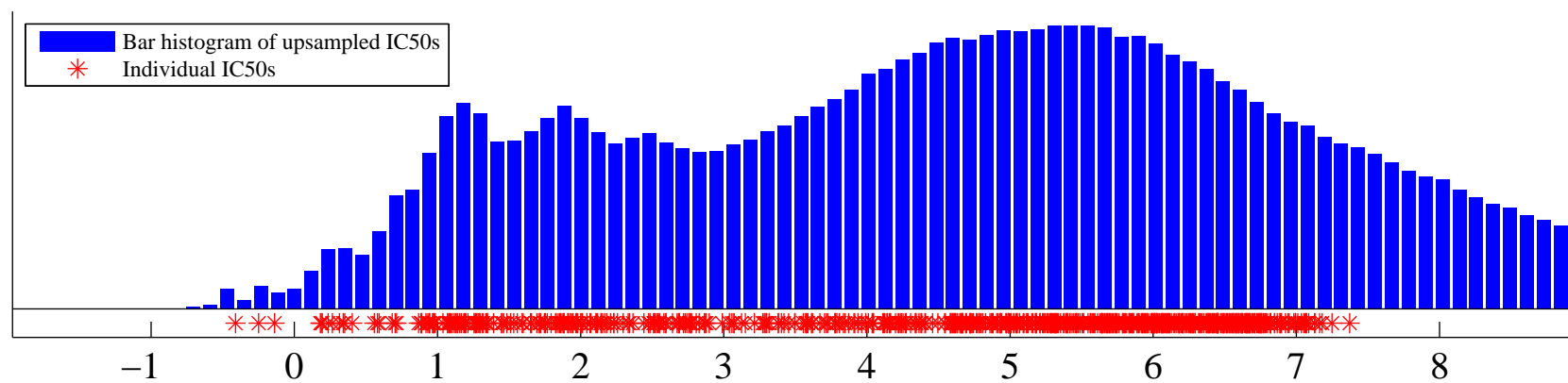
Model name	1 1		1 2		1 3		1 4		2 1		2 2		3 1		4 1	
K	M		M		M		M		M		M		M		M	
Logic formula	L7-9		L7-9 & L10-14		L7-9 & L10-14 &		L7-9 & L10-14 &		L7-9		[L7-9 & L13-7] [L7-11 & L10-14]		L7-9		L7-9	
TP FP	131 57	0.74	89 7	0.97	89 7	0.97	89 7	0.97	131 57	0.74	130 24	0.89	131 57	0.74	131 57	0.74
FN TN	21 163	0.7	63 213	0.93	63 213	0.93	63 213	0.93	21 163	0.7	22 196	0.84	21 163	0.7	21 163	0.7
Recall	0.86		0.59		0.59		0.59		0.86		0.86		0.86		0.86	



