
On Modelability and Generalizability: Are Machine Learning Models for Drug Synergy Exploiting Artefacts and Biases in Available Data?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Synergy models are useful tools for exploring drug combinatorial search space and
2 identifying promising sub-spaces for in vitro/vivo experiments. Here, we report
3 that distributional biases in the training-validation-test sets used for predictive
4 modeling of drug synergy can explain much of the variability observed in model
5 performances (up to 0.22 $\Delta AUPRC$). We built 145 classification models spanning
6 4,577 unique drugs and 75,276 pair-wise drug combinations extracted from Drug-
7 Comb, and examined spurious correlations in both the input feature and output
8 label spaces. We posit that some synergy datasets are easier to model than others
9 due to factors such as synergy spread, class separation, chemical structural diversity,
10 physicochemical diversity, combinatorial tests per drug, and combinatorial label
11 entropy. We simulate distribution shifts for these dataset attributes and report that
12 the drug-wise homogeneity of combinatorial labels most influences modelability
13 ($0.16 \pm 0.06 \Delta AUPRC$). Our findings imply that seemingly high-performing
14 drug synergy models may not generalize well to broader medicinal space. We
15 caution that the synergy modeling community's efforts may be better expended in
16 examining data-specific artefacts and biases rigorously prior to model building.

17 1 Introduction

18 For complex, multifactorial diseases such as cancer, combination therapies offer the possibility of
19 enhanced efficacies [19], with reduced effective doses and associated host toxicities [9], as well
20 as a strategy for slowing the evolved drug resistance commonly observed in monotherapies [32].
21 It is, however, more challenging to perform clinical trials for combination therapies [22] and the
22 large number of possible drug combinations renders exhaustive testing by brute-force heuristics
23 infeasible. Machine learning is a useful tool for exploring the vast drug combinatorial search space
24 and identifying promising sub-spaces for in vitro/vivo experiments.

25 Currently, research in the field of predictive modeling for drug synergy is largely focused on model
26 generation and the optimization of performance metrics such as AUC (which overestimates model
27 performance on imbalanced datasets [30, 15]), rather than the context in which models are generated
28 and deployed. Model improvements are not reported in tandem with descriptive statistics characteriz-
29 ing the quality and modelability of datasets. Nair et al. [18] proffer that a limitation of their dataset
30 is that drug combination screens are generally discordant across independent studies. There is no
31 consensus definition for drug synergy [17, 29] and the experimental endpoints modeled are often
32 proxies of drug response that can be easily measured in a high-throughput fashion but lack clinical
33 relevance or even reproducibility [20].

34 Biases have been reported in datasets used for model generation in adjacent research fields, such as
35 PDBBind and CASF for the prediction of ligand-protein binding affinities [27]. In a systematic review
36 of 41 genomic machine learning studies, Barnett et al. [2] investigated which components of a study
37 contributed to improvements in model performance and whether reported improvements represent a
38 true improvement or an unaddressed bias inflating performance. They found that data leakage due
39 to feature selection and the number of hyperparameter optimizations were significantly associated
40 with an increase in reported model performance. In a review of 62 machine learning studies on the
41 detection and prognostication of COVID-19 using chest radiographs and chest computed tomography
42 images, Roberts et al. [26] found that none of the models identified were of potential clinical use due
43 to biases in either the methodology or underlying data.

44 Previous studies on drug synergy prediction have not examined artefacts and biases in dataset
45 composition. To the best of our knowledge, no attempt has been made to quantify the sensitivity of
46 synergy models to underlying distributions in either input feature or output label spaces. Alsherbiny
47 et al. [1] note that the source of drug combination screening data, i.e. NCI-ALMANAC [8] versus
48 ONEIL [21], has a more significant impact on model performance than feature engineering. Similarly,
49 Rani et al. [25] note that synergy models built using NCI-ALMANAC tend to outperform those built
50 using ONEIL. Here, we report that distributional biases in the datasets used for predictive modeling of
51 drug synergy explain much of the variability observed in model performances (up to 0.22 $\Delta AUPRC$).
52 We built 145 binary classification models using drug combination screens extracted from DrugComb
53 [35] spanning 4,577 unique drugs and 75,276 pair-wise drug combinations. We characterize the
54 central tendencies and dispersions of various dataset attributes, and subsequently simulate distribution
55 shifts to demonstrate that model performance can improve or deteriorate depending on the direction
56 of attribute shift.

57 **2 Methodology**

58 **2.1 Synergy Definition**

59 We use the Bliss Independence model [3], one of several synergy reference models [17, 29], to
60 qualify and quantify the expected additive or null response of administering a drug combination.
61 Operating under assumptions of statistical independence between drugs (i.e., the modes of action of
62 constituent drugs in a combination differ), symmetry in drug interactions, no variability in responses,
63 and continuous dose-response relationships, Bliss excess is defined mathematically as:

$$E_{Bliss} = E_{AB} - (E_A + E_B - E_A \times E_B)$$

64 where E_{AB} is the observed effect of the drug combination, and E_A and E_B are the observed
65 individual effects of drugs A and B, respectively. $E_{Bliss} = 0$ is the threshold for additivity, while
66 $E_{Bliss} > 0$ indicates synergy and $E_{Bliss} < 0$ indicates antagonism.

67 **2.2 Data Collection and Pre-Processing**

68 Drug pair synergy data targeting 142 cancer cell lines and 3 malarial parasites was extracted from
69 DrugComb v1.5 [35]. Thirty-three percent of drug-drug-cell line tuples were replicate experiments,
70 which we deduplicated by computing the geometric mean synergy score across replicate samples.
71 Thirty-nine percent ($N = 306,282$) of the combination-cell line tuples were sourced from NCI-
72 ALMANAC [8] and twenty-five percent ($N = 198,722$) were sourced from FRIEDMAN [12], with
73 the remainder sourced from twenty-two other combination screens including ONEIL [21] (twelve
74 percent; $N = 92,208$) and CLOUD [14] (five percent; $N = 40,160$). In total, 75,276 pair-wise drug
75 combinations comprising 4,577 unique drugs were obtained for 145 cell-line synergy endpoints
76 defined by the Bliss Independence model. We selected the top and bottom fifteen percent of each
77 cell-line dataset’s distribution of Bliss synergy scores to obtain balanced classes after filtering out
78 additive samples.

