

1 **Supplementary Information for**

2 **Extrapolation Validation (EV): A Universal Validation Method**

3 **for Mitigating Machine Learning Extrapolation Risks**

4 Mengxian Yu¹, Yin-Ning Zhou², Qiang Wang¹, Fangyou Yan^{1*}

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22 **Text S1** Statistical parameters.

23 Statistical parameters are defined as follows:

24 Standard deviation (σ):
$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (y - \frac{1}{n} \sum_{i=1}^n y)^2} \quad (\text{S1})$$

25 Coefficient of determination (R^2):
$$R^2 = 1 - \frac{\sum_{i=1}^n (y_{i,\text{exp}} - y_{i,\text{cal}})^2}{\sum_{i=1}^n (y_{i,\text{exp}} - \bar{y}_{\text{exp}})^2} \quad (\text{S2})$$

26 Root mean squared error (RMSE):
$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_{i,\text{exp}} - y_{i,\text{cal}})^2} \quad (\text{S3})$$

27 Mean absolute error (MAE):
$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_{i,\text{exp}} - y_{i,\text{cal}}| \quad (\text{S4})$$

28 where, n is the numbers of samples.

29

30 **Text S2** Development of mathematical relationships models.

31 To establish ML models for multiple algorithms, the scikit-learn¹ version 1.3.0 package for
32 Python² 3.11 was adopted to develop the MLR, LASSO, Ridge, GPR, MLP, AdaBoost, RF, GBDT,
33 SVM, and KNN models, and the xgboost³ version 1.7.6 package was applied to develop the XGBoost
34 model. The independent variable data were standardized for modelling. The initial 33 ML models for
35 the three datasets were established based on the default hyperparameters of a total of 11 ML
36 algorithms within the scikit-learn¹ and xgboost³ packages.

37 The phenomenon of extrapolation-failure was initially identified by analyzing the performance of
38 the initial model training set, test (I) set, test (F) set, and test (B) set by R^2 . Subsequently, the
39 hyperparameters of 33 models were optimized (Figs. S1~S3) by exhaustive enumeration to observe
40 the influence of hyperparameters on the performance of different algorithms. For selecting the model
41 with the best extrapolation ability, the model with the smallest average value of MAE_{training} , $MAE_{\text{test(F)}}$
42 and $MAE_{\text{test(B)}}$ ($(MAE_{\text{training}} + MAE_{\text{test(F)}} + MAE_{\text{test(B)}})/3$) statistical parameters within the tuned
43 hyperparameters was taken as the optimal model (Tables S4, S6 and S8). To be clear, During the
44 development of machine learning models containing the parameter “random_state”, this
45 parameter is always set to 1. The R^2 of the optimal models for linear univariate, linear multivariate,
46 and nonlinear multivariate data relationships of 11 ML algorithms further confirm the hypothesis that
47 ML models involving tree algorithms may not have extrapolation ability.

48

49 **Table S1** Statistical results of the initial hyperparametric ML models for data with linear univariate

50 functional relationships.

Model	RMSE _{training}	RMSE _{test(F)}	RMSE _{test(I)}	RMSE _{test(B)}	R^2_{training}	$R^2_{\text{test(F)}}$	$R^2_{\text{test(I)}}$	$R^2_{\text{test(B)}}$
MLR	0.0000	0.0000	0.0000	0.0000	1.0000	1.0000	1.0000	1.0000
Lasso	1.0000	2.9627	0.6703	2.8991	1.0000	1.0000	1.0000	1.0000
Ridge	0.5754	1.7048	0.3857	1.6682	1.0000	1.0000	1.0000	1.0000
SVM	98.4925	500.1807	40.1207	488.1564	0.8805	0.8202	0.9738	0.8397
GPR	0.0001	256.9459	0.0001	6.2930	1.0000	0.3048	1.0000	0.9983
MLP	691.4459	1146.0182	670.2223	200.2035	0.3231	1.0000	0.6214	1.0000
AdaBoost	6.0552	249.8868	4.5307	239.7142	0.9988	0.0000	0.9986	0.0000
XGBoost	0.3888	232.0081	0.4183	220.4561	1.0000	0.0000	1.0000	0.0000
RF	0.3940	232.7778	0.3881	220.9704	1.0000	0.0000	1.0000	0.0000
KNN	0.3651	235.2828	0.0000	223.7365	1.0000	0.0000	1.0000	0.0000
GBDT	1.0646	232.1061	1.1254	220.5793	1.0000	0.0000	0.9999	0.0000

52 **Table S2** Statistical results of the initial hyperparametric ML models for data with linear multivariate

53 functional relationships.

Model	$RMSE_{\text{training}}$	$RMSE_{\text{test(F)}}$	$RMSE_{\text{test(I)}}$	$RMSE_{\text{test(B)}}$	R^2_{training}	$R^2_{\text{test(F)}}$	$R^2_{\text{test(I)}}$	$R^2_{\text{test(B)}}$
MLR	0.0000	0.0000	0.0000	0.0000	1.0000	1.0000	1.0000	1.0000
Lasso	1.0004	1.9373	0.6647	5.3968	1.0000	1.0000	1.0000	1.0000
Ridge	0.7173	7.8683	0.5561	158.4471	1.0000	1.0000	1.0000	0.0560
SVM	44.3835	297.7298	12.1532	294.7730	0.9212	0.9175	0.9899	0.5270
GPR	0.0001	465.6339	0.0000	102.1042	1.0000	0.9062	1.0000	0.6258
MLP	394.8063	576.1380	396.1819	336.3715	0.0783	0.9997	0.0910	0.6676
AdaBoost	3.8536	151.9961	2.9115	147.9956	0.9986	0.0000	0.9984	0.0000
XGBoost	0.2246	138.9467	0.2347	137.0528	1.0000	0.0000	1.0000	0.0000
RF	0.1703	139.4360	0.1737	137.3924	1.0000	0.0000	1.0000	0.0000
KNN	0.2211	140.9398	0.0001	139.0768	1.0000	0.0000	1.0000	0.0000
GBDT	0.6509	139.0176	0.6796	137.1404	1.0000	0.0000	0.9999	0.0000

55 **Table S3** Statistical results of the initial hyperparametric ML models for data with nonlinear

56 multivariate functional relationships.

