
Forecasting labels under distribution-shift for machine-guided sequence design

Lauren Berk Wheelock^{*1}, Stephen Malina^{*}, Jeffrey Gerold^{*}, and Sam Sinai^{*2}

^{*}Dyno Therapeutics, 343 Arsenal Street, Suite 101, Watertown, MA, US

¹lauren.wheelock@dynotx.com

²sam.sinai@dynotx.com

Abstract

1 The ability to design and optimize biological sequences with specific functionalities
2 would unlock enormous value in technology and healthcare. In recent years,
3 machine learning-guided sequence design has progressed this goal significantly,
4 though validating designed sequences in the lab or clinic takes many months and
5 substantial labor. It is therefore valuable to assess the likelihood that a designed
6 set contains sequences of the desired quality (which often lies outside the label
7 distribution in our training data) *before* committing resources to an experiment.
8 Forecasting, a prominent concept in many domains where feedback can be delayed
9 (e.g. elections), has not been used or studied in the context of sequence design. Here
10 we propose a method to guide decision-making that forecasts the performance of
11 high-throughput libraries (e.g. 10^5 unique variants) based on estimates provided by
12 models, providing a posterior for the distribution of labels in the library. We show
13 that our method outperforms baselines that naively use model scores to estimate
14 library performance, which are the only tool available today for this purpose.

15 1 Introduction

16 Biological sequence design has long been of interest to practitioners in many domains, from agricul-
17 ture to therapeutics. For decades, sequences were designed through two means (i) Labor-intensive
18 rational design where expert human knowledge would generate a handful of candidate sequences
19 [1], (ii) High-throughput directed evolution approaches that utilize biological evolution to optimize
20 sequences towards a desired property [2]. Recently, the ability to synthesize DNA in high-throughput,
21 together with the wide adoption of high-capacity machine learning models, has opened a new
22 path that can combine the benefits of rational design (high quality), and directed evolution (high
23 throughput) [3, 4, 5, 6, 7]. In this setting, libraries containing up to 10^5 sequences are designed
24 using machine learning algorithms. Machine learning methods are used to score, optimize, and
25 filter sequences before committing to experiments [8, 9, 5, 10, 11]. In recent years, increasingly
26 nuanced perspectives on how to improve our trust in the output of machine learning models and
27 paired optimization procedures have evolved [12, 13, 10, 14, 15, 16, 17, 18]. Using these methods,
28 sequences targeting different objectives can be synthesized (e.g. transcription factor binding or other
29 regulatory sequences [19, 20]) in a library that can be measured in the desired context. However,
30 especially with products or traits of high complexity (e.g. *in-vivo* studies of proteins [21]), the overall
31 cost required to validate designs can be prohibitive. Therefore, even with model evaluations and
32 calibration of uncertainty around samples, there remains a gap in our ability to *forecast* the probability
33 of success: be it reaching a certain maximum performance, or finding a certain number of variants
34 above a minimum desired performance. This is distinct from attempting to predict the performance
35 of single sequences, in that it focuses on predicting the right-tail distribution of the performance our
36 entire library. In many settings, we would want to know whether the experiment has a high chance of

37 finding a (generally rare) high-performing sequence overall. Forecasts can help us decide whether to
 38 commit to a certain design and can save large costs. Forecasts can also inform other decisions such
 39 as deciding whether to repeat the design procedure for a library, deciding among libraries designed
 40 for different targets, or estimating the final price of developing a drug.

41 Forecasting is ubiquitous in domains with delayed feedback such as elections [22, 23]. The related
 42 topic of label shift [24, 25] classically relies on the “anticausal” assumption that the distribution
 43 of inputs given labels is constant across training and test sets - an assumption that is invalid in
 44 the case of design. More generally, domain adaptation has been studied in biological sequence
 45 design [26] but does not directly address forecasting and calibrating distributions under covariate
 46 and label shift. Recent work in conformal prediction directly tackles the problem of the kind of
 47 covariate shift that arises in sequence design settings [16], but its use requires known probability
 48 distributions over training and test sets, and nonzero prior probabilities on the entire support, meaning
 49 it cannot be applied to most libraries that were not designed with this approach in mind. To our
 50 knowledge, there are currently no methods that are suitable for forecasting library performance in the
 51 sequence design setting. This setting presents an interesting and somewhat unique challenge. For
 52 every designed sequence we can obtain scores for the expected performance, possibly from multiple
 53 models. However we are often aiming to make sequences that have a significantly higher score than
 54 anything observed in our training data, i.e. distribution shift is by design. Our challenge is to find the
 55 right balance between trusting our models’ predictions out-of-distribution and betting that our new
 56 designs would provide us with better-than-observed sequences.

57 2 Forecasting method overview

58 We start with labeled training data $(\mathbf{S}^0, \mathbf{Y}^0)$, where $\mathbf{S}^0 = \{s_i^0\}$ is a set of biological sequences and
 59 $\mathbf{Y}^0 = \{y_i^0\}$ is a set of continuous-valued labels, generally a fitness measurement in the sequence
 60 design setting such as packaging or transduction efficiency rates. Our goal is to forecast a distribution
 61 of labels \mathbf{Y}^1 for an unlabeled set \mathbf{S}^1 . That is, we are not concerned with the accuracy of each pair
 62 (s_i^1, y_i^1) , but only the overall distribution of \mathbf{Y}^1 , and in particular, in the right tail of \mathbf{Y}^1 , which
 63 indicates the maximum quality achieved by the set of sequences. To create our forecast, we have at
 64 our disposal a set of J regression models trained on $(\mathbf{S}^0, \mathbf{Y}^0)$, which produce test set predictions m_{ij}
 65 for sequence i by model j .

66 Naively, we could form ensembled point estimates $\hat{y}_i = \sum_j m_{ij}/J$ for each sequence and predict the
 67 distribution of \mathbf{Y}^1 to be the distribution of \hat{y}_i . There are three main disadvantages to this approach,
 68 which inspire different aspects of our forecasting method. We address them briefly below, and give a
 69 more thorough treatment with a complete algorithm in Section A.1.

70 **1) An implicit unimodal Gaussian assumption** Empirically, model ensembles tend towards unimodal
 71 Gaussian score distributions which do not empirically fit experimental data from designed sequences
 72 well. At several points in the experimental pipeline variants may “drop out,” failing to produce enough
 73 signal to reliably approximate a label (for example, due to failure of a protein to fold). This results in
 74 a multimodal distribution at both the population level, and implicitly, at the level of each sequence’s
 75 posterior. Thus, we seek to model each sequence as a bimodal Gaussian Mixture Model (GMM), and
 76 learn the parameters for each sequence’s posterior from its model scores. Moreover, while we use a
 77 GMM, our method could in theory be applied using a range of more complex distributions, with the
 78 only constraint being our ability to sample from them.