79 2.3 Dataset Attributes and Metrics

80 **Synergicity** Synergicity measures the degree to which a given drug is associated with synergistic
81 combinatorial labels: it is defined in this work, as in previous work [34], as the fraction of combi-
82 nations for which individual drugs have been labelled synergistic as opposed to antagonistic. At
83 the cell-line dataset level, the interquartile range or H-spread was used to capture the bimodality
84 of synergicity distributions and test the hypothesis that cell-line datasets with drugs found pri-
85 marily in antagonistic-only combinations ($\text{synergicity} = 0$) and synergistic-only combinations
86 ($\text{synergicity} = 1$) are easier to model with higher AUPRC scores.

87 **Combinatorial Label Entropy** Combinatorial label entropy measures the level of disorder or
88 heterogeneity of combinatorial labels. It is defined mathematically as Shannon entropy:

$$H(X) = - \sum_{i=1}^n P(x_i) \log_2(P(x_i))$$

89 where $H(X)$ is the Shannon entropy of a discrete random variable X and $P(x_i)$ is the probability of
90 outcome x_i occurring in the system. The sum is taken over all n possible outcomes x_i . In our case,
91 $H(X)$ has range $[0, 1]$ and measures how homogeneous the combinatorial labels associated with a
92 given drug are: if a drug occurs predominantly in drug combinations labelled synergistic-only or
93 antagonistic-only, then its combinatorial label entropy is low (close to 0); if a drug occurs in drug
94 combinations labelled synergistic approximately half of the time and antagonistic approximately half
95 of the time, then its combinatorial label entropy is high (close to 1).

96 **Feature Similarity** Feature similarity in chemical structural and physicochemical spaces was
97 defined in two steps: cosine similarity computed pair-wise amongst all drugs tested per cell line,
98 followed by the cell-line fraction of pair-wise similarities above 0.15. Mathematically, the cosine
99 similarity between two feature vectors A and B is defined as:

$$\text{cosine_similarity}(A, B) = \frac{A \cdot B}{\|A\| \cdot \|B\|}$$

100 **Non-Additivity** A drug’s tendency for non-additivity when combined was scored as the median
101 absolute distance from Bliss additivity across combinations. This measure was used to test the
102 hypothesis that a drug’s combinatorial label entropy decreases with its tendency for non-additivity
103 in combinations. In other words, non-additivity thus defined was used to test whether the degree of
104 synergism or antagonism achieved by a drug was associated with the consistency or homogeneity of
105 its combinatorial labels.

106 2.4 Model Generation and Evaluation

107 We formulate drug synergy prediction as a supervised classification task: we construct one binary
108 model per cell-line dataset, resulting in a total of 145 binary models, to predict synergistic versus
109 antagonistic class labels for drug-drug pairs using the CRAN "randomForest" [13, 24] implementation
110 of the traditional random forest learner by Breiman [4] under default hyperparameter optimizations.
111 Given that the focus of this work is the influence of dataset composition on model performance,
112 and not the influence of model architecture on model performance, we required a single learner
113 to serve as our baseline before and after shifting attribute distributions. We deliberately chose a
114 decision tree ensemble learner as our baseline due to its computational efficiency on high-dimensional
115 data, adequate interpretability and explainability, as well as state-of-the-art model performance on
116 balanced and minority classes [6]. We constructed two sets of drug features: structural 2048-bit
117 Morgan fingerprints (with radius 3) and 43-element long physicochemical profiles of all available
118 molecular descriptors on RDKit [11]. Feature vectors were concatenated for each drug-drug pair
119 in both permutations. Our 80%-20% train-test split strategy was drug-pair-stratified with five-fold
120 cross-validation. To evaluate model performance, we computed Area under the Precision-Recall
121 curve (AUPRC), which is less sensitive to class imbalance and thus more practically relevant and

122 actionable than Area under the Receiver Operating Characteristic curve (AUROC) [30, 15]. The mean
123 AUPRC across all models ($n = 145$) was 0.76 ± 0.09 . For our categorical analyses, we categorized
124 cell-line models with AUPRC greater than or equal to 0.8 as high-performing ($n = 50$), and cell-line
125 models with AUPRC less than 0.8 as low-performing ($n = 95$).

126 2.5 Simulating Distribution Shifts in Dataset Attributes

127 We simulated distribution shifts in dataset attributes by sub-sampling each cell-line dataset. For
128 originally high-performing models, we selected subsets of drugs with high combinatorial label
129 entropy (upper 15%), few combinatorial tests per drug (lower 15%), low physicochemical similarity
130 to other drugs (lower 15%), and low structural similarity to other drugs (lower 15%). Conversely, for
131 originally low-performing models, we selected subsets of drugs with low combinatorial label entropy
132 (lower 15%), many combinatorial tests per drug (upper 15%), high physicochemical similarity to
133 other drugs (upper 15%), and high structural similarity to other drugs (upper 15%). This simulated
134 shifts in attribute distributions such that high-performing models now resembled low-performing
135 models, and vice versa. Cell-line models with insufficient drugs remaining were discarded, yielding
136 103 models for structural similarity, 109 models for physicochemical similarity, 117 models for
137 combinatorial tests per drug, and 91 models for combinatorial label entropy per drug. The simulations
138 were run for each of the dataset attributes identified individually, as well as pair-wise, but the latter
139 yielded datasets too small for model generation. To distinguish change in model performance due to
140 shifting bias versus reduction in dataset size, models were trained, validated, and tested on shifted
141 and non-shifted subsets of comparable size for each cell line.

142 3 Results

143 3.1 Synergy Spread and Class Separation

144 We first analyzed the effect of dataset span, measured as standard deviation of Bliss synergy scores,
145 and class separation, measured as difference in mean Bliss synergy scores of antagonistic vs synergistic
146 classes, on cell-line model performance, measured as AUPRC. The results are shown in Figure 1.
147 It can be seen that high-performing cell-line models tended to exhibit broader synergy spread with
148 difference in means between high- and low-performing models of 15.4–24.1 (95% CI) Bliss synergy
149 units (Welch’s two-sample $t = 9.13$, $df = 71.3$, $p = 1.26e-13$). This is consistent with the relationship
150 between potency span and achievable model performance reported by Brown et al. [5] in the context
151 of predicting binding affinity of small-molecule ligands for protein targets. High-performing cell-line
152 models also tended to exhibit greater class separation in synergy space with difference in means
153 between high- and low-performing models of 12.9–17.6 (95% CI) Bliss synergy units (Welch’s
154 two-sample $t = 13.1$, $df = 94.4$, $p < 2.20e-16$). Easier class splits may inflate model performance,
155 particularly on AUROC [30, 15] but also AUPRC: DeepSynergy, for instance, defined the top 10%
156 of combinations as the synergistic or positive class and modeled the remainder as the negative class
157 [23]. Our findings show that both synergy spread and class separation influence modelability.