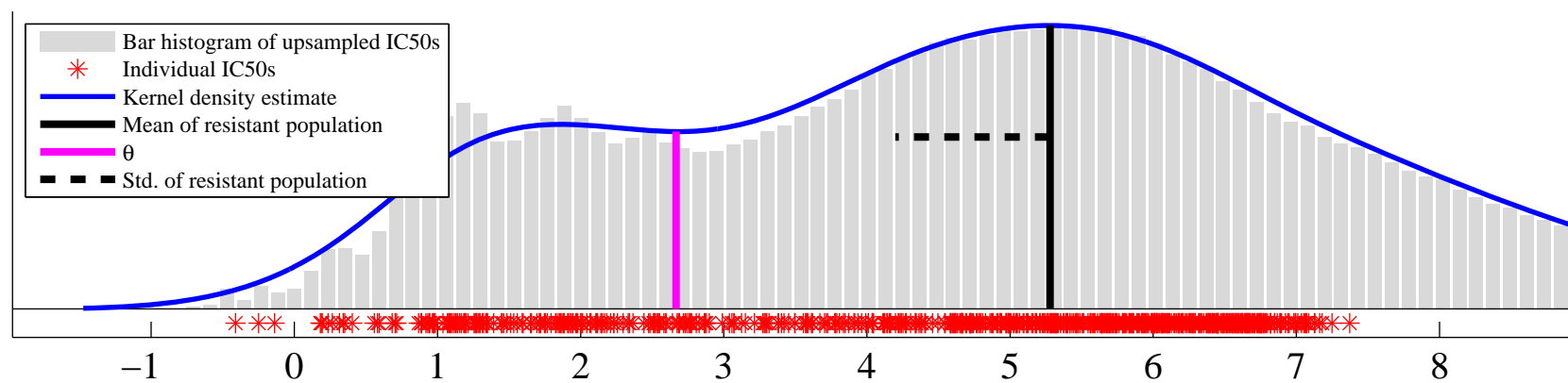
Distribution of IC50s for Nutlin-3a

a**b**

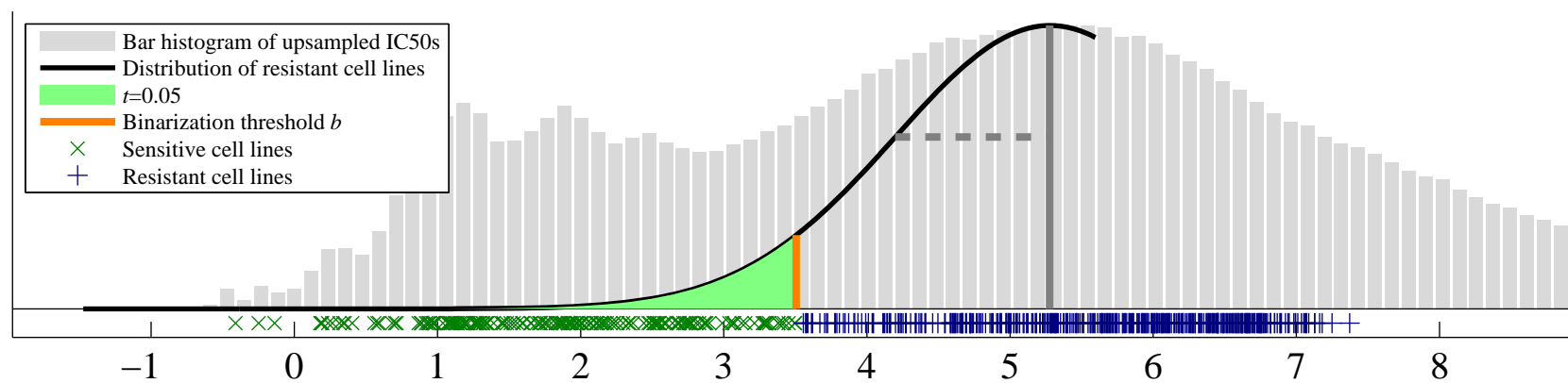
Upsampled distribution

**c**

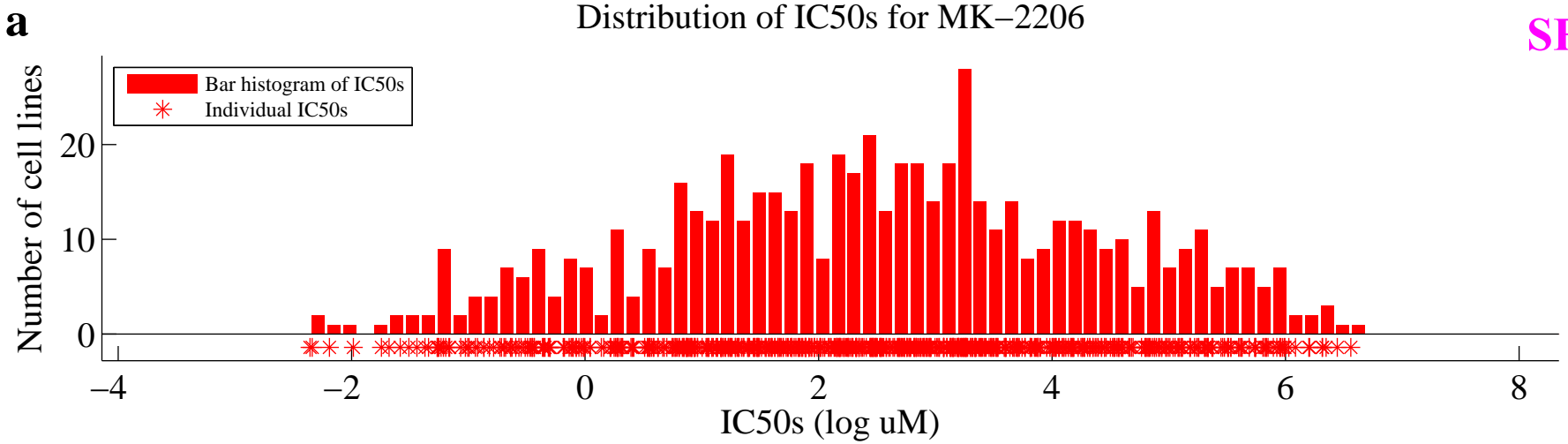
Model of population of resistant cell lines