79 **2) Distribution (covariate and label) shift** Typically, the sequence set \mathbf{S}^1 is designed with model-
 80 guided exploration strategies informed by $(\mathbf{S}^0, \mathbf{Y}^0)$, with the objective of producing sequences that
 81 outperform the best sequences in \mathbf{S}^0 . This results in both significant distribution (covariate) shift,
 82 because the sequences \mathbf{S}^1 are reaching into untested areas of sequence space, and label shift, since we
 83 anticipate that \mathbf{Y}^1 will dominate \mathbf{Y}^0 , both on average and among each sets top-performers. (While
 84 there has been some important recent work on prediction in this design setting [16], this work assumes
 85 a shift in distribution within a consistent domain between \mathbf{S}^0 and \mathbf{S}^1 and no label shift.) To address
 86 distribution and label shift, we start by applying non-parametric non-linear calibration techniques
 87 to produce a “conservative” forecast that still allows for some label shift due to model uncertainty.
 88 We then consider scenarios with some trust placed in raw model scores to allow for some amount of
 89 extrapolation to regions further from our training set.

90 **3) Point estimates to posteriors** The point estimates that arise from model ensembles do not
91 provide a posterior for \mathbf{Y}^1 (nor do the model score variances, directly), and consequently these
92 tend to underestimate the frequency of events that are rare at the sequence level, but common at the
93 population level, such as the occurrence of high-valued sequences in the library. In our method, we
94 simulate draws of the entire library from the sequence-level posteriors to produce both expected
95 distributions as well as the frequency of rare events, which we interpret as posterior probabilities.

96 **3 Experimental results**

97 We validate our method by conducting experiments on four datasets - a set of simulated RNA binding
98 landscapes that allow us to access ground truth values for every sequence (and repeat multiple
99 experiments) as well as three experimentally-measured assays of protein fitness landscapes. These are
100 a viral protein packaging landscape, an experimentally measured IgG-Fc binding dataset for Protein
101 G's GB1 domain, and an experimentally measured GFP fluorescence dataset. Of the experimentally
102 measured landscapes, the viral protein assay is precisely conducted in the manner that forecast is
103 intended to be used. The GFP and GB1 landscapes are not conducted in this way, but we still expect
104 (and observed) forecast to outperform model point estimates.

105 **3.1 Experimental design**

106 For each of these experiments, we generated training and test sets \mathbf{S}^0 , \mathbf{Y}^0 and \mathbf{S}^1 , \mathbf{Y}^1 , and models
107 \mathbf{M} . We used the training set, models, and unlabeled test data to generate forecasts, and evaluated
108 them against realized distributions of \mathbf{Y}^1 using two key tools: a 2-sample Kolmogorov-Smirnov
109 statistic measuring distribution fit of top percentiles, and confidence interval coverage for counts of
110 points measured above a fixed threshold value. Both evaluate fit in the right tail of the distribution,
111 which is both the most challenging region to predict, and the one that is most critical to sequence
112 design applications. We present results of forecasts on all four datasets in Figure 1, with one dataset
113 per row of figures (RNA, AAV, GB1, and GFP respectively). While these experiments demonstrate
114 the efficacy of the forecasting procedure in a variety of experimental settings, we acknowledge that
115 additional experiments in follow-on studies could help clarify the method's strengths and weaknesses.

116 **3.2 Discussion of results**

117 In the left column of plots in Figure 1 we report on the Kolmogorov-Smirnov fit between the true
118 distribution of scores \mathbf{Y}^1 and either the forecast or the set of point estimates μ_i from model ensemble
119 scores. In figure 1a we plot the mean and 95% confidence interval across landscapes and starts
120 (see Sec A.7.1 for details), while the other experiments have a single landscape each. Since we are
121 primarily concerned with fit in the right tail, we limited each distribution to values above a percentile
122 threshold, and varied that threshold between 50% and 95%. Across this range, the forecasting method
123 improved distributional fit compared to ensemble point estimates, and while point estimates typically
124 decayed towards 1.0 (the theoretical worst-case upper bound of the statistic), the forecast consistently
125 maintained some predictive power even in the top 5th percentile. A sharp covariate shift in the AAV
126 capsid design sets \mathbf{S}^0 and \mathbf{S}^1 account for the problem difficulty in figures 1e and 1f, though even in
127 this problem our method directionally improves upon model estimates.

128 In the right column of figures we focus more closely on the right tail of the label distribution of \mathbf{Y}^1 ,
129 reporting the forecast's confidence interval for the number of sequence we can expect to find above a
130 threshold value compared to the estimate from model ensembles and the true counts. Since Figure
131 1b encompasses several landscapes and seeds, we set one threshold per landscape/seed at the 99th
132 percentile, while in the remaining experiments with a single landscape we evaluate accuracy over
133 a range of thresholds. Here we see that the forecast gives confidence intervals that include the true
134 count of sequences above a high threshold some of the time, and always improves upon the ensemble
135 estimates, which for most datasets severely underpredicts the prevalence of top-performers.

136 A key item of interest for sequence design is the performance of the top variants. We report ensemble,
137 forecast, and true values for the 99th percentile, mean of the top percentile, and maximum value for
138 each experiment in Figures 2 and 3 in the Appendix. These results echo the conclusions from Figure
139 1, showing highly accurate predictions on the AAV and GFP landscapes, and directionally correct
140 adjustments on GB1 and RNA landscapes.