Model	RMSE _{training}	RMSE _{test(F)}	RMSE _{test(I)}	RMSE _{test(B)}	R^2_{training}	$R^2_{\text{test(F)}}$	$R^2_{\text{test(I)}}$	$R^2_{\text{test(B)}}$
MLR	0.0000	0.0038	0.0000	2.4726	1.0000	1.0000	1.0000	0.9996
Lasso	1.0100	0.7533	0.7541	57.6987	1.0000	1.0000	1.0000	0.9246
Ridge	0.7746	7.3083	0.5963	126.6090	1.0000	1.0000	1.0000	0.0378
SVM	63.2524	362.8745	20.7764	381.9859	0.8992	0.9132	0.9830	0.5085
GPR	0.0001	589.5278	0.0000	143.2726	1.0000	0.9078	1.0000	0.6095
MLP	535.0079	778.6396	534.5355	340.9841	0.0708	0.9998	0.0839	0.6845
AdaBoost	4.9844	181.0609	4.1720	197.7776	0.9985	0.0000	0.9979	0.0000
XGBoost	0.2602	166.3100	0.2871	182.8309	1.0000	0.0000	1.0000	0.0000
RF	0.2070	166.9202	0.2073	183.3232	1.0000	0.0000	1.0000	0.0000
KNN	0.2753	168.7369	0.0004	185.4687	1.0000	0.0000	1.0000	0.0000
GBDT	0.8094	166.4499	0.8436	182.9998	1.0000	0.0000	0.9999	0.0000

58 **Table S4** Optimal hyperparameters of the optimal ML models for data with linear univariate
 59 functional relationships.

Model	Parameters
Lasso	max_iter = 100
Ridge	alpha = 0.1, max_iter = 100
SVM	C = 51, gamma = 0.001, kernel = 'linear'
GPR	kernel=2**2 * DotProduct(sigma_0=4)
MLP	activation = 'identity', hidden_layer_sizes = (210, 10), solver = 'lbfgs'
AdaBoost	learning_rate = 0.11, n_estimators = 200
XGBoost	learning_rate = 0.19, max_depth = 7, n_estimators = 600
RF	max_depth = 9, n_estimators = 10
KNN	n_neighbors = 1
GBDT	learning_rate = 0.05, max_depth = 5, n_estimator = 600

61 **Table S5** Statistical results of the optimal ML models for data with linear univariate functional

62 relationships.

Model	$RMSE_{\text{training}}$	$RMSE_{\text{test(F)}}$	$RMSE_{\text{test(I)}}$	$RMSE_{\text{test(B)}}$	R^2_{training}	$R^2_{\text{test(F)}}$	$R^2_{\text{test(I)}}$	$R^2_{\text{test(B)}}$
MLR	0.0000	0.0000	0.0000	0.0000	1.0000	1.0000	1.0000	1.0000
Lasso	0.1000	0.2963	0.0670	0.2899	1.0000	1.0000	1.0000	1.0000
Ridge	0.0577	0.1710	0.0387	0.1673	1.0000	1.0000	1.0000	1.0000
SVM	0.0579	0.1716	0.0388	0.1679	1.0000	1.0000	1.0000	1.0000
GPR	0.0000	0.0000	0.0000	0.0000	1.0000	1.0000	1.0000	1.0000
MLP	0.0000	0.0000	0.0000	0.0000	1.0000	1.0000	1.0000	1.0000
AdaBoost	10.1147	259.6758	7.4628	249.4922	0.9966	0.0000	0.9959	0.0000
XGBoost	0.0035	231.8071	0.0039	220.2600	1.0000	0.0000	1.0000	0.0000
RF	0.5990	232.6736	0.6088	220.4324	1.0000	0.0000	1.0000	0.0000
KNN	0.0000	231.8060	0.0000	220.2589	1.0000	0.0000	1.0000	0.0000
GBDT	0.0000	231.8060	0.0000	220.2589	1.0000	0.0000	1.0000	0.0000

64 **Table S6** Optimal hyperparameters of the optimal ML models for data with linear multivariate

65 functional relationships.

Model	Parameters
Lasso	alpha = 0.3, max_iter = 100, tol = 0.0002
Ridge	alpha = 1.9, max_iter = 100
SVM	C = 191, gamma = 0.001, kernel = 'linear'
GPR	kernel=2**2 * DotProduct(sigma_0=3)
MLP	activation = 'identity', alpha = 0.0005, hidden_layer_sizes = (110, 10), solver = 'lbfgs'
AdaBoost	learning_rate = 0.13, n_estimators = 200
XGBoost	learning_rate = 0.13, max_depth = 7, n_estimators = 900
RF	max_depth = 9, max_features = 4, n_estimators = 10
KNN	n_neighbors = 1
GBDT	learning_rate = 0.05, max_depth = 5, n_estimator = 500

67 **Table S7** Statistical results of the optimal ML models for data with linear multivariate functional

68 relationships.

Model	$RMSE_{\text{training}}$	$RMSE_{\text{test(F)}}$	$RMSE_{\text{test(I)}}$	$RMSE_{\text{test(B)}}$	R^2_{training}	$R^2_{\text{test(F)}}$	$R^2_{\text{test(I)}}$	$R^2_{\text{test(B)}}$
MLR	0.0000	0.0000	0.0000	0.0000	1.0000	1.0000	1.0000	1.0000
Lasso	0.3048	0.2207	0.2088	3.9420	1.0000	1.0000	1.0000	0.9978
Ridge	1.2107	10.8737	0.9283	113.7655	0.9999	1.0000	0.9999	0.0008
SVM	0.0700	0.9778	0.0758	37.3265	1.0000	1.0000	1.0000	0.7957
GPR	0.0000	0.0000	0.0000	0.0001	1.0000	1.0000	1.0000	1.0000
MLP	0.0006	0.0244	0.0005	5.2772	1.0000	1.0000	1.0000	0.9967
AdaBoost	5.7320	156.6871	4.0441	153.0914	0.9970	0.0000	0.9968	0.0000
XGBoost	0.0036	138.8537	0.0040	136.9603	1.0000	0.0000	1.0000	0.0000
RF	0.3001	139.3735	0.2986	137.0648	1.0000	0.0000	1.0000	0.0000
KNN	0.0000	138.8526	0.0000	136.9591	1.0000	0.0000	1.0000	0.0000
GBDT	0.0000	138.8526	0.0000	136.9591	1.0000	0.0000	1.0000	0.0000

70 **Table S8** Optimal hyperparameters of the optimal ML models for data with nonlinear multivariate