158 3.2 Synergicity and Entropy of Combinatorial Labels

159 We then analyzed the effect of combinatorial label homogeneity on model performance (Sub-Figures
160 2A-B). It can be seen that the cell-line H-spread of synergicity, defined as the fraction of combinations
161 for which individual drugs have been labelled synergistic as opposed to antagonistic, is positively
162 correlated with cell-line model performance, measured as AUPRC (Spearman’s $\rho = 0.539$, $p =$
163 $1.77e-10$). Conversely, the cell-line arithmetic mean heterogeneity of combinatorial labels, measured
164 as Shannon entropy for individual drugs, is negatively correlated with cell-line model performance,
165 measured as AUPRC (Pearson’s $r = -0.691$, $p < 2.20e-16$). The more bimodal a cell line’s drug
166 synergicity distribution, the more homogeneous its drug-wise combinatorial labels and the easier
167 to predict combinations unseen during training with at least one seen-before drug. Our findings
168 imply that cell lines comprising drugs with homogeneous combinatorial labels, i.e., drugs occurring

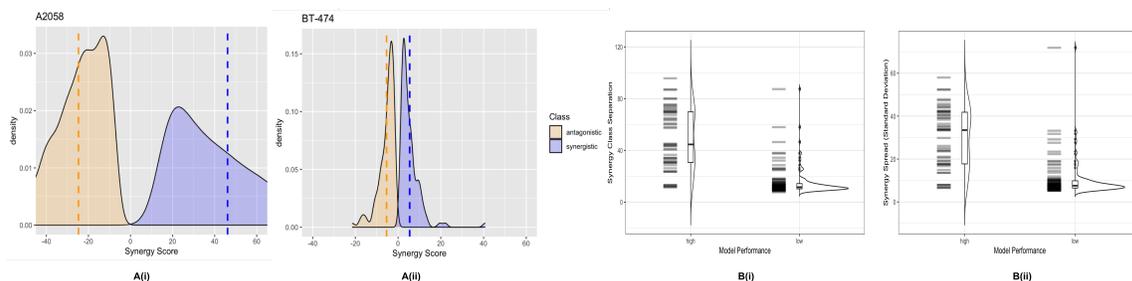


Figure 1: **Panel A.** Distribution of Bliss synergy scores for the best-performing cell-line model, **A(i)**, and the worst-performing cell-line model, **A(ii)**. **Panel B.** Each barcode line in the violin plots represents one cell-line model. Differences in synergy class means, **B(i)**, and standard deviations of overall synergy distributions, **B(ii)**, for all cell-line models binned into high versus low AUPRCs

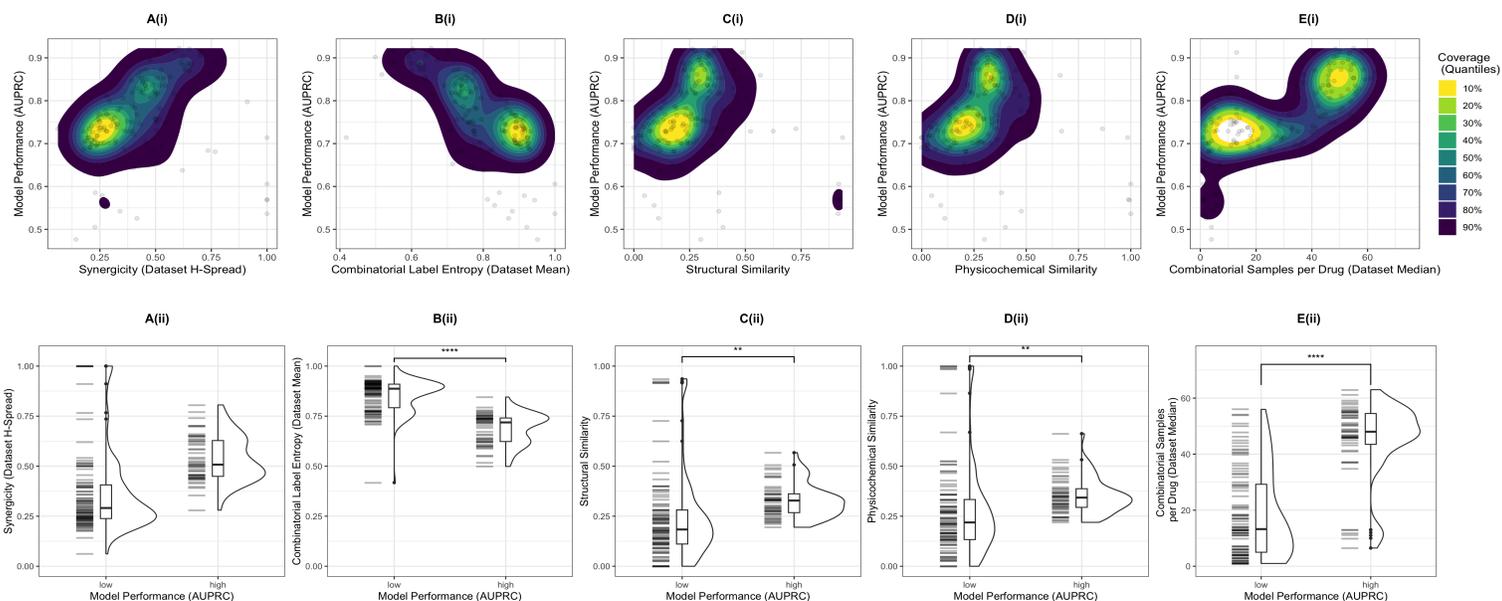


Figure 2: **Panel A.** Density and violin plots of cell-line H-spread of the fraction of combinations for which individual drugs have been labelled synergistic (dubbed synergicity) and cell-line model performance (Spearman's $\rho = 0.539$). **Panel B.** Density and violin plots of cell-line mean combinatorial label entropy and cell-line model performance (Pearson's $r = -0.691$). High-performing cell-line models exhibited lower diversity spanning 3.91%–13.8% (95% *CI*) higher cosine similarity in structural space with Spearman's $\rho = 0.359$ (**Panel C**) and 2.28%–12.9% (95% *CI*) higher cosine similarity in physicochemical space with Spearman's $\rho = 0.327$ (**Panel D**), as well as 17.1–31.0 (95% *CI*) more combinations tested per drug with Pearson's $r = 0.504$ (**Panel E**). Each dot in the density plots (upper panels) and each barcode line in the violin plots (lower panels) represents one cell-line model.