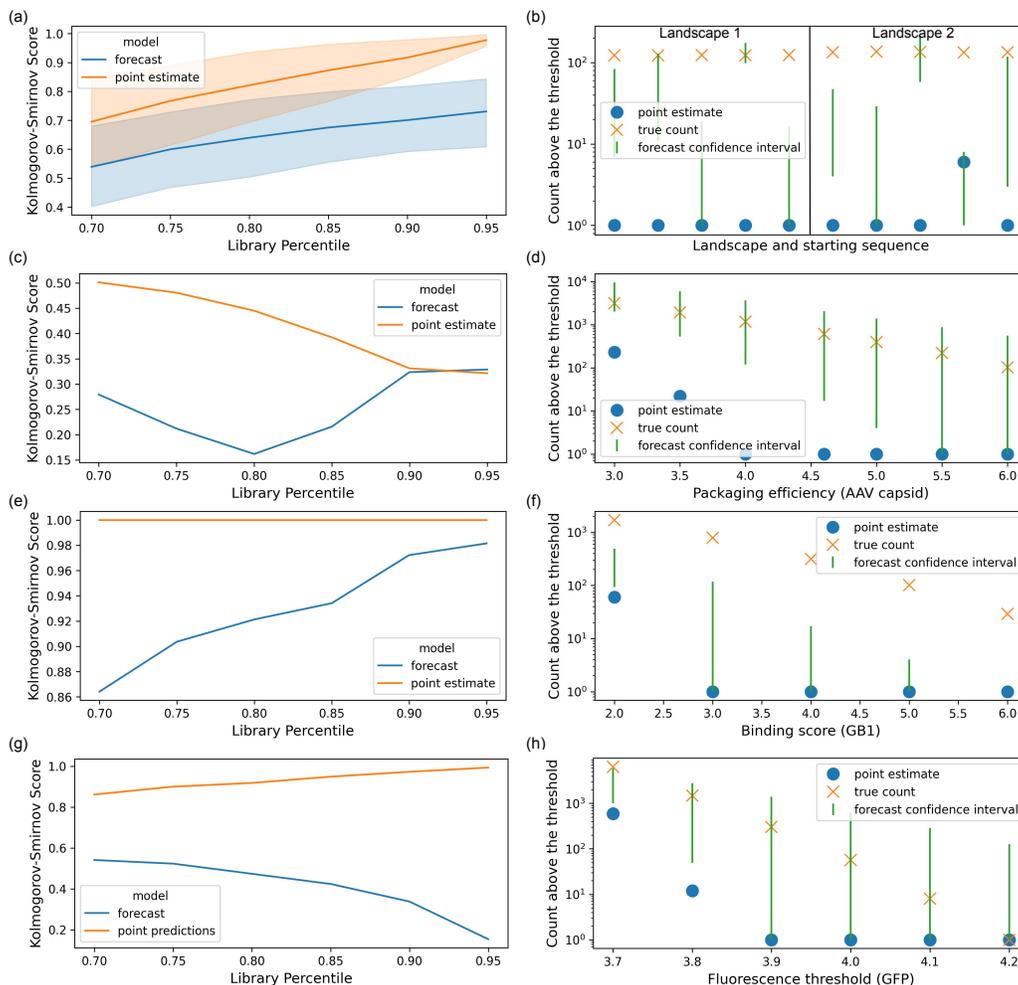


Figure 1: Comparison between forecast and point predictions in describing the right tail statistics of designed libraries. a) Distance from true (right-tail) distribution as measured by KS two-sample score for ensemble point estimate and our forecast on two RNA landscapes and five distinct designs per landscape (10 total forecasts) b) Top centile confidence interval coverage for RNA landscape 1 (14 nt) and RNA landscape 2 (50 nt) c,e,g) Distance from true distribution fit for ensemble point estimate and forecast d,f,h) Confidence interval coverage based on the number of samples above a certain measured performance c,d) For the AAV capsid design problem e,f) For the GB1 binding landscape g,h) For the GFP fluorescence landscape.

141 **4 Conclusion**

142 In this paper we argued for the relevance and impact of forecasting in the sequence design setting.
 143 We developed a novel approach for forecasting label distributions under covariate and label shift
 144 that occurs during model-guided design. Our approach can be used on any machine-guided library
 145 design for which we have regression models. We applied these methods in simulated and real-world
 146 sequence design settings and showed near-universal improvement (and never worse than the naive
 147 approach) in our ability to predict the shape of the right tail and counts of top performers. This work
 148 enables valuable estimation of the quality of designed libraries of biological sequences long before
 149 experimental results can be produced, which provides essential feedback to the designer. We hope
 150 that by defining this problem framework, and showcasing an approach to address it, we inspire further
 151 development for improving distributional forecasting in the model-guided sequence design setting.

References

- 152
- 153 [1] HW Hellinga. Rational protein design: combining theory and experiment. *Proceedings of the*
154 *National Academy of Sciences*, 94(19):10015–10017, 1997.
- 155 [2] Frances H Arnold. Design by directed evolution. *Accounts of chemical research*, 31(3):125–131,
156 1998.
- 157 [3] Bruce J Wittmann, Kadina E Johnston, Zachary Wu, and Frances H Arnold. Advances in
158 machine learning for directed evolution. *Current opinion in structural biology*, 69:11–18, 2021.
- 159 [4] Kevin K Yang, Zachary Wu, and Frances H Arnold. Machine-learning-guided directed evolution
160 for protein engineering. *Nature methods*, 16(8):687–694, 2019.
- 161 [5] Drew H Bryant, Ali Bashir, Sam Sinai, Nina K Jain, Pierce J Ogden, Patrick F Riley, George M
162 Church, Lucy J Colwell, and Eric D Kelsic. Deep diversification of an AAV capsid protein by
163 machine learning. *Nature Biotechnology*, 39(6):691–696, 2021.
- 164 [6] Ethan C Alley, Grigory Khimulya, Surojit Biswas, Mohammed AlQuraishi, and George M
165 Church. Unified rational protein engineering with sequence-based deep representation learning.
166 *Nature methods*, 16(12):1315–1322, 2019.
- 167 [7] Chloe Hsu, Hunter Nisonoff, Clara Fannjiang, and Jennifer Listgarten. Combining evolutionary
168 and assay-labelled data for protein fitness prediction. *bioRxiv*, 2021.
- 169 [8] Ratul Chowdhury, Nazim Bouatta, Surojit Biswas, Christina Floristean, Anant Kharkare,
170 Koushik Roye, Charlotte Rochereau, Gustaf Ahdrizt, Joanna Zhang, George M Church, et al.
171 Single-sequence protein structure prediction using a language model and deep learning. *Nature*
172 *Biotechnology*, pages 1–7, 2022.
- 173 [9] Jung-Eun Shin, Adam J Riesselman, Aaron W Kollasch, Conor McMahon, Elana Simon, Chris
174 Sander, Aashish Manglik, Andrew C Kruse, and Debora S Marks. Protein design and variant
175 prediction using autoregressive generative models. *Nature communications*, 12(1):1–11, 2021.
- 176 [10] David H Brookes, Hahnbeom Park, and Jennifer Listgarten. Conditioning by adaptive sampling
177 for robust design. *arXiv preprint arXiv:1901.10060*, 2019.
- 178 [11] Sam Sinai, Nina Jain, George M Church, and Eric D Kelsic. Generative aav capsid diversification
179 by latent interpolation. *bioRxiv*, 2021.
- 180 [12] Richard Fox. Directed molecular evolution by machine learning and the influence of nonlinear
181 interactions. *Journal of theoretical biology*, 234(2):187–199, 2005.
- 182 [13] Philip A Romero, Andreas Krause, and Frances H Arnold. Navigating the protein fitness land-
183 scape with gaussian processes. *Proceedings of the National Academy of Sciences*, 110(3):E193–
184 E201, 2013.
- 185 [14] Aviral Kumar and Sergey Levine. Model inversion networks for model-based optimization.
186 *Advances in Neural Information Processing Systems*, 33:5126–5137, 2020.
- 187 [15] Christof Angermueller, David Belanger, Andreea Gane, Zelda Mariet, David Dohan, Kevin
188 Murphy, Lucy Colwell, and D Sculley. Population-based black-box optimization for biological
189 sequence design. In *International Conference on Machine Learning*, pages 324–334. PMLR,
190 2020.
- 191 [16] Clara Fannjiang, Stephen Bates, Anastasios Angelopoulos, Jennifer Listgarten, and Michael I
192 Jordan. Conformal prediction for the design problem. *arXiv preprint arXiv:2202.03613*, 2022.
- 193 [17] Erik Nijkamp, Jeffrey Ruffolo, Eli N Weinstein, Nikhil Naik, and Ali Madani. Progen2:
194 exploring the boundaries of protein language models. *arXiv preprint arXiv:2206.13517*, 2022.
- 195 [18] Leo Feng, Padideh Nouri, Aneri Muni, Yoshua Bengio, and Pierre-Luc Bacon. Designing
196 biological sequences via meta-reinforcement learning and bayesian optimization. *arXiv preprint*
197 *arXiv:2209.06259*, 2022.