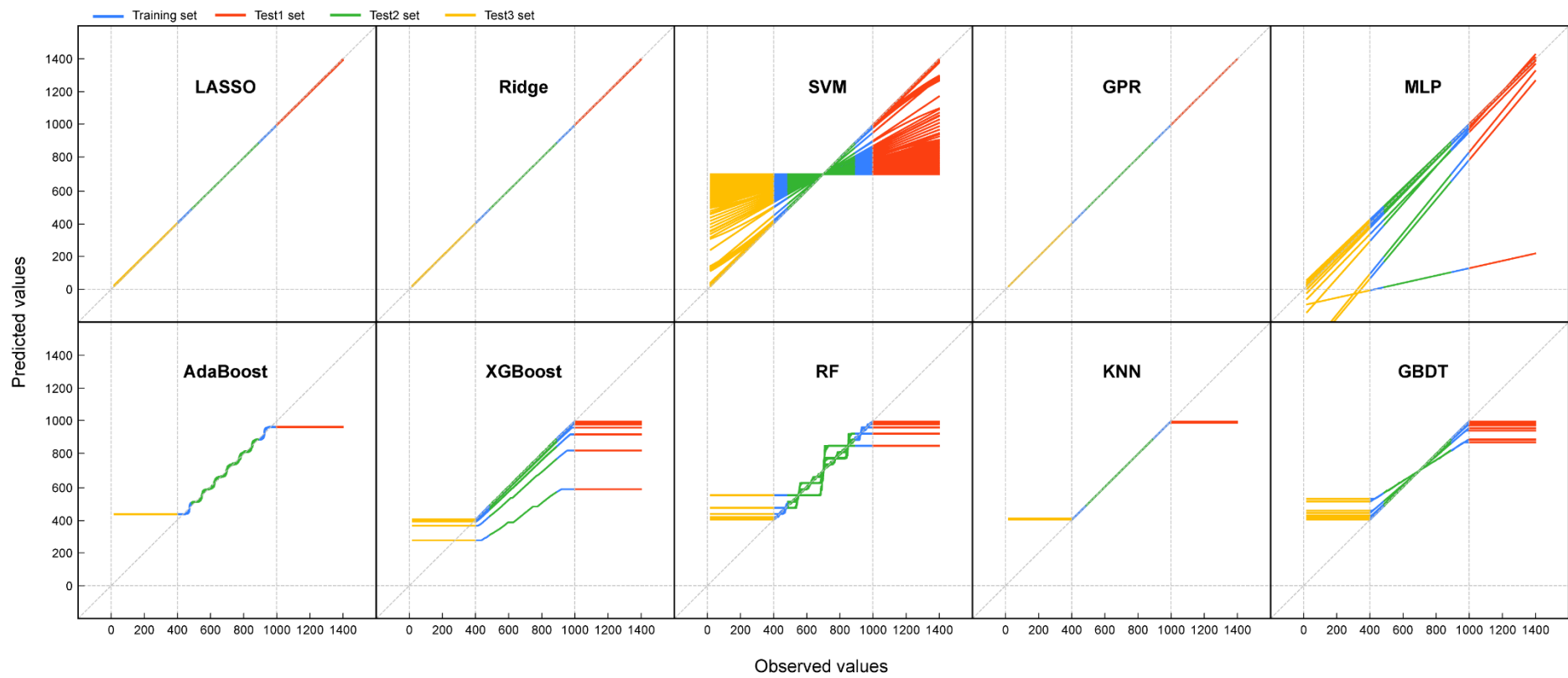
71 functional relationships.

Model	Parameters
Lasso	alpha = 1.7, max_iter = 600
Ridge	alpha = 1.9, max_iter = 100
SVM	C = 101, gamma = 0.011, kernel = 'sigmoid'
GPR	kernel=1**2 * DotProduct(sigma_0=2)
MLP	activation = 'identity', alpha = 0.0007, hidden_layer_sizes = (210, 310), solver = 'lbfgs'
AdaBoost	learning_rate = 0.15, n_estimators = 200
XGBoost	learning_rate = 0.17, max_depth = 7, n_estimators = 700
RF	max_depth = 9, max_features = 5, n_estimators = 10
KNN	n_neighbors = 1
GBDT	learning_rate = 0.19, max_depth = 7, n_estimator = 200

73 **Table S9** Statistical results of the optimal ML models for data with nonlinear multivariate functional

74 relationships.

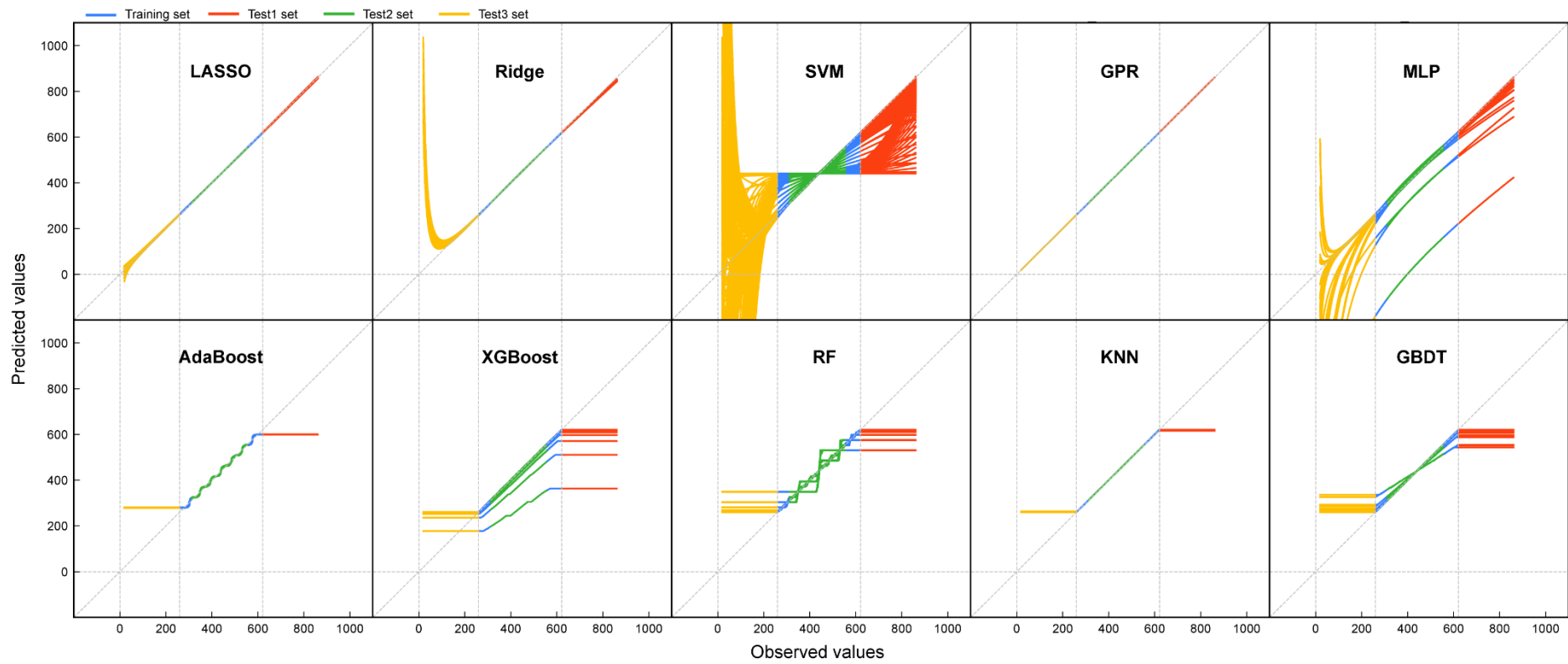
Model	RMSE _{training}	RMSE _{test(F)}	RMSE _{test(I)}	RMSE _{test(B)}	R^2_{training}	$R^2_{\text{test(F)}}$	$R^2_{\text{test(I)}}$	$R^2_{\text{test(B)}}$
MLR	0.0000	0.0038	0.0000	2.4726	1.0000	1.0000	1.0000	0.9996
Lasso	1.7006	3.0380	1.1622	4.7048	1.0000	1.0000	1.0000	0.9990
Ridge	1.3402	10.6818	1.0273	76.6308	0.9999	1.0000	1.0000	0.4168
SVM	0.5614	8.3121	0.2512	24.8271	1.0000	1.0000	1.0000	0.9576
GPR	0.0000	0.0038	0.0000	2.4724	1.0000	1.0000	1.0000	0.9996
MLP	0.0038	0.0901	0.0034	0.3379	1.0000	1.0000	1.0000	1.0000
AdaBoost	7.1512	187.5066	5.1517	204.2875	0.9970	0.0000	0.9966	0.0000
XGBoost	0.0033	166.2165	0.0035	182.7721	1.0000	0.0000	1.0000	0.0000
RF	0.3807	166.8446	0.3883	182.9056	1.0000	0.0000	1.0000	0.0000
KNN	0.0000	166.2155	0.0000	182.7710	1.0000	0.0000	1.0000	0.0000
GBDT	0.0000	166.2155	0.0000	182.7710	1.0000	0.0000	1.0000	0.0000