**d**

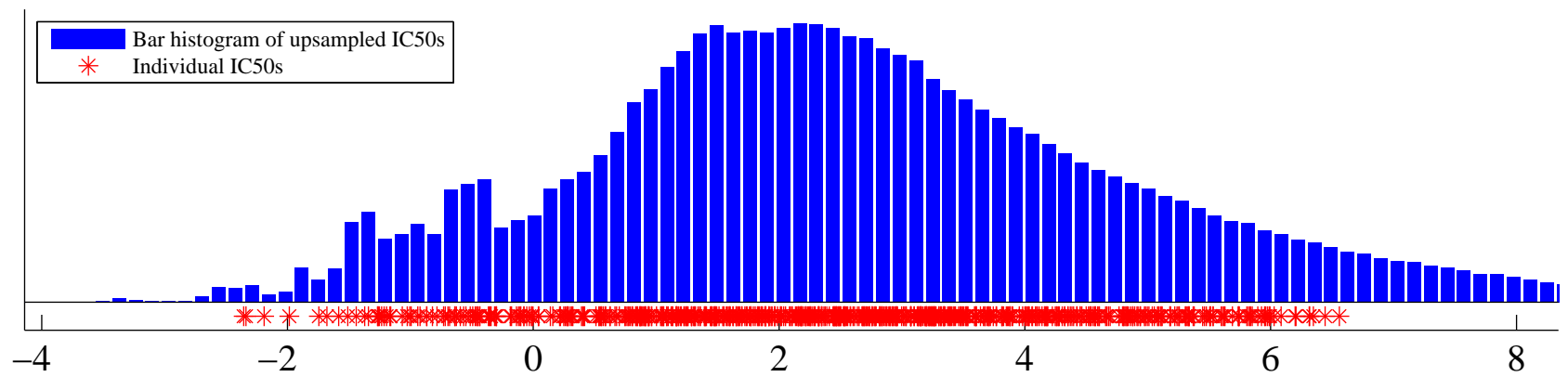
Binarization threshold



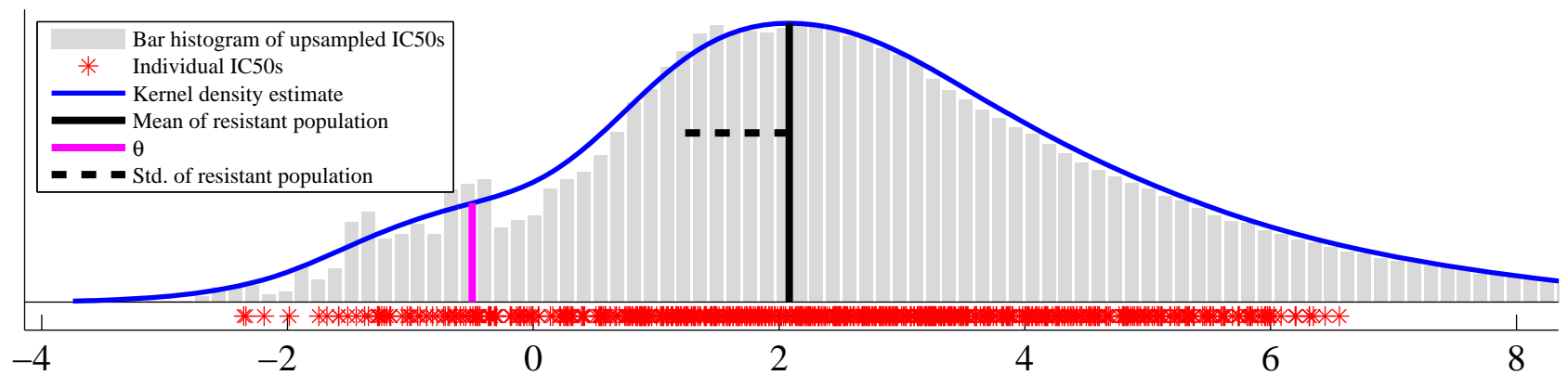
Distribution of IC50s for MK-2206