169 primarily in antagonistic-only combinatorial labels and synergistic-only combinatorial labels, tend to
 170 be easier to model with higher AUPRC scores.

171 3.3 Structural Diversity, Physicochemical Diversity, Combinatorial Tests Per Drug

172 We then analyzed the effects of drug diversity in structural Morgan fingerprint and physicochemical
 173 spaces, both measured as fraction of drugs in a cell-line dataset with pair-wise cosine similarity
 174 above a defined threshold, on cell-line model performance, measured as AUPRC. Panel C of Figure 2
 175 shows that the dataset attribute, compound structural similarity, is positively correlated with model

176 performance (Spearman’s $\rho = 0.359$, $p = 1.012e-05$): high-performing cell-line models exhibited
 177 3.91%–13.8% (95% CI) higher pair-wise cosine similarity between drugs in Morgan fingerprint
 178 space than low-performing cell-line models (Welch’s two-sample $t = 3.54$, $df = 132.64$, $p = 0.0005$).
 179 Similarly, Panel D of Figure 2 shows that the dataset attribute, compound physicochemical similarity,
 180 is positively correlated with model performance (Spearman’s $\rho = 0.327$, $p = 6.282e-05$): high-
 181 performing cell-line models exhibited 2.28%–12.9% (95% CI) higher pair-wise cosine similarity
 182 between drugs in physicochemical space than low-performing cell-line models (Welch’s two-sample t
 183 $= 2.83$, $df = 131.33$, $p = 0.005$). Summarily, the breadth of compound structural and physicochemical
 184 spaces both appear to influence modelability, which one might expect as it is easier to model a smaller
 185 space with greater overlap between train and validation/test sets. We subsequently investigated the
 186 relationship between cell-line model performance, measured as AUPRC, and number of combinatorial
 187 tests per drug. It can be seen in Panel E of Figure 2 that this dataset attribute is positively correlated
 188 with model performance (Pearson’s $r = 0.504$, $p = 1.24e-10$). High-performing cell-line models
 189 comprized 17.1-31.0 (95% CI) more combinations tested per drug than low-performing cell-line
 190 models (Welch’s two-sample $t = 6.86$, $df = 141.19$, $p = 1.99e-10$), which one might expect as it
 191 is easier to model a smaller space with fewer distinct drugs tested in more combinations. These
 192 findings imply that seemingly high-performing drug synergy models do not generalize well to broader
 193 medicinal space.

194 3.4 Simulating Distribution Shifts in Dataset Attributes

195 To test whether the differences in model performance observed across cell lines was due to underlying
 196 data modelability versus biological variability, we simulated shifts in dataset attribute distributions
 197 and compared resulting changes in model performance ($\Delta AUPRC$). We selected subsets of drug-drug
 198 samples to shift distributions for low-performing cell-line models to resemble high-performing cell-
 199 line models, and vice versa. The simulations were run for each of the dataset attributes identified
 200 individually, as well as pair-wise, but the latter yielded datasets too small for model generation. The
 201 results are summarized in Figure 3.

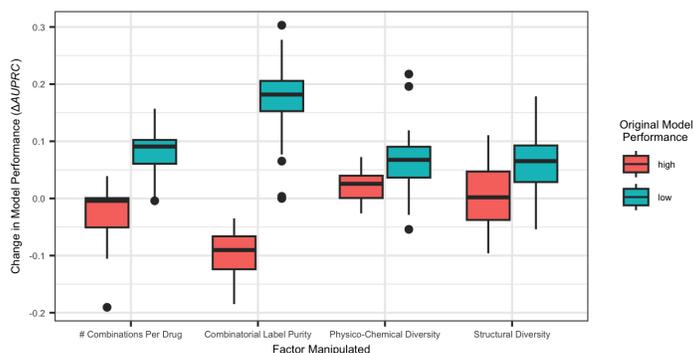


Figure 3: Change in model performance, $\Delta AUPRC$, after simulating distribution shifts for each dataset attribute individually. Attribute distributions for previously low-performing cell-line models were shifted to resemble attribute distributions for high-performing cell-line models, and vice versa. Performance improved for previously low-performing models (blue) under all simulations, albeit to varying degrees ($+0.06 \pm 0.04 \Delta AUPRC$ for physicochemical diversity versus $+0.18 \pm 0.05 \Delta AUPRC$ for combinatorial label entropy). Performance deteriorated most noticeably for previously high-performing models (red) following shifts in distributions for combinatorial label entropy ($-0.10 \pm 0.04 \Delta AUPRC$).

202 It can be seen that subsetting data points that result in greater class separation, broader synergy spread,
 203 lower structural diversity, lower physicochemical diversity, higher number of combinatorial tests per
 204 drug, and lower combinatorial label entropy generally increased model performance. Conversely,
 205 subsetting data points that result in smaller class separation, narrower synergy spread, lower number
 206 of combinatorial tests per drug, and higher combinatorial label entropy generally decreased model

207 performance. In other words, simulating shifts in attribute distributions tended to boost model
 208 performance for originally low-performing models, and tended to degrade model performance for
 209 originally high-performing models. This suggests that the differences observed in model performance
 210 across cell lines was likely due to differences in dataset composition and not due to inherent biological
 211 variation. Of the dataset attributes identified and manipulated, combinatorial label entropy most
 212 influenced modelability, increasing the performance of originally low-performing models by $+0.18 \pm$
 213 $0.05 \Delta AUPRC$, which is comparable to the original difference in mean performance between high-
 214 versus low-performers ($0.15 \Delta AUPRC$). It is important to note that factors are not decoupled in these
 215 simulations as shifting one attribute distribution in isolation was not feasible; shifting one distribution
 216 simultaneously shifted other distributions to varying degrees since we must also consider how dataset
 217 attributes are correlated with each other. To contextualize these findings, we refer to improvements
 218 over state-of-the-art models reported in drug synergy literature, such as $+0.04 \Delta AUPRC$ by Preuer
 219 et al. [23] and Wang et al. [31].

220 3.5 Synergy, Lipophilicity, and Model Performance

221 We then analyzed whether mechanistic insights reported in drug synergy literature, particularly the
 222 relationship between synergy and lipophilicity [34], influence modelability. Figure 4A shows
 223 that, for the well-characterized cell line MCF7, a drug’s lipophilicity (CrippenClogP) is positively
 224 correlated with its synergy, measured as the fraction of combinations for which the drug has been
 225 experimentally labelled synergistic as opposed to antagonistic, particularly in the region most relevant
 226 for drug discovery, i.e., CrippenClogP interval (1,6). Figure 4B shows the correlation between
 227 lipophilicity and synergy for all cell lines plotted against model performance (Spearman’s $\rho =$
 228 -0.351 , $p = 1.575e-05$): high-performing models evidently do not rely on the positive correlation
 229 between lipophilicity and synergy reported here and in literature [34] for predictions.