76

77

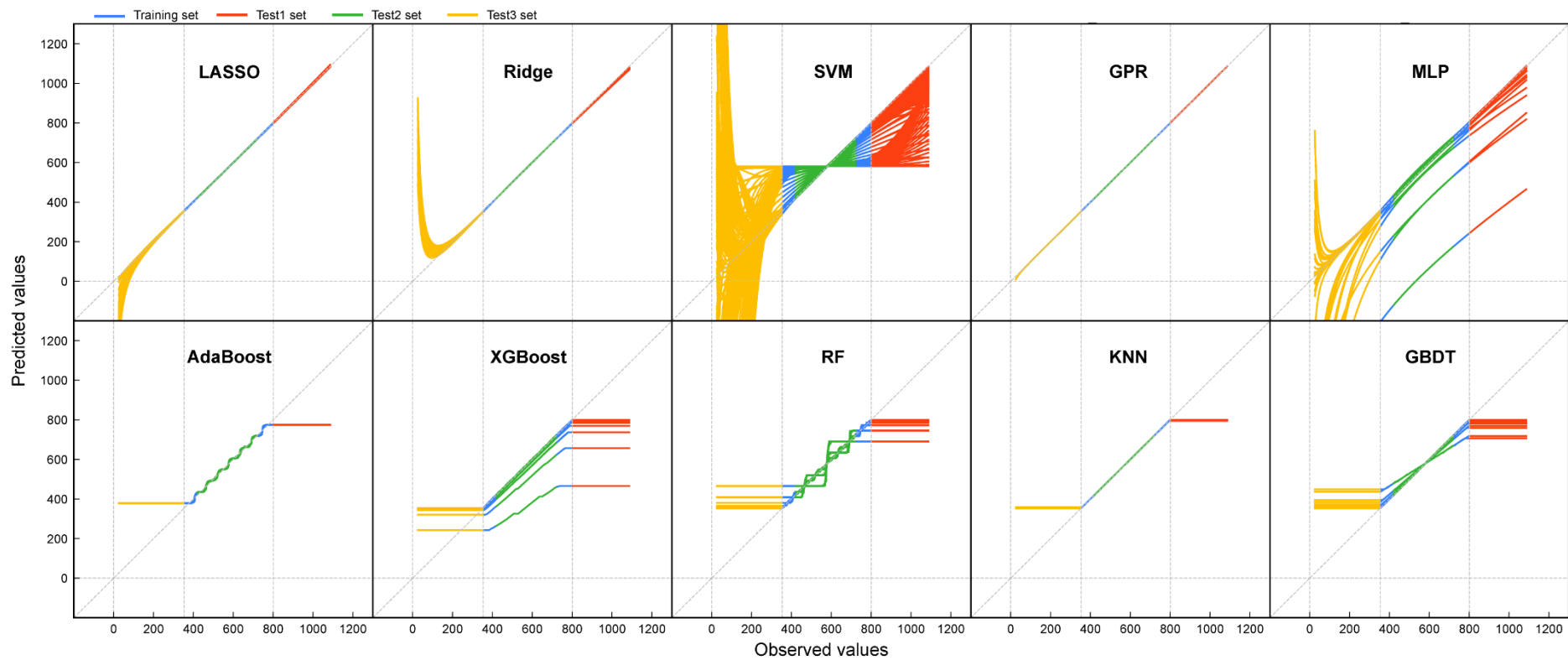
Fig. S1 Hyperparametric adjustment process for linear univariate ML models.



78

79

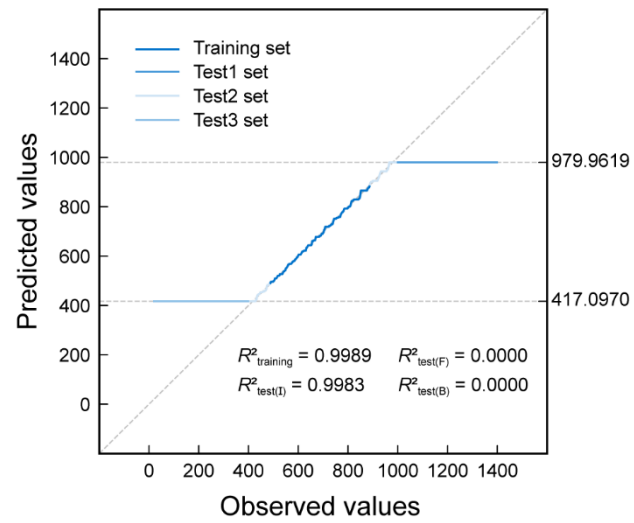
Fig. S2 Hyperparametric adjustment process for linear multivariate ML models.



80

81

Fig. S3 Hyperparametric adjustment process for nonlinear multivariate ML models.