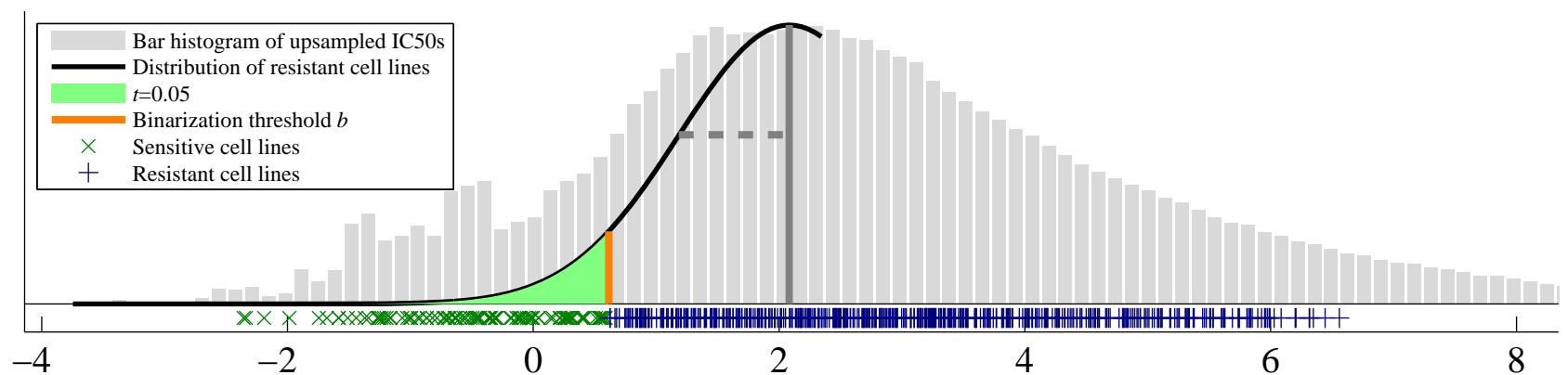
Upsampled distribution



Model of population of resistant cell lines

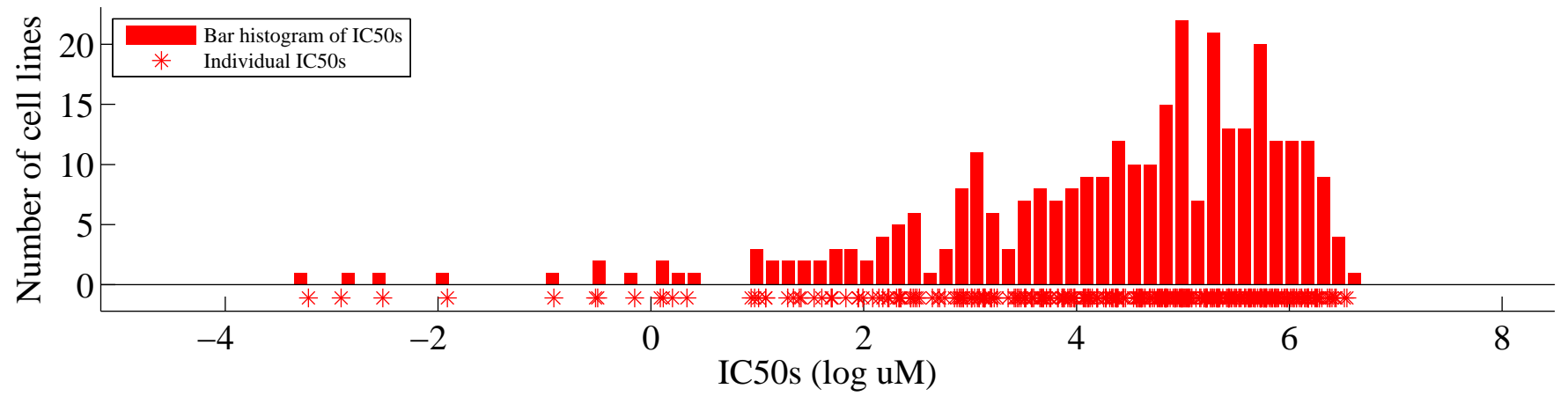


Binarization threshold