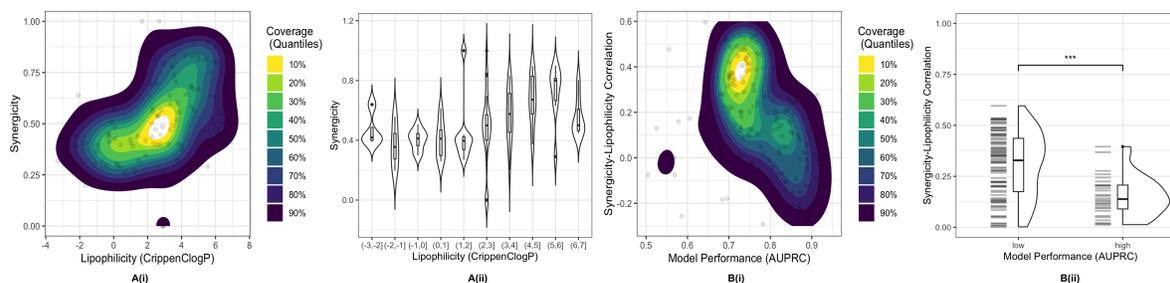


Figure 4: **Panel A.** A drug’s lipophilicity (CrippenClogP) is correlated with its synergy in the MCF7 cell-line dataset, particularly for drug-like molecules in CrippenClogP interval (1,6). **Panel B.** Correlation between lipophilicity and synergy plotted as a function of model performance for all cell-line datasets. High-performing models evidently do not rely on the correlation between lipophilicity and synergy reported here and in literature for predictions.

230 3.6 Non-Additivity, Combinatorial Label Homogeneity, Drug Similarity

231 We considered the dependence of combinatorial label homogeneity, an output dataset attribute, on
 232 various input dataset attributes, such as drug similarity. It can be seen in Appendix Figure 7 that
 233 cell-line drug similarity in physicochemical (Pearson’s $r = 0.480$) and structural (Pearson’s $r =$
 234 0.514) spaces correlate with combinatorial label homogeneity. A drug is more likely to behave
 235 generally synergistically or generally antagonistically, or rather elicit mostly synergistic-only or
 236 antagonistic-only labels, when combined with similar drugs, since similar drugs hit similar pathways
 237 exhibiting homogeneous synergistic *or* antagonistic effect. Different drugs hit different pathways
 238 exhibiting heterogeneous synergistic *and* antagonistic effect: synergy with some drugs and antagonism
 239 with other drugs depending on pathway hit [16]. We then considered the relationship between a
 240 drug’s combinatorial label homogeneity and its tendency for non-additivity, defined in this work

241 as median absolute distance from Bliss additivity across combinations. The correlation between
 242 these attributes varied across cell-line models and tended to increase with dataset modelability or
 243 increasing model performance in AUPRC (Pearson’s $r = 0.378$, Figure 5A). High-performing cell-line
 244 models comprized drugs exhibiting a stronger correlation between combinatorial label homogeneity
 245 and non-additivity with a 95% CI [0.091,0.241] higher Pearson correlation coefficient (PCC) than
 246 low-performing cell-line models (Welch’s two-sample $t = 4.39$, $df = 114.7$, $p < 0.00002$). 19.4%
 247 of cell-line datasets exhibited PCCs between combinatorial label homogeneity and non-additivity
 248 ≥ 0.5 . Of these, 75% had model performances AUPRC ≥ 0.8 . Figure 5B shows one such cell-line
 249 dataset, namely the skin epithelial-like cell line IST-MEL1, with AUPRC ≥ 0.9 and PCC between
 250 combinatorial label homogeneity and non-additivity $r = 0.643$. In other words, drugs that elicited
 251 close-to-additive effects when combined tended to have low combinatorial label homogeneity, while
 252 drugs that elicited highly synergistic or highly antagonistic effects when combined tended to have
 253 high combinatorial label homogeneity. These findings imply that combinatorial label homogeneity
 254 could function as a crude proxy for non-additivity in some contexts, yielding greater modelability.

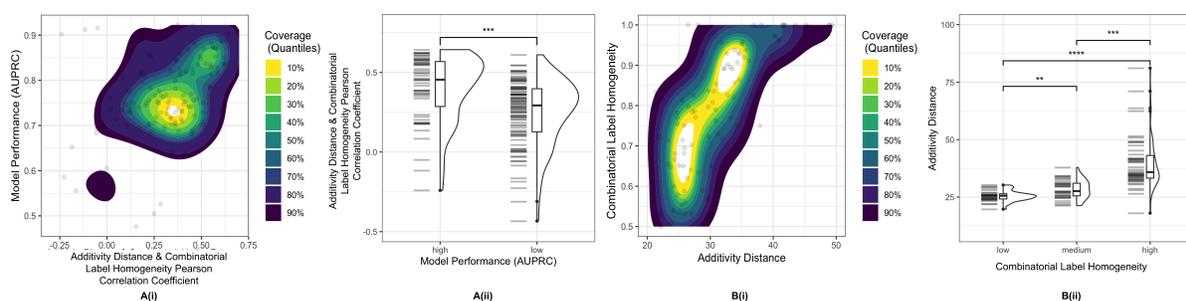


Figure 5: **Panel A.** Each dot in the density plot and each barcode line in the violin plot represents one cell-line model. **Panel A(i).** Model performance (AUPRC) tended to increase with increasing strength of correlation between combinatorial label homogeneity and degree of non-additivity (Pearson’s $r = 0.378$). **Panel A(ii).** High-performing cell-line models spanned drugs with a stronger correlation between combinatorial label homogeneity and degree of non-additivity: 95% CI [0.091,0.241] difference in mean PCCs. **Panel B.** Combinatorial label homogeneity versus degree of non-additivity for the IST-MEL1 cell line with AUPRC ≥ 0.9 (Pearson’s $r = 0.643$).