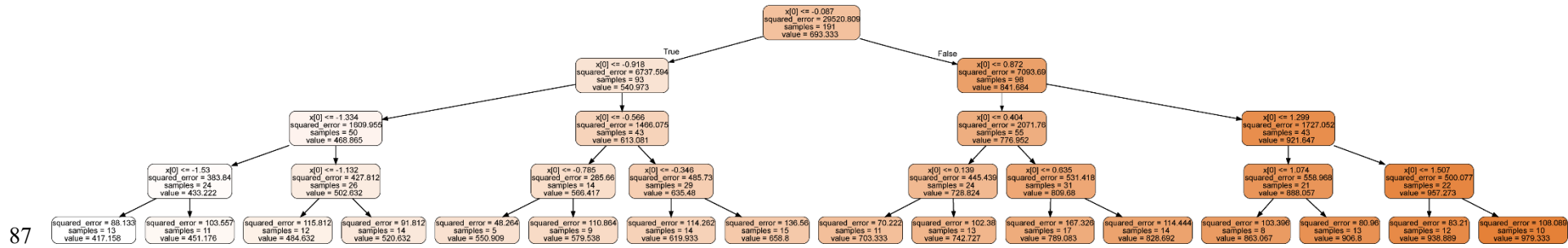
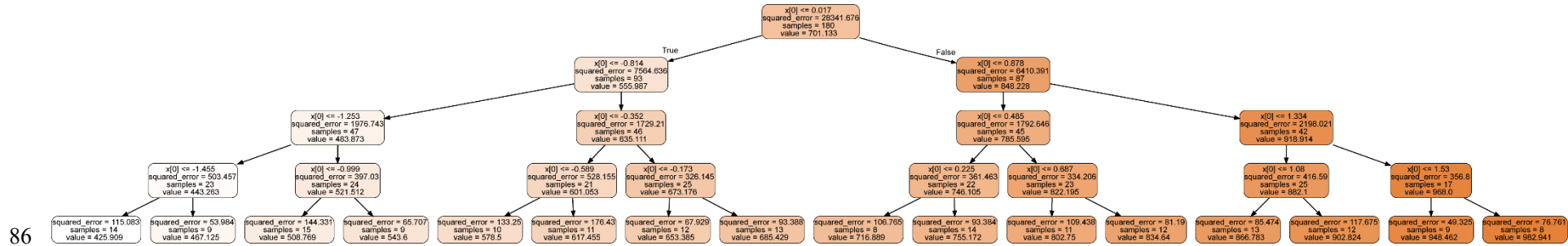
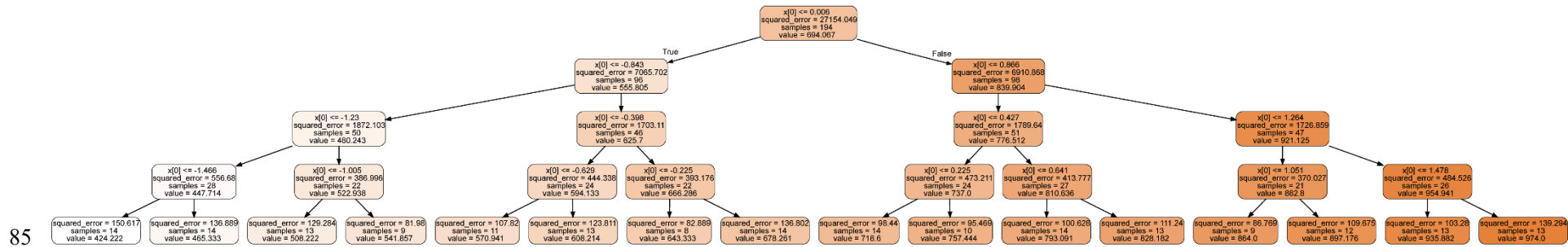


82

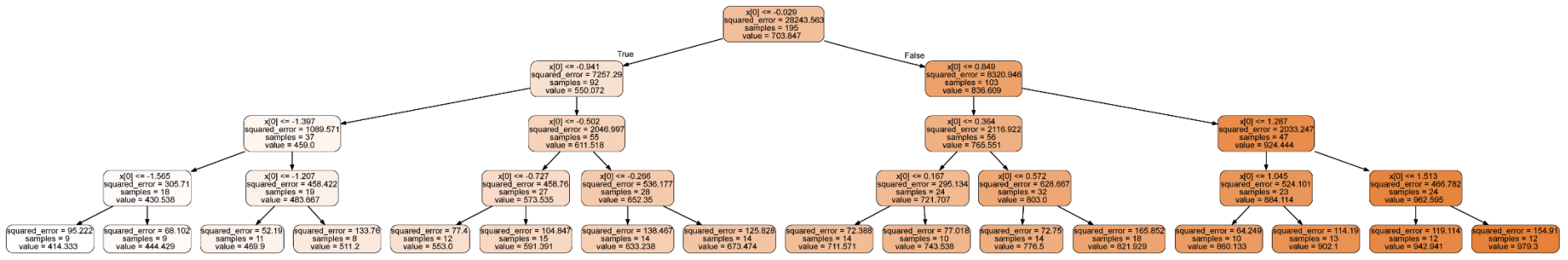
83 **Fig. S4** Example Random Forest (RF) model predicted values vs. observed values. Note: The RF model with linear univariate relationships dataset which has 10

84

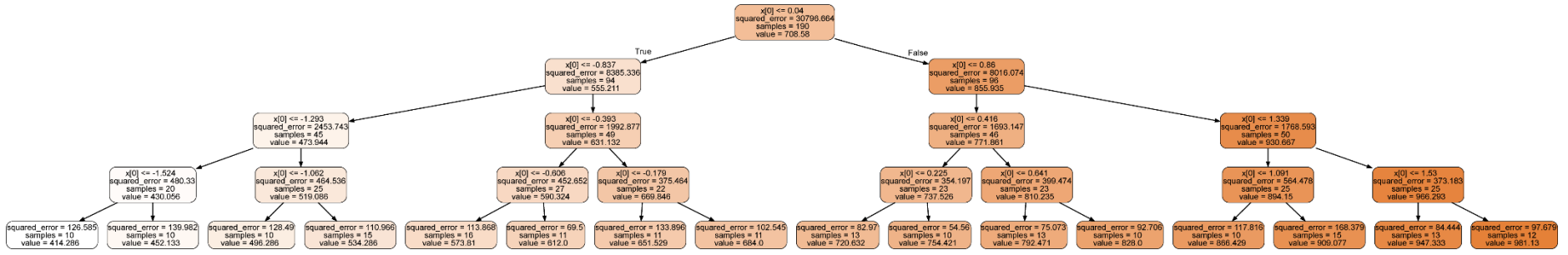
DTs with a depth of 4 each ($n_{\text{estimators}} = 10$, $\text{max_depth} = 4$).