- 198 [19] Luis A Barrera, Anastasia Vedenko, Jesse V Kurland, Julia M Rogers, Stephen S Gisselbrecht,
199 Elizabeth J Rossin, Jaie Woodard, Luca Mariani, Kian Hong Kock, Sachi Inukai, et al. Survey
200 of variation in human transcription factors reveals prevalent dna binding changes. *Science*,
201 351(6280):1450–1454, 2016.
- 202 [20] Eeshit Dhaval Vaishnav, Carl G de Boer, Jennifer Molinet, Moran Yassour, Lin Fan, Xian
203 Adiconis, Dawn A Thompson, Joshua Z Levin, Francisco A Cubillos, and Aviv Regev. The
204 evolution, evolvability and engineering of gene regulatory dna. *Nature*, 603(7901):455–463,
205 2022.
- 206 [21] Pierce J Ogden, Eric D Kelsic, Sam Sinai, and George M Church. Comprehensive AAV
207 capsid fitness landscape reveals a viral gene and enables machine-guided design. *Science*,
208 366(6469):1139–1143, 2019.
- 209 [22] Patrick Hummel and David Rothschild. Fundamental models for forecasting elections. *Re-*
210 *searchDMR.com/HummelRothschild_FundamentalModel*, 2013.
- 211 [23] Wei Wang, David Rothschild, Sharad Goel, and Andrew Gelman. Forecasting elections with
212 non-representative polls. *International Journal of Forecasting*, 31(3):980–991, 2015.
- 213 [24] Zachary Lipton, Yu-Xiang Wang, and Alexander Smola. Detecting and correcting for label shift
214 with black box predictors. In *International conference on machine learning*, pages 3122–3130.
215 PMLR, 2018.
- 216 [25] Marco Saerens, Patrice Latinne, and Christine Decaestecker. Adjusting the outputs of a classifier
217 to new a priori probabilities: a simple procedure. *Neural computation*, 14(1):21–41, 2002.
- 218 [26] Shireen Abestesh, Ilke Akartuna, Alexander Brown, Megan Cramer, Farhan Damani, Jeff
219 Gerold, Jorma Gorns, Jeff Jones, Helene Kuchwara, Jamie Kwasnieski, et al. Efficient design
220 of optimized AAV capsids using multi-property machine learning models trained across cells,
221 organs and species. In *MOLECULAR THERAPY*, volume 29, pages 12–12. CELL PRESS,
222 2021.
- 223 [27] Nobuyuki Otsu. A threshold selection method from gray-level histograms. *IEEE transactions*
224 *on systems, man, and cybernetics*, 9(1):62–66, 1979.
- 225 [28] Jerome Friedman and Robert Tibshirani. The monotone smoothing of scatterplots. *Technomet-*
226 *rics*, 26(3):243–250, 1984.
- 227 [29] Steven Henikoff and Jorja G Henikoff. Performance evaluation of amino acid substitution
228 matrices. *Proteins: Structure, Function, and Bioinformatics*, 17(1):49–61, 1993.
- 229 [30] Sam Sinai, Richard Wang, Alexander Whatley, Stewart Slocum, Elina Locane, and Eric D
230 Kelsic. Adalead: A simple and robust adaptive greedy search algorithm for sequence design.
231 *arXiv preprint arXiv:2010.02141*, 2020.
- 232 [31] Nikolaus Hansen and Andreas Ostermeier. Completely derandomized self-adaptation in evolu-
233 tion strategies. *Evolutionary computation*, 9(2):159–195, 2001.
- 234 [32] Colin A Smith and Tanja Kortemme. Predicting the tolerated sequences for proteins and protein
235 interfaces using rosettabackrub flexible backbone design. *PloS one*, 6(7):e20451, 2011.
- 236 [33] Sebastian Kmiecik and Andrzej Kolinski. Folding pathway of the b1 domain of protein g
237 explored by multiscale modeling. *Biophysical journal*, 94(3):726–736, 2008.
- 238 [34] Nicholas C Wu, Lei Dai, C Anders Olson, James O Lloyd-Smith, and Ren Sun. Adaptation in
239 protein fitness landscapes is facilitated by indirect paths. *Elife*, 5:e16965, 2016.
- 240 [35] Karen S Sarkisyan, Dmitry A Bolotin, Margarita V Meer, Dinara R Usmanova, Alexander S
241 Mishin, George V Sharonov, Dmitry N Ivankov, Nina G Bozhanova, Mikhail S Baranov,
242 Onuralp Soylemez, et al. Local fitness landscape of the green fluorescent protein. *Nature*,
243 533(7603):397–401, 2016.

244 A Supplemental Material

245 A.1 Detailed description of the forecasting method

246 To generate forecasts, we first transform model predictions for each sequence into parameters for
 247 their posterior distributions. We then draw from those sequence-level posterior distributions to form
 248 simulations of the library, generating a library-level posterior. This library-level posterior reflects
 249 our epistemic uncertainty about the ground-truth performance of each sequence given our model
 250 predictions and the aleatoric uncertainty of our measurements given this ground-truth. That is, our
 251 predictions are more like prediction intervals than confidence intervals, and we do not generate a
 252 posterior for ground-truth values. Our objective is to infer sequence-level posterior parameters from
 253 model predictions in a way that is both well-calibrated to the training data and allows for test set
 254 performance of the sequences to differ from or exceed training set performance due to distribution
 255 shift.