255 4 Conclusions

256 In this work, we qualify and quantify various synergy dataset attributes influencing modelability:
 257 synergy spread, class separation, chemical structural diversity, physicochemical diversity, combina-
 258 torial tests per drug, and combinatorial label entropy. We simulate shifts in distributions of these
 259 attributes and report that combinatorial label entropy improved and degraded model performance
 260 most, depending on the direction of attribute shift. It is important to note that the attributes were
 261 not decoupled in our simulations as shifting one attribute distribution in isolation was not feasible;
 262 shifting one distribution simultaneously shifted other distributions to varying degrees. Overall, our
 263 findings imply that model performance is highly sensitive to distributional biases in available data.
 264 We find that distributional biases in the training-validation-test sets used for predictive modeling of
 265 drug synergy can explain up to 0.22 $\Delta AUPRC$ of the difference observed in model performances.
 266 For comparison, we refer to performance improvements over state-of-the-art models reported in drug
 267 synergy literature, such as 0.04 $\Delta AUPRC$ by Preuer et al. [23] and Wang et al. [31]. We caution
 268 that the synergy modeling community’s efforts may be better expended in examining data-specific
 269 artefacts and biases rigorously prior to model building. We recommend that synergy modelers
 270 characterize the applicability domain wherein models can be expected to work reliably and report
 271 explicitly the statistical biases underlying datasets used for model generation.

272 References

- 273 [1] Alsherbiny, M. A., Radwan, I., Moustafa, N., Bhuyan, D. J., El-Waisi, M., Chang, D. and Li, C. G. [2023],
274 'Trustworthy Deep Neural Network for Inferring Anticancer Synergistic Combinations', *IEEE Journal of*
275 *Biomedical and Health Informatics* **27**(4), 1691–1700.
- 276 [2] Barnett, E., Onete, D., Salekin, A. and Faraone, S. V. [2022], 'Genomic Machine Learning Meta-regression:
277 Insights on Associations of Study Features with Reported Model Performance', *medRxiv* pp. 2022–01.
- 278 [3] Bliss, C. I. [1939], 'The Toxicity of Poisons Applied Jointly', *Annals of applied biology* **26**(3), 585–615.
- 279 [4] Breiman, L. [2001], 'Random forests', *Machine learning* **45**, 5–32.
- 280 [5] Brown, S. P., Muchmore, S. W. and Hajduk, P. J. [2009], 'Healthy skepticism: assessing realistic model
281 performance', *Drug discovery today* **14**(7-8), 420–427.
- 282 [6] Chen, J., Wu, L., Liu, K., Xu, Y., He, S. and Bo, X. [2023], 'EDST: a decision stump based ensemble
283 algorithm for synergistic drug combination prediction', *BMC bioinformatics* **24**(1), 1–21.
- 284 [7] Cortés-Ciriano, I. and Bender, A. [2016], 'How consistent are publicly reported cytotoxicity data? Large-
285 scale statistical analysis of the concordance of public independent cytotoxicity measurements', *ChemMed-*
286 *Chem* **11**(1), 57–71.
- 287 [8] Holbeck, S. L., Camalier, R., Crowell, J. A., Govindharajulu, J. P., Hollingshead, M., Anderson, L. W.,
288 Polley, E., Rubinstein, L., Srivastava, A., Wilsker, D. et al. [2017], 'The National Cancer Institute
289 ALMANAC: a comprehensive screening resource for the detection of anticancer drug pairs with enhanced
290 therapeutic activity', *Cancer research* **77**(13), 3564–3576.
- 291 [9] Jia, J., Zhu, F., Ma, X., Cao, Z. W., Li, Y. X. and Chen, Y. Z. [2009], 'Mechanisms of drug combinations:
292 interaction and network perspectives', *Nature reviews Drug discovery* **8**(2), 111–128.
- 293 [10] Khetan, R., Curtis, R., Deane, C. M., Hadsund, J. T., Kar, U., Krawczyk, K., Kuroda, D., Robinson, S. A.,
294 Sormanni, P., Tsumoto, K. et al. [2022], Current advances in biopharmaceutical informatics: guidelines,
295 impact and challenges in the computational developability assessment of antibody therapeutics, in 'MAbs',
296 Vol. 14, Taylor & Francis, p. 2020082.
- 297 [11] Landrum, G., Tosco, P., Kelley, B., Sriniker, Gedeck, NadineSchneider, Vianello, R., Ric, Dalke, A., Cole,
298 B., AlexanderSavelyev, Swain, M., Turk, S., N. D., Vaucher, A., Kawashima, E., Wójcikowski, M., Probst,
299 D., Godin, G., Cosgrove, D., Pahl, A., JP, Francois Berenger, Strets123, JLVarjo, O'Boyle, N., Fuller, P.,
300 Jensen, J. H., Sfoma, G. and DoliathGavid [2020], 'rdkit/rdkit: 2020_03_1 (q1 2020) release'.
301 **URL:** <https://zenodo.org/record/3732262>
- 302 [12] *Landscape of targeted anti-cancer drug synergies in melanoma identifies a novel BRAF-VEGFR/PDGFR*
303 *combination treatment, author=Friedman, Adam A and Amzallag, Arnaud and Pruteanu-Malinici, Iulian*
304 *and Baniya, Subash and Cooper, Zachary A and Piris, Adriano and Hargreaves, Leeza and Igras, Vivien*
305 *and Frederick, Dennie T and Lawrence, Donald P and others* [2015], *PloS one* **10**(10), e0140310.
- 306 [13] Liaw, A. and Wiener, M. [2002], 'Classification and Regression by randomForest', *R News* **2**(3), 18–22.
307 **URL:** <https://cran.r-project.org/package=randomForest>
- 308 [14] Licciardello, M. P., Ringler, A., Markt, P., Klepsch, F., Lardeau, C.-H., Sdelci, S., Schirghuber, E., Müller,
309 A. C., Caldera, M., Wagner, A. et al. [2017], 'A combinatorial screen of the CLOUD uncovers a synergy
310 targeting the androgen receptor', *Nature chemical biology* **13**(7), 771–778.
- 311 [15] Lobo, J. M., Jiménez-Valverde, A. and Real, R. [2008], 'AUC: a misleading measure of the performance
312 of predictive distribution models', *Global ecology and Biogeography* **17**(2), 145–151.
- 313 [16] Martin, Y. C., Kofron, J. L. and Traphagen, L. M. [2002], 'Do structurally similar molecules have similar
314 biological activity?', *Journal of medicinal chemistry* **45**(19), 4350–4358.
- 315 [17] Meyer, C. T., Wooten, D. J., Lopez, C. F. and Quaranta, V. [2020], 'Charting the fragmented landscape of
316 drug synergy', *Trends in pharmacological sciences* **41**(4), 266–280.
- 317 [18] Nair, N. U., Greninger, P., Zhang, X., Friedman, A. A., Amzallag, A., Cortez, E., Sahu, A. D., Lee, J. S.,
318 Dastur, A., Egan, R. K. et al. [2023], 'A landscape of response to drug combinations in non-small cell lung
319 cancer', *Nature Communications* **14**(1), 3830.
- 320 [19] Narayan, R. S., Molenaar, P., Teng, J., Cornelissen, F. M., Roelofs, I., Menezes, R., Dik, R., Lagerweij, T.,
321 Broersma, Y., Petersen, N. et al. [2020], 'A cancer drug atlas enables synergistic targeting of independent
322 drug vulnerabilities', *Nature communications* **11**(1), 2935.