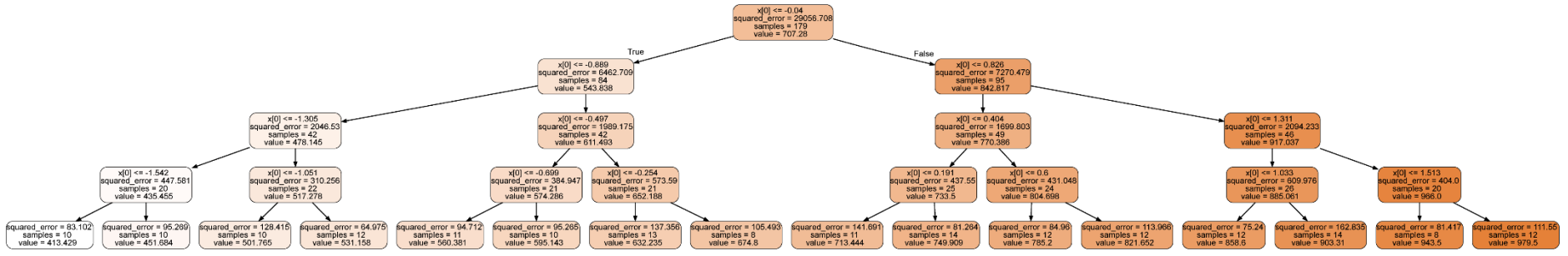
88

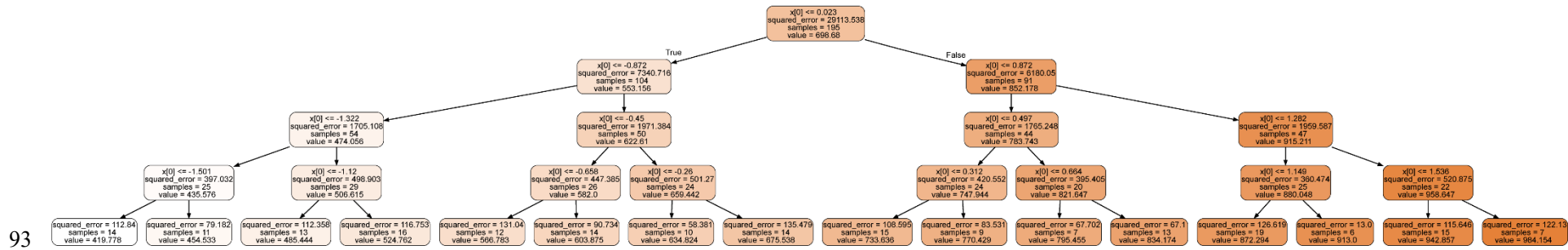
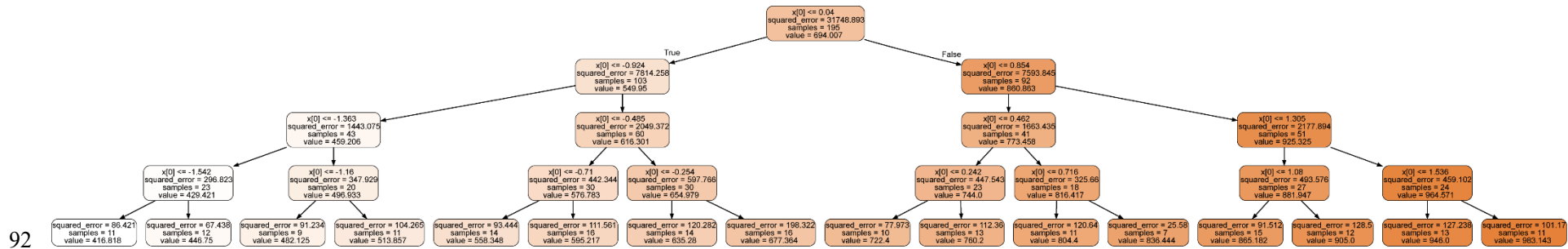
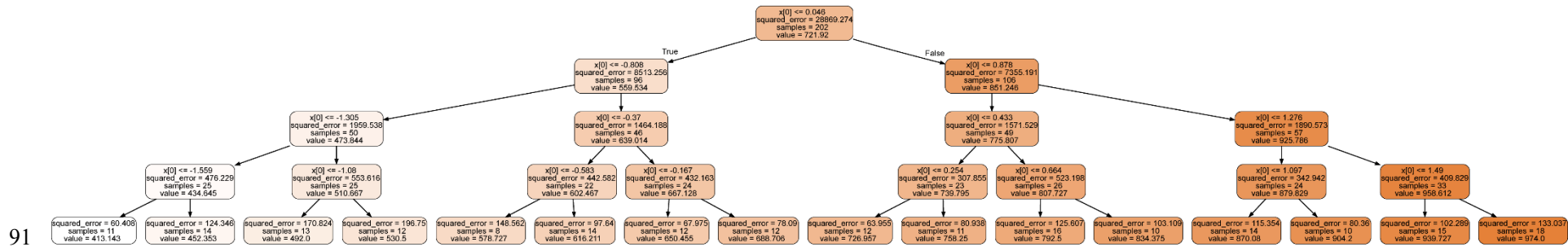