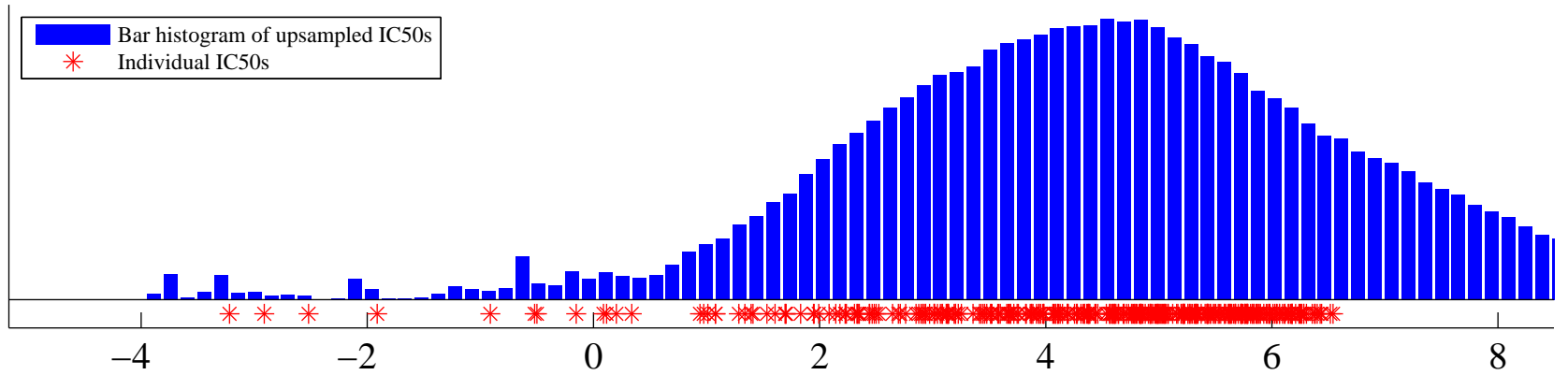


a

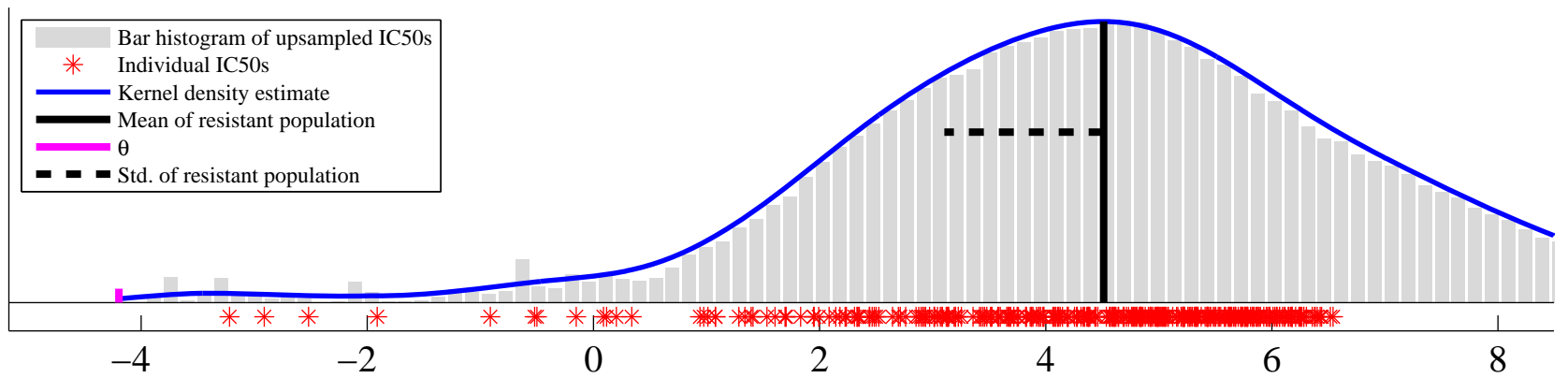
Distribution of IC50s for Erlotinib

**b**

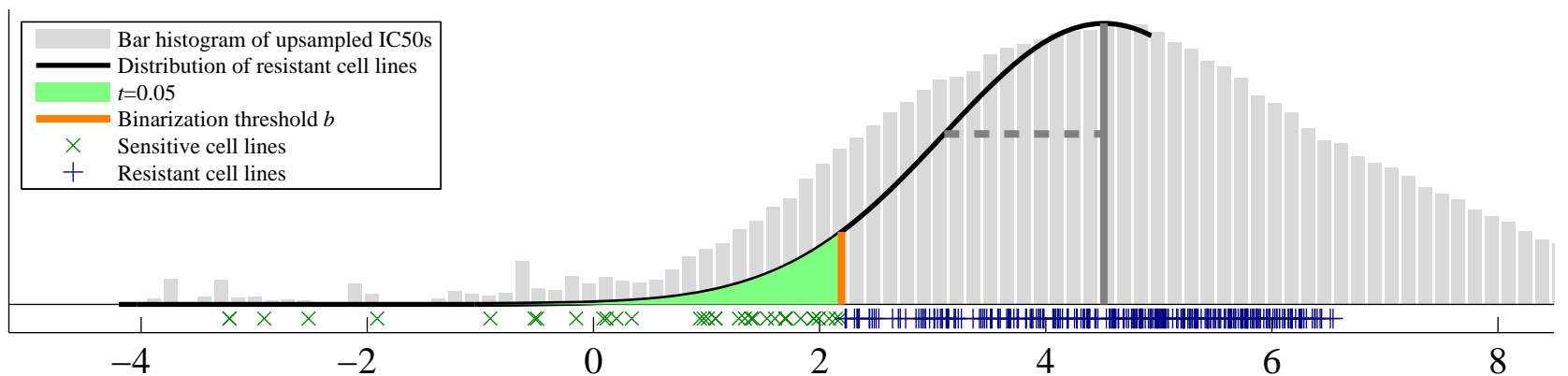