- 323 [20] Niepel, M., Hafner, M., Mills, C. E., Subramanian, K., Williams, E. H., Chung, M., Gaudio, B., Barrette,
324 A. M., Stern, A. D., Hu, B. et al. [2019], 'A multi-center study on the reproducibility of drug-response
325 assays in mammalian cell lines', *Cell systems* **9**(1), 35–48.
- 326 [21] O'Neil, J., Benita, Y., Feldman, I., Chenard, M., Roberts, B., Liu, Y., Li, J., Kral, A., Lejnine, S., Loboda, A.
327 et al. [2016], 'An unbiased oncology compound screen to identify novel combination strategies', *Molecular*
328 *cancer therapeutics* **15**(6), 1155–1162.
- 329 [22] Pemovska, T., Bigenzahn, J. W. and Superti-Furga, G. [2018], 'Recent advances in combinatorial drug
330 screening and synergy scoring', *Current opinion in pharmacology* **42**, 102–110.
- 331 [23] Preuer, K., Lewis, R. P., Hochreiter, S., Bender, A., Bulusu, K. C. and Klambauer, G. [2018], 'DeepSynergy:
332 predicting anti-cancer drug synergy with Deep Learning', *Bioinformatics* **34**(9), 1538–1546.
- 333 [24] R Core Team [2023], *R: A Language and Environment for Statistical Computing*, R Foundation for
334 Statistical Computing, Vienna, Austria.
335 **URL:** <https://www.R-project.org/>
- 336 [25] Rani, P., Dutta, K. and Kumar, V. [2023], 'Performance evaluation of drug synergy datasets using computa-
337 tional intelligence approaches', *Multimedia Tools and Applications* pp. 1–27.
- 338 [26] Roberts, M., Driggs, D., Thorpe, M., Gilbey, J., Yeung, M., Ursprung, S., Aviles-Rivero, A. I., Etmann,
339 C., McCague, C., Beer, L. et al. [2021], 'Common pitfalls and recommendations for using machine
340 learning to detect and prognosticate for COVID-19 using chest radiographs and CT scans', *Nature Machine*
341 *Intelligence* **3**(3), 199–217.
- 342 [27] Scantlebury, J., Vost, L., Carbery, A., Hadfield, T. E., Turnbull, O. M., Brown, N., Chenthamarakshan, V.,
343 Das, P., Grosjean, H., von Delft, F. et al. [2022], 'A Step Towards Generalisability: Training a Machine
344 Learning Scoring Function for Structure-Based Virtual Screening', *bioRxiv* pp. 2022–10.
- 345 [28] Shim, M., Lee, S.-H. and Hwang, H.-J. [2021], 'Inflated prediction accuracy of neuropsychiatric biomarkers
346 caused by data leakage in feature selection', *Scientific Reports* **11**(1), 7980.
- 347 [29] Tang, J., Wennerberg, K. and Aittokallio, T. [2015], 'What is synergy? The Saariselkä agreement revisited',
348 *Frontiers in pharmacology* **6**, 181.
- 349 [30] *The relationship between Precision-Recall and ROC curves, author=Davis, Jesse and Goodrich, Mark*
350 [2006], in 'Proceedings of the 23rd International Conference on Machine Learning', pp. 233–240.
- 351 [31] Wang, T., Wang, R. and Wei, L. [2023], 'AttenSyn: An Attention-Based Deep Graph Neural Network for
352 Anticancer Synergistic Drug Combination Prediction', *Journal of Chemical Information and Modeling* .
- 353 [32] Worthington, R. J. and Melander, C. [2013], 'Combination approaches to combat multidrug-resistant
354 bacteria', *Trends in biotechnology* **31**(3), 177–184.
- 355 [33] Wu, L., Wen, Y., Leng, D., Zhang, Q., Dai, C., Wang, Z., Liu, Z., Yan, B., Zhang, Y., Wang, J. et al.
356 [2022], 'Machine learning methods, databases and tools for drug combination prediction', *Briefings in*
357 *bioinformatics* **23**(1), bbab355.
- 358 [34] Yilancioglu, K., Weinstein, Z. B., Meydan, C., Akhmetov, A., Toprak, I., Durmaz, A., Iossifov, I., Kazan,
359 H., Roth, F. P. and Cokol, M. [2014], 'Target-independent prediction of drug synergies using only drug
360 lipophilicity', *Journal of chemical information and modeling* **54**(8), 2286–2293.
- 361 [35] Zheng, S., Aldahdooh, J., Shadbahr, T., Wang, Y., Aldahdooh, D., Bao, J., Wang, W. and Tang, J. [2021],
362 'DrugComb update: a more comprehensive drug sensitivity data repository and analysis portal', *Nucleic*
363 *acids research* **49**(W1), W174–W184.