89

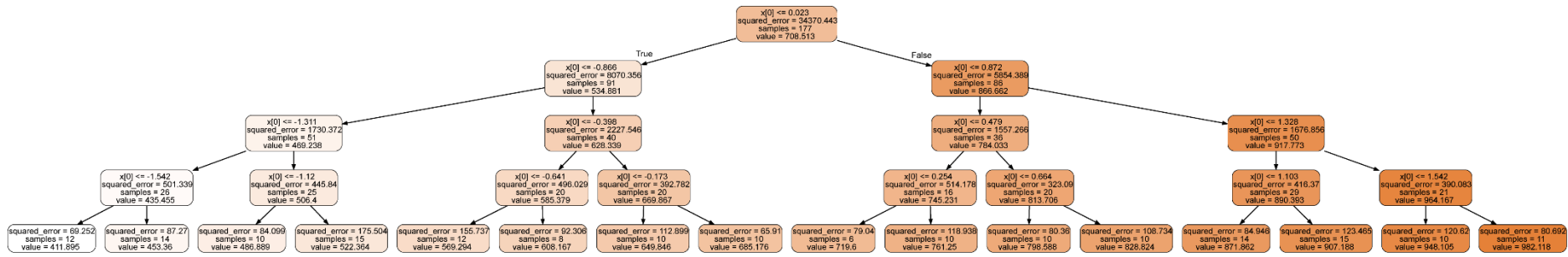


90





94



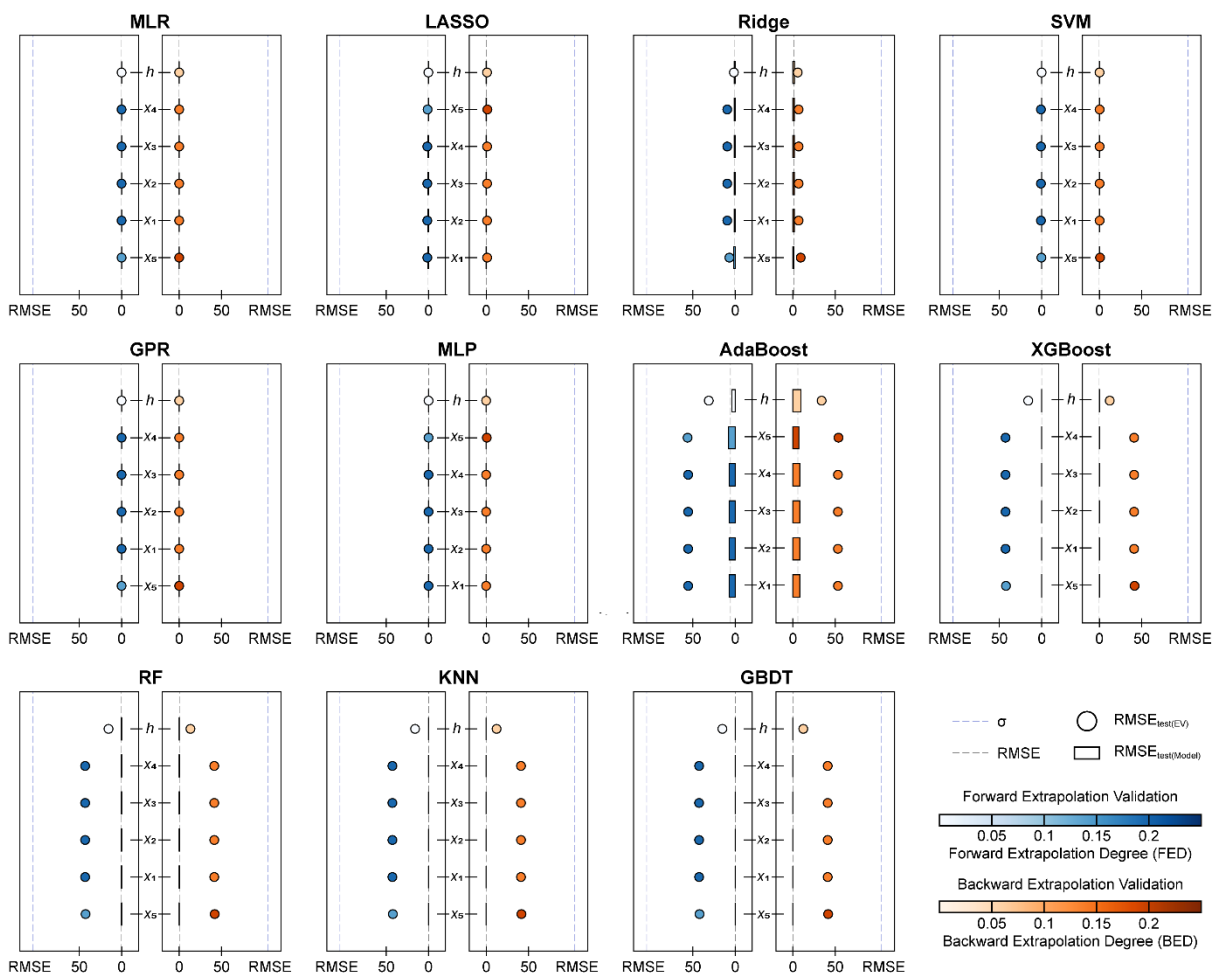
95

Fig. S5 10 decision trees (DTs) of example random forest (RF) model visualization.

96 **Text S3** Extrapolation validation of mathematical relationships models

97 Applying the EV method to evaluate the extrapolation ability of the three functional relations
98 optimal model established by 11 ML algorithms. Serialized extrapolation was performed on the *h* and
99 independent variables of the training set, and the training (EV) and test (EV) sets were re-divided in
100 the ratio of 8:2, the ED was calculated and the extrapolation ability of the model was evaluated. (Figs.
101 4a, b and c, Fig. S6 and S7).

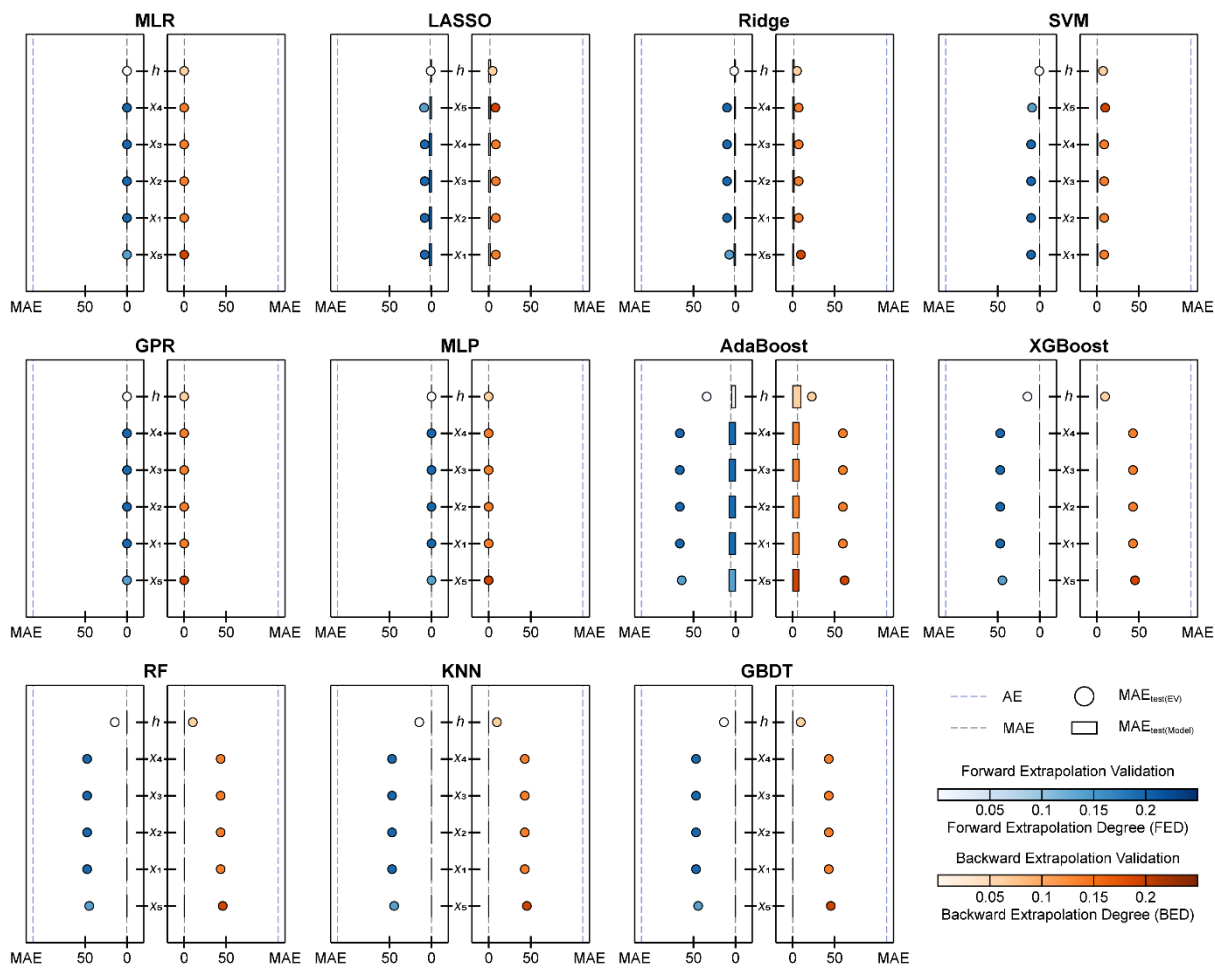
102



103

104 **Fig. S6** Extrapolation validation (EV) results of MLR, LASSO, Ridge, SVM, MLP, AdaBoost, RF, KNN,

105 GBDT and GPR Model for ML models developed with data from linear multivariate.



106

107 **Fig. S7** Extrapolation validation (EV) results of MLR, LASSO, Ridge, SVM, MLP, AdaBoost, RF, KNN,

108

GBDT and GPR Model for ML models developed with data from nonlinear multivariate.