256 A.2 Fitting sequence label posteriors

257 While this forecasting framework can be applied to learn posteriors for many distributions families,
 258 informed by empirical library label distributions from historical experiments, we believe the natural
 259 distribution for sequence labels in our data is a Gaussian mixture model (GMM) with two modes:
 260 one for “functional” sequences and one for “non-functional” or “broken” sequences. (To arrive at
 261 this conclusion, we considered a number of alternative distribution families, including varying the
 262 skew and kurtosis of each mode of the GMM, but did not find sufficient improvements to justify the
 263 increased model complexity.)

264 To model a sequence s_i using a GMM, we assume there is a probability of functionality p_i , a mean
 265 and variance in the functional mode $(\mu_i^+, (\sigma_i^+)^2)$, and a mean and variance in the non-functional
 266 mode $(\mu_i^-, (\sigma_i^-)^2)$ that parameterize normal distributions \mathcal{N} so that

$$Y_i \sim p_i \mathcal{N}(\mu_i^+, \sigma_i^+) + (1 - p_i) \mathcal{N}(\mu_i^-, \sigma_i^-). \quad (1)$$

267 In contrast, our predictive models only provide point estimates m_{ij} for each sequence. We assume
 268 that the true, multimodal distribution of each sequence can be summarized with two degrees of
 269 freedom (a mean μ_i and standard deviation σ_i) and that these two parameters independently generate
 270 model scores, mixture model parameters, and measurement values. Explicitly, we assume the model
 271 predictions $m_{ij} \sim \mathcal{N}(\mu_i, \sigma_i)$, so that, given a set of model predictions, we infer μ_i and σ_i to be the
 272 models’ sample mean and variance ($\mu_i = 1/J \sum_j m_{ij}$ and $\sigma_i^2 = 1/J \sum_j (m_{ij} - \mu_i)^2$). We further
 273 assume there are independent relationships between μ_i and the set (p_i, μ_i^+, μ_i^-) , and between σ_i and
 274 the pair σ_i^+ and σ_i^- .

275 Since the GMM parameters $(p_i, \mu_i^+, \mu_i^-, \sigma_i^+, \sigma_i^-)$ are unique to each sequence, we cannot infer them
 276 in the usual manner using \mathbf{S}^0 as a training set. Instead, we need to further model and learn the
 277 relationship between the pair μ_i, σ_i and the GMM parameters. Specifically, we start by identifying
 278 the value y_{mid} that we use to separate the two modes of \mathbf{Y} . This value can either be set manually,
 279 using expert knowledge, or automatically by analyzing the distribution \mathbf{Y}^0 . We found that Otsu’s
 280 method [27], which finds the separating point that minimizes intra-class variance, provided robust
 281 values of y_{mid} on our data. We can then divide our set \mathbf{S}^0 into two halves across the boundary:
 282 $\mathbf{S}^{0+} = \{s_i \mid y_i^0 \geq y_{mid}\}$ and $\mathbf{S}^{0-} = \{s_i \mid y_i^0 < y_{mid}\}$. This provides us with separate training sets
 283 for the functional parameters μ_i^+, σ_i^+ and the non-functional parameters μ_i^-, σ_i^- .

284 A.3 About isotonic regression

285 We will run isotonic regression to find the best monotonic piece-wise linear fit to this data. Explicitly,
 286 isotonic regression operates on a dataset of pairs of scalars $\{x_i, y_i\}$ and produces a non-parametric
 287 model represented by data-prediction pairs (x_i, \hat{y}_i) that seek to minimize the least squares error
 288 $\sum_i (\hat{y}_i - y_i)^2$ subject to a monotonicity constraint $y_i \leq y_j \forall i, j$ s.t. $x_i < x_j$. This results in a
 289 quadratic program, though it is easily solvable exactly by sorting y_i and iteratively averaging pairs of
 290 “violators” of monotonicity, making training efficient and deterministic [28].

291 As a point of notation, we will use the abbreviation IR_y to refer to an isotonic regression model
 292 trained to predict y defined by the pairs (x_i, \hat{y}_i) , and $IR_y(x'_i)$ to be the prediction of this model given

293 input x'_j . To compute $IR_y(x'_j)$ on a new data point, first we check to see if $x'_j = x_i$ for some i and
 294 if so we predict the corresponding $IR_y(x'_j) = \hat{y}_i$. Otherwise we sort x_i and find a consecutive pair
 295 such that $x_i < x'_j < x_{i+1}$ and predict $IR_y(x'_j)$ by linearly interpolating between \hat{y}_i and \hat{y}_{i+1} . If x'_j
 296 is less than all x_i , or more than all x_i , $IR_y(x'_j)$ is set to the min and max values of \hat{y}_i respectively.
 297 We note here that this necessarily means that isotonic regression will never produce labels outside the
 298 range of the training labels - we will address this in a few ways in the coming sections.

299 A.4 Inferring parameters with isotonic regression

300 We assume a non-linear monotonic relationship between the model ensemble mean for a sequence
 301 $\mu_i = 1/J \sum_j \mu_{ij}$ and the probability that the sequence i will be functional (p_i). To infer p_i , we
 302 train an isotonic regression model IR_p that, given μ_i , aims to predict the indicator I_i^+ which is 1
 303 if $s_i > y_{mid}$ and 0 otherwise. Effectively, given a new input μ_i , this model returns the fraction of
 304 sequences that are functional out of the training samples with similar mean values, and interprets this
 305 rate as the probability that the sequence i will be functional.

306 Inferring the mean parameters μ_i^+ and μ_i^- is more straight-forward: we build isotonic regression
 307 models IR_{μ^+} and IR_{μ^-} to predict y_i from μ_i , but restrict the training set to \mathbf{S}^{0+} and \mathbf{S}^{0-} respectively.
 308 This gives us calibration to the conditional distributions for being functional and non-functional
 309 respectively.

310 To infer the variances σ_i^{2+} and σ_i^{2-} , we first form squared residuals of labels given model ensemble
 311 means $res_i = (y_i - \mu_i)^2$ and build isotonic models IR_{σ^+} and IR_{σ^-} relating the model variance
 312 to these residuals res_i . As with μ_i^+ and μ_i^- , we compute σ_i^{2+} and σ_i^{2-} by training models on the
 313 disjoint training sets \mathbf{S}^{0+} and \mathbf{S}^{0-} respectively. The complete algorithm for inferring the GMM
 314 parameters is described in Algorithm 1.