Upsampled distribution

**c**

Model of population of resistant cell lines

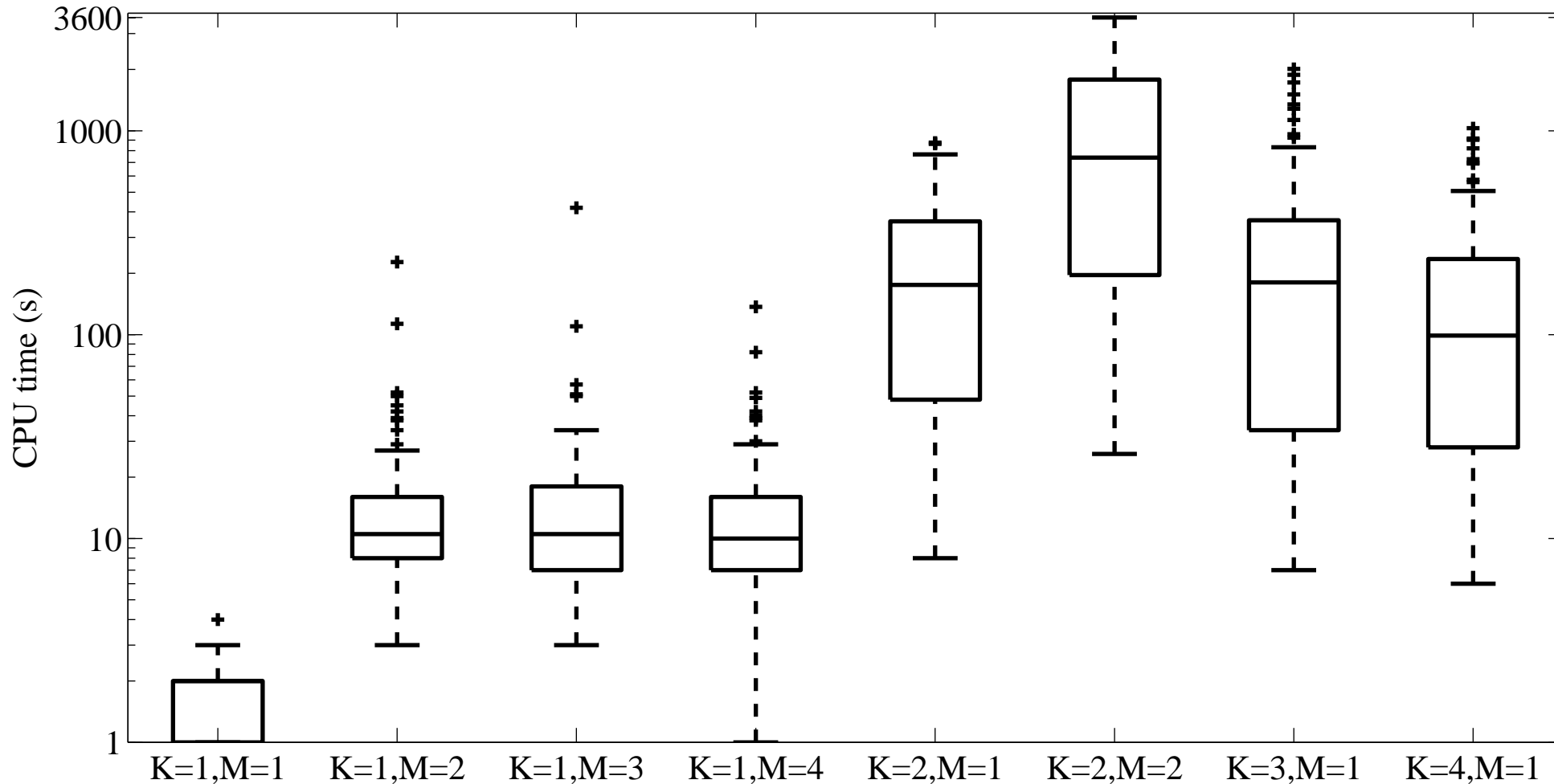
**d**

Binarization threshold



CPU time necessary to find optimal solution as a function of the model complexity

SF16



x_1 x_2 x_3 x_4 x_5 x_6 x_7 x_8 x_9 x_{10}

y



SF17

$$y = x_5 \bar{x}_8 + x_3 \bar{x}_6 x_7$$