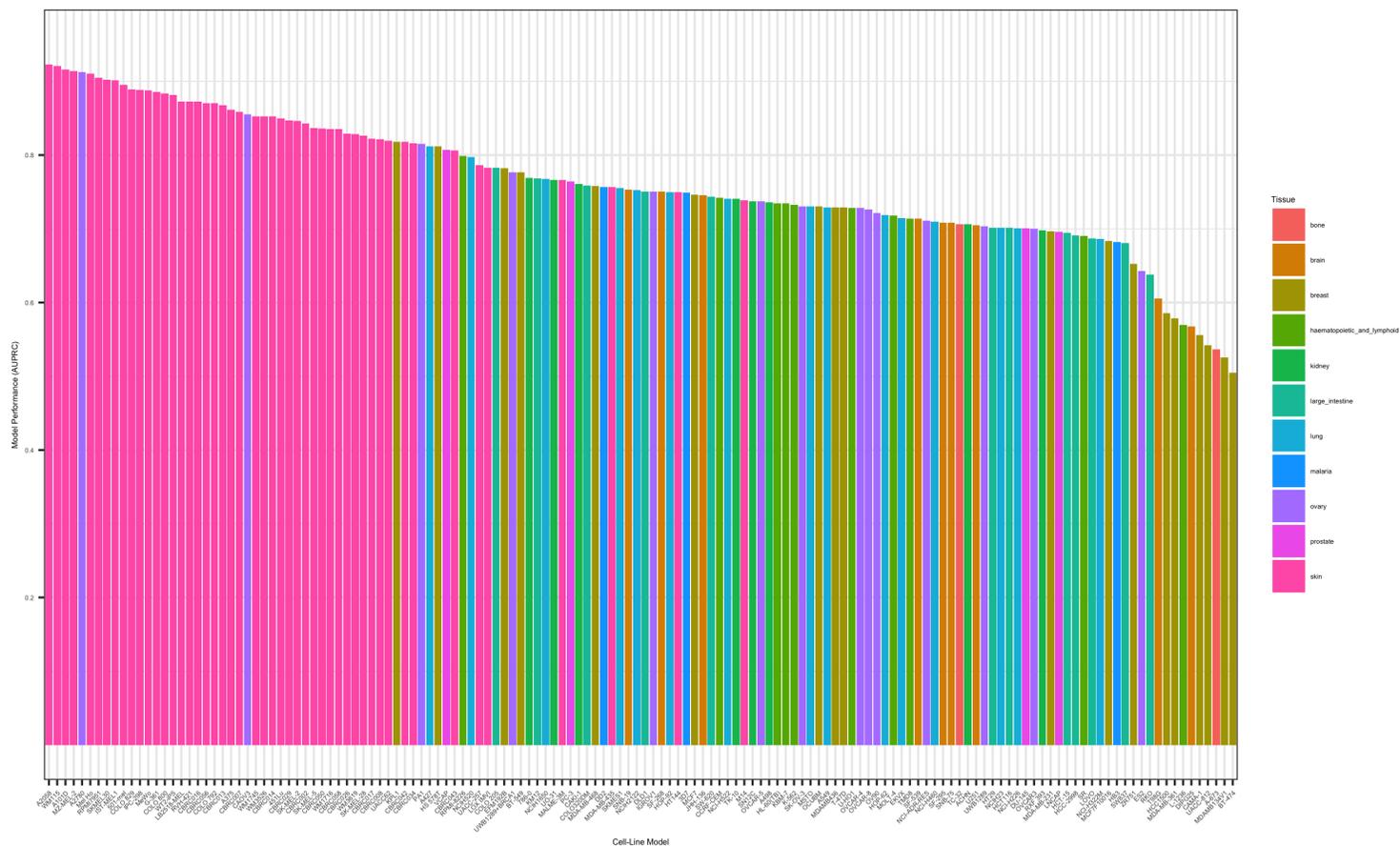


Figure 6: AUPRC performances for all cell-line models investigated in this study.

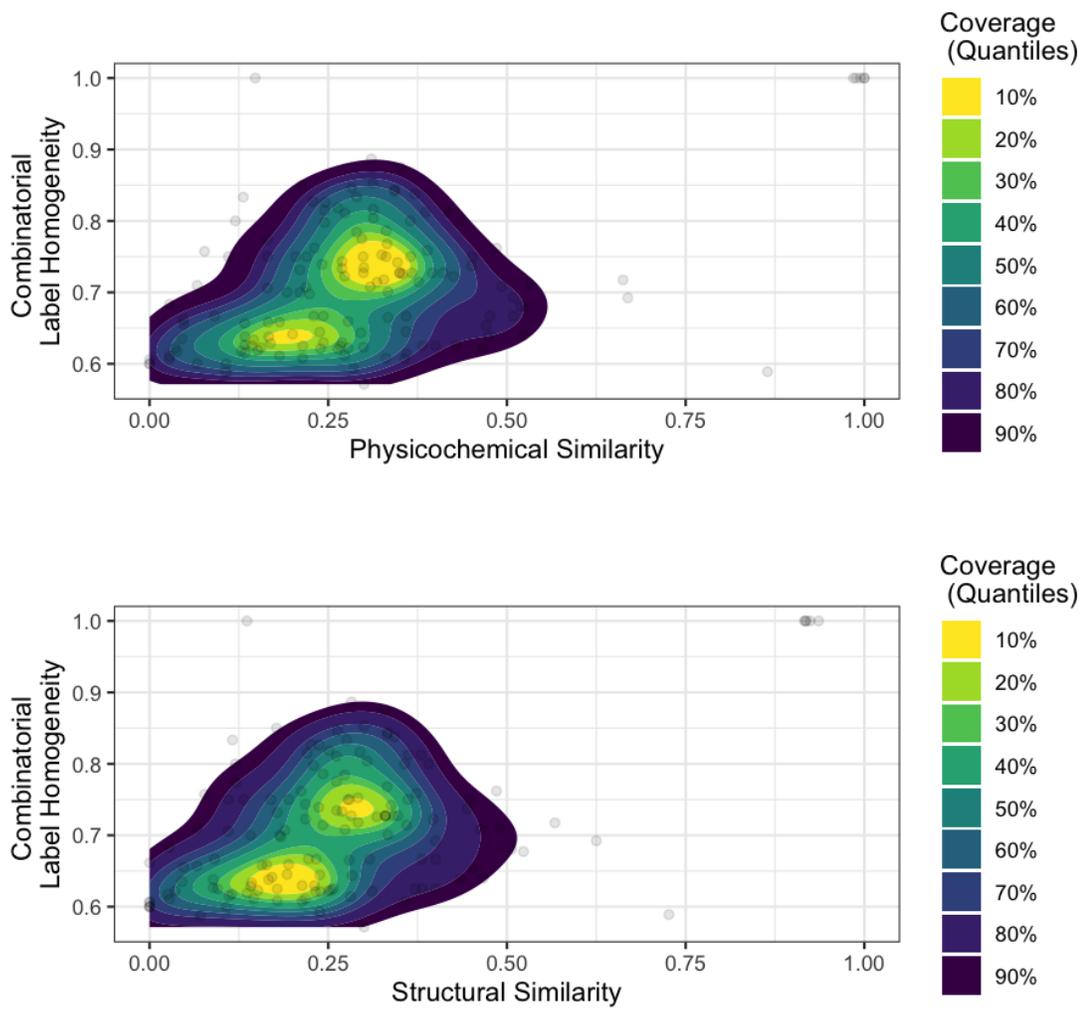


Figure 7: Drug similarity in physicochemical (upper) and structural (lower) spaces correlate with combinatorial label homogeneity.