315 **Applying the forecast to non-ensembles** While the presentation of our method assumes access to
 316 an ensemble of models, we note that thus far the only information we have used from the ensembles
 317 is the ensemble mean and variance (μ_i, σ_i^2) for each feature. Therefore, as an alternative, any single
 318 model that itself outputs an expected value and uncertainty (which includes many neural networks)
 319 can stand alone in providing the input (μ_i, σ_i^2) to forecasting calibration. The only technique that does
 320 not generalize from ensembles to models-with-uncertainty is the “optimistic model de-ensembling”
 321 technique discussed in the Section A.6.

322 A.5 Simulating the posterior distribution

323 Given the parameters generated by Algorithm 1, we can draw samples \hat{y}_i^1 for each sequence, and
 324 aggregate them into draws for the entire distribution $\hat{\mathbf{Y}}^1$. We can then treat the set of simulated
 325 values of $\hat{\mathbf{Y}}^1$ as a posterior distribution and query this distribution to determine the frequency of
 326 distribution-level events. By computing metrics on $\hat{\mathbf{Y}}^1$ and considering their distributions across
 327 simulations, we can arrive at empirical confidence intervals for metrics such as the count of sequences
 328 that perform above some threshold value, as we see in Figures 1b,d,f,h.

329 A.6 Tuning the forecast from conservative to optimistic

330 We can further refine this basic algorithm using additional techniques that allow us to diversify our
 331 approach over degrees of trust in our training set.

332 **Semi-calibrated regression** Our main calibration tool, isotonic regression, aggressively limits
 333 predicted labels to be within the range of training values. To allow for some distribution shift, we can
 334 gradually transition from calibrated predictions towards the center of the distribution of \mathbf{S}^1 towards
 335 uncalibrated, out-of-distribution values towards the limits of the distribution, in a technique we call
 336 “semi-calibration.”

337 Let $P_Y(y)$ be the percentile of the value y from among the empirical distribution Y . That is, $P(y)$
 338 is the fraction of $y \in Y$ with $y_i < y$. In our case, we consider the distribution of model ensemble
 339 means on our training set \mathbf{S}^0 , that is, the set $M = \{1/J \sum_j m_{ij} \mid s_i \in \mathbf{S}^0\}$. Then given a new
 340 sequence s_i we can compute its model ensemble prediction μ_i as well as its functional isotonicly

341 calibrated mean μ_i^+ , and evaluate where its model ensemble falls relative to the training distribution
 342 by computing the percentile $P_M(\mu_i)$. Finally, for any temperature-like coefficient $0 < q \leq 1$, we
 343 define our semi-calibrated mean $\tilde{\mu}_i(q)$ to be

$$\tilde{\mu}_i(q) = (qP_M(\mu_i))(\mu_i) + (1 - qP_M(\mu_i))(IR_{\mu_i^+}(\mu_i)). \quad (2)$$

344 Thus, lower values will be completely calibrated to the training set, while higher values will be a mix
 345 of calibrated and uncalibrated values. Note that we only produce this correction for functional mean
 346 values μ_i^+ , as we expect non-functional values to be fully in the training set distribution. This leads
 347 to an update to the model from Equation 1:

$$\hat{y}_i \sim p_i \mathcal{N}(\tilde{\mu}_i(q), \sigma_i^+) + (1 - p_i) \mathcal{N}(\mu_i^-, \sigma_i^-). \quad (3)$$

348 **Correcting for covariate shift** In addition to model score distribution shift, we also see covariate
 349 shift that creates model score bias. In our context, we consider edit distance to wild type the primary
 350 such covariate, though the method easily generalizes to more complex distance metrics (such as
 351 BLOSUM [29]), as well as other quantitative side-information. To correct for this shift, we form
 352 *signed* residuals from the training set between the *calibrated* μ_i^+ values and the true values y_i (i.e.
 353 $(\mu_i^+ - y_i)$). (If we are also applying the semi-calibration technique from the last paragraph, we use
 354 $(\tilde{\mu}_i - y_i)$ instead.) We can regress those residuals on edit distance (ED_i) using either isotonic or
 355 linear regression, and apply this correction back to the mean prediction μ_i^+ . We can also apply this
 356 approach to adjust the probability parameter p_i , encoding the understanding that sequences are less
 357 likely to be functional at higher distances from the wild-type. That is, we compute:

$$\begin{aligned} res_i &= \mu_i^+ - y_i \\ IR^{ED} &= \text{IR model trained on } \{(ED_i, res_i)\} \\ \mu_i^{ED} &= \mu_i^+ - IR^{ED}(ED_i) \\ \hat{y}_i &\sim p_i \mathcal{N}(\mu_i^{ED}, \sigma_i^+) + (1 - p_i) \mathcal{N}(\mu_i^-, \sigma_i^-). \end{aligned} \quad (4)$$

358 **Optimistic model de-ensembling** So far, we have assumed model scores m_{ij} are drawn from a
 359 Gaussian distributions parameterized by μ_i, σ_i . Alternatively, we could assume that each model j
 360 represents a distinct distribution, and that μ_i are drawn from these distributions with equal probability.
 361 Using this approach, in each simulation we first randomly select one model independently for each
 362 sequence and use that model’s prediction as the sequence’s expected value: $\mu_i m_{ij}$. This can result in
 363 more optimistic forecasts when scores have high inter-model variance.

364 **Hedging against calibration assumptions** Together, these calibration techniques create a menu of
 365 options that allow us to build forecasts that range from conservative to optimistic given input data.
 366 Given a set of calibration strategies, we can simulate instances of \mathbf{Y}^1 , and by aggregating simulations
 367 across frameworks, we can form a posterior for the distribution of \mathbf{Y}^1 that captures our uncertainties
 368 at the sequence level, the model level, and the overall forecasting approach level.

369 A.7 Descriptions of experimental data

370 A.7.1 Simulated RNA landscapes

371 Our first set of experiments investigates the performance of forecasting approach using FLEXS
 372 [30], a simulation environment for sequence design which gives access to ground-truth and model-
 373 approximated fitness landscapes. We study design problems on two RNA landscapes with a hidden
 374 binding target of size 14 and 50 nucleotides. Each training set \mathbf{S}^0 was constructed by mutating a
 375 sequence from a starting seed (5 seeds for each landscape) with between 1 and 3 mutations per
 376 sequence on average. We trained four kinds of predictive models on one-hot encoded sequences:
 377 linear regressions, convolutional neural networks, multi-layer perceptions, and random forest models.
 378 We used four exploration algorithms to design sequences using these models \mathbf{S}^1 : CMA-ES [31],
 379 CbAS [10], Adalead [30], and random sampling.