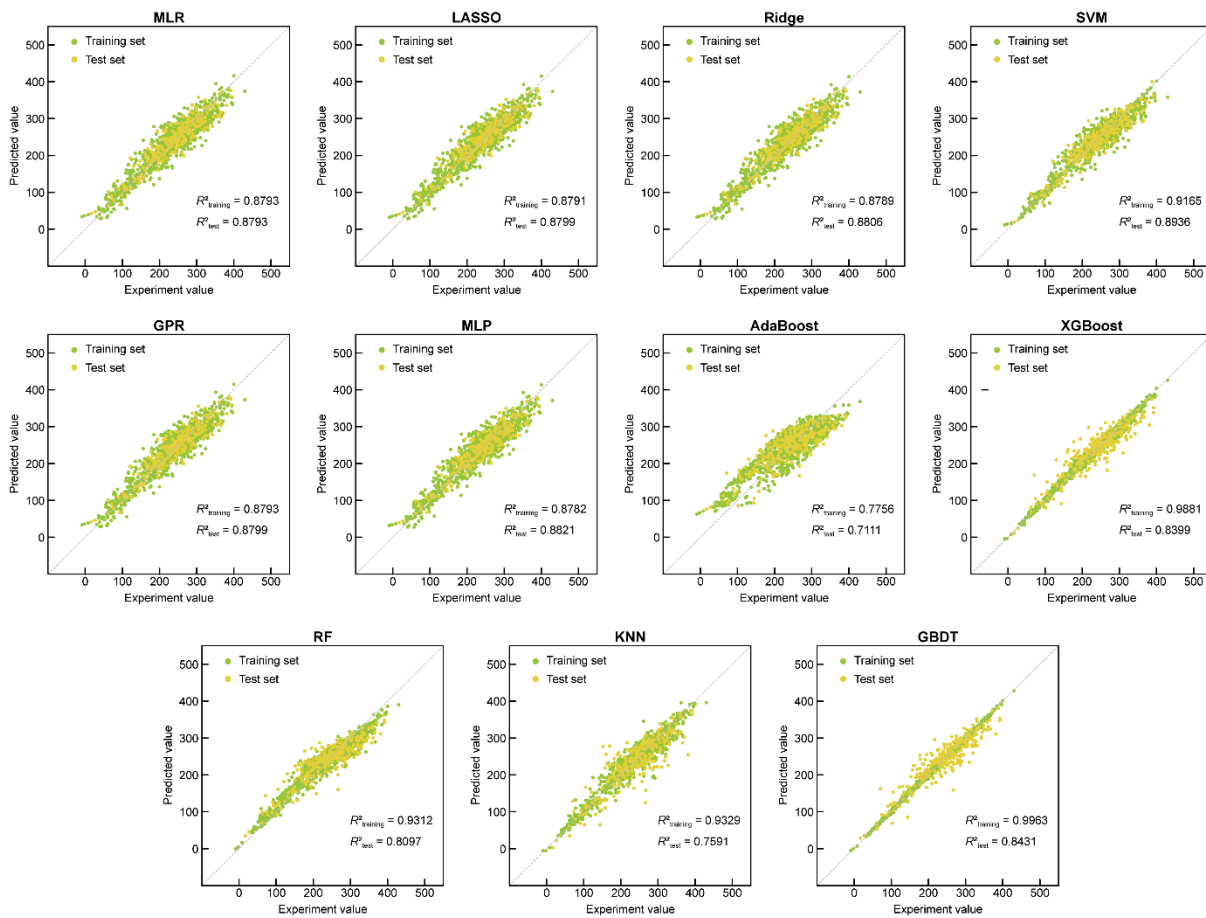
109 **Text S4** Development of PI- T_g models

110 The 29 norm indexes (I) required for MLR models developed in the literature were calculated.
111 With these 29 I s as the independent variables and the PI glass transition temperature (T_g) data
112 provided in the literature as the dependent variables, 10 ML algorithmic models for LASSO, Ridge,
113 GPR, MLP, AdaBoost, RF, GBDT, SVM, KNN and XGBoost were developed using scikit-learn¹
114 version 1.3.0 and xgboost³ version 1.7.6. The I data were standardized for modelling. The ratio of the
115 training set to the test set in the modeling is completely consistent with the settings in the literature.
116 Using R^2 and MAE as statistical parameters to evaluate the established model, the optimal model
117 hyperparameters were determined with the best average performance between the training and test
118 sets. (Table S10 and Fig. S8). To be clear, During the development of machine learning models
119 containing the parameter “random_state”, this parameter is always set to 1.

120

121 **Table S10** Hyperparameters of the PI- T_g ML model developed in this study.

Model	Parameters
Lasso	alpha = 0.042, tol = 0.0009
Ridge	alpha = 0.99, max_iter = 100
SVM	C = 750, epsilon = 1.3, gamma = 0.009, kernel = 'rbf'
GPR	kernel=kernel=2**2 * DotProduct(sigma_0=2) + WhiteKernel(noise_level=0.9)
MLP	activation = 'identity', alpha = 0.0001, hidden_layer_sizes = (150, 50), solver = 'adam'
AdaBoost	learning_rate = 0.9, n_estimators = 750
XGBoost	learning_rate = 0.09, max_depth = 4, n_estimators = 280
RF	max_depth = 9, max_features = 9, n_estimators = 150
KNN	n_neighbors = 2
GBDT	learning_rate = 0.08, max_depth = 5, n_estimator = 250



123

124

Fig. S8 PI- T_g ML models predicted vs. experimental values.

125

126 **Text S5** Extrapolation validation of PI- T_g models.

127 Evaluation of the extrapolation ability of 11 PI- T_g models by using the EV method. Serialized
128 extrapolation was performed on the h and Is of the training set, and the training (EV) and test (EV)
129 sets were re-divided in the ratio of 8:2, the ED was calculated and the extrapolation ability of the
130 model was evaluated.

131

132 **References**

133 1. Scikit-learn Scikit-learn. <https://scikit-learn.org/stable/>.

134 2. Python Python. www.python.org/.

135 3. xgboost xgboost. <https://xgboost.readthedocs.io/en/stable/>.

136 4. joblib joblib. <https://joblib.readthedocs.io/en/stable/>.

137