380 A.7.2 In-vitro AAV packaging assay

381 Bryant et al. [5] quantitatively assay the packaging efficiency of 200K viral capsid variants, modified
 382 in a 28 amino-acid region of the protein. The experiment was designed in two steps, where first a

383 smaller set of training examples were assayed and used to train *classification* models. Then these
384 models were used to design a set of variants that were optimized for the probability of packaging.
385 The paper uses classification models and greedy optimization to generate the second batch. The
386 existence of this first training set, and the distribution-shifted designed set is exactly the setting we
387 have devised the our forecasting method for, and where its performance can be most meaningfully
388 evaluated.

389 For this data set we retrain *regression* models on the first set of sequences (S^0) using five independently
390 seeded convolutional neural networks, and five recurrent neural networks. We used these models to
391 generate ensembled point estimates as well as forecast the distribution of packaging efficiencies on
392 the model-designed set of sequences (S^1).

393 **A.7.3 Protein G GB1 IgG-Fc binding domain**

394 A region of four amino acids in Protein G (GB1) is known to be critical for IgG-Fc binding and has
395 been used extensively as a tool for evaluating sequence landscape prediction and design tasks [32, 33].
396 Our data is sourced from experiments by Wu et al. [34]. We created S^0 by selecting sequences
397 with performance below that of the wild-type, combined with a small fraction of sequences between
398 wild-type and the median performance in the set, leaving all other sequences for S^1 . We trained five
399 independently seeded random forest models and five multi-layer perceptrons on S^0 (forgoing more
400 sophisticated models due to the short length of the variable sequence), and used these to generate
401 point estimates and forecast predictions for S^1 .

402 **A.7.4 Green Fluorescent Protein**

403 The green fluorescent protein of *Aequorea victoria* (avGFP) provides an additional fitness landscape
404 for studying sequence prediction and design. We used a dataset of 540,250 protein variants and their
405 associated fluorescence level [35]. We followed the same procedure as GB1 for splitting S^0 and S^1 ,
406 putting sequences with fluorescence below $3x$ log-fluorescence of the wild-type and a small portion
407 of sequences up to the median fluorescence above this threshold into S^0 , and the rest into S^1 . We
408 then trained the same models as in our AAV experiments - five independently seeded convolutional
409 neural networks, and five recurrent neural networks, and generated point estimates and forecasts for
410 S^1 using these models.

411 **A.8 Computational details**

412 **A.8.1 Compute**

413 All of our experiments were run using a single server with a single GPU running in GCP (Google
414 Cloud Platform). We used an Nvidia V100 for training models on the GFP landscape and an Nvidia
415 K80 for the other three experiments' model training.

416 **A.8.2 Hyperparameters**

417 Across all of our experiments, we used five model architectures: convolutional neural networks
418 (CNNs), recurrent neural networks (RNNs), multi-layer perceptrons (MLPs), linear models, and
419 random forests. Linear models and random forests were initialized with default parameters using
420 the sklearn library. CNNs used 32 filters, 64 filters, and 256 filters with 1, 2, and 2 convolutional
421 layers followed by 1, 2, and 2 hidden layers of width 32, 64, and 64 for the AAV, RNA, and GFP
422 experiments respectively. RNNs used embeddings of size 32 combined with 1 and 2 recurrent layers,
423 then followed by 1 hidden layer of size 56 and 128 for AAV and GFP respectively. MLPs used 1 and
424 3 hidden layers of width 50 and 32 for GB1 and RNA experiments respectively. All three model
425 architectures were trained used Adam with a learning rate of $1e-3$ across experiments.

426 **A.8.3 Licenses**

427 FLEXS is open source and Apache licensed. All other code was written for this project in python
428 using common packages that use BSD, PSFL, Apache, and MIT licenses.

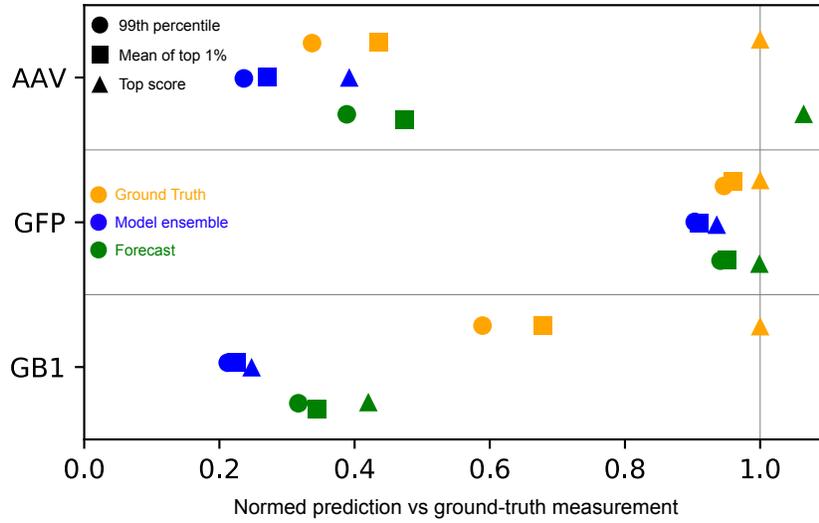


Figure 2: The ensemble, forecast, and measured values for the 99th percentile, mean of the top percentile, and maximum value for the AAV, GFP, and GB1 experiments, normalized to the maximum ground-truth measured value for each experiment.

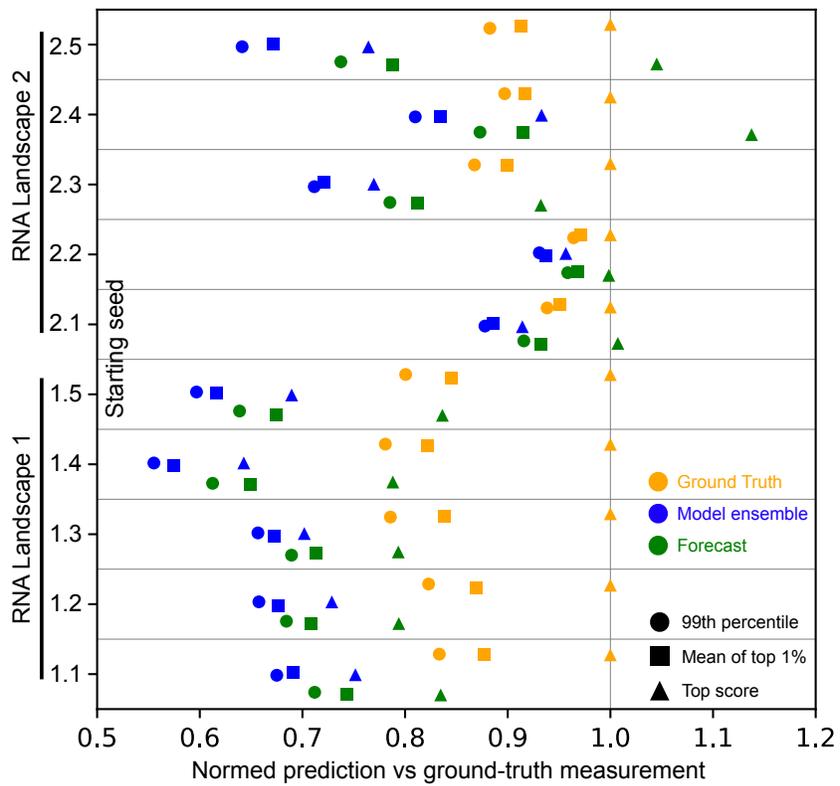


Figure 3: The ensemble, forecast, and true values for the 99th percentile, mean of the top percentile, and maximum value for the RNA experiments

Table 1: Relative experimental difficulty due to model score-based covariate and label shift, as measured by the K-S score between distributions of training and test ensemble means, and between measurement distributions from among top-scoring variants, respectively

Experiment	Model score covariate shift	Model score-based label shift
AAV	0.332	0.159
GB1	0.960	0.916
GFP	0.888	0.541
Landscape 1 Start 1	0.176	0.283
Landscape 1 Start 2	0.162	0.227
Landscape 1 Start 3	0.329	0.518
Landscape 1 Start 4	0.579	0.575
Landscape 1 Start 5	0.214	0.521
Landscape 2 Start 1	0.371	0.606
Landscape 2 Start 2	0.402	0.094
Landscape 2 Start 3	0.283	0.830
Landscape 2 Start 4	0.285	0.794
Landscape 2 Start 5	0.381	0.747

429 A.9 Plots of forecast accuracy on top performers

430 A.10 Quantifying distribution shift

431 As an attempt to quantify the difficulty of each forecasting problem, we computed metrics of
 432 covariate shift and label shift. Covariate shift measures the change in distribution in covariate space
 433 (our sequences and covariates associated with those sequences), while label shift measures a change
 434 in the conditional distribution of the outcome given those covariates. For this preliminary analysis,
 435 we restricted our study to using model ensemble scores as the main covariate of interest. It would
 436 also be reasonable to apply this to edit distances, or higher-dimensional covariates.

437 To measure model score covariate shift, we can apply a 2-sample Kolmogorov-Smirnov test to the
 438 entire distributions of model scores for \mathbf{S}^0 and \mathbf{S}^1 . This gives us a measure on a common scale from
 439 no shift (0) to completely disjoint supports (1).

440 Measuring model score-based label shift precisely is challenging in our setting, since our data
 441 regularly violates the common assumption for label shift research that the test set output support is a
 442 subset of the training set support, so we cannot calculate ratios between the density functions. Instead,
 443 we again use the 2-sample K-S test, this time comparing distributions of \mathbf{Y}^0 and \mathbf{Y}^1 but conditioned
 444 on high model scores (defined as the 90th percentile of the training set distribution and above).

445 We report these metrics in Table 1. We note that the AAV experiment, where the forecast performed
 446 especially well, had a lesser degree of covariate and label shift compared to other experiments.
 447 At the other extreme, the GB1 experiment had extreme covariate and label shift, and while the
 448 forecasting method improved upon the ensemble prediction directionally, the forecast produced very
 449 low confidence interval coverage for this experiment. This suggests a possible connection between
 450 shift scores and forecasting difficulty. On the other hand, we can look at the RNA experiments
 451 and consider one landscape at a time, which allows us to potentially isolate the relationship between
 452 these covariate shift metrics and forecasting performance. Here, however, there does not appear to
 453 be any clear relationship between either type of distribution shift and the accuracy of the forecast.
 454 Therefore, while the AAV and GB1 results suggest a possible connection, further experiments will
 455 be needed to validate these metrics as a useful tool for quantifying forecasting difficulty.

456 A.11 Detailed forecasting algorithm

457 See Algorithm 1 for a complete description of the forecasting algorithm described in Section A.1
 458 (excluding the extensions in Section A.5).

Algorithm 1 Inferring Gaussian mixture model parameters from a set of normally distributed model scores

Input: a training set ($\mathbf{S}^0, \mathbf{Y}^0$) and test set (\mathbf{S}^1) with model values m_{ij} for each $s_i \in \mathbf{S}^0 \cup \mathbf{S}^1$.

Returns: $(p_i, \mu_i^+, \mu_i^-, \sigma_i^+, \sigma_i^-)$ for each $i \in \mathbf{S}^1$

Learn cutoff value y_{mid} from \mathbf{Y}^0 using Otsu's method

for s_i **in** S^0 **do**

 Compute $\mu_i = \frac{\sum_j m_{ij}}{J}$ (model ensemble means)

 Compute $\sigma_i^2 = \frac{\sum_j (m_{ij} - \mu_i)^2}{J}$ (model ensemble variance)

 Compute $res_i^2 = (y_i - \mu_i)^2$ (squared residuals of model ensemble means)

 Compute $I_i^+ = 1$ if $Y_i^0 < y_{mid}$ and 0 otherwise

end for

Define $S^{0+} = \{i \mid Y_i^0 \geq y_{mid}\}$ (training subset for “functional” sequences)

Define $S^{0-} = \{i \mid Y_i^0 < y_{mid}\}$ (training subset for “broken” sequences)

Train isotonic model IR_p on pairs (μ_i, I_i^+) for $i \in S^0$

Train isotonic model IR_{μ^+} on (μ_i, y_i) for $i \in S^{0+}$

Train isotonic model IR_{μ^-} on (μ_i, y_i) for $i \in S^{0-}$

Train isotonic model IR_{σ^+} on (σ_i^2, res_i^2) for $i \in S^{0+}$

Train isotonic model IR_{σ^-} on (σ_i^2, res_i^2) for $i \in S^{0-}$

for s_i **in** S^1 **do**

 Compute $\mu_i = \frac{\sum_j m_{ij}}{J}$ (model ensemble means)

 Compute $\sigma_i^2 = \frac{\sum_j (m_{ij} - \mu_i)^2}{J}$ (model ensemble variance)

 Compute $(p_i, \mu_i^+, \mu_i^-, \sigma_i^+, \sigma_i^-) = (IR_p(\mu_i), IR_{\mu^+}(\mu_i), IR_{\mu^-}(\mu_i), IR_{\sigma^+}(\sigma_i), IR_{\sigma^-}(\sigma_i))$

end for

return $\{(p_i, \mu_i^+, \mu_i^-, \sigma_i^+, \sigma_i^-)\}$
