

## A Algorithmic Details

Figure 4 summarizes how IBUG generates a probabilistic prediction for a given input instance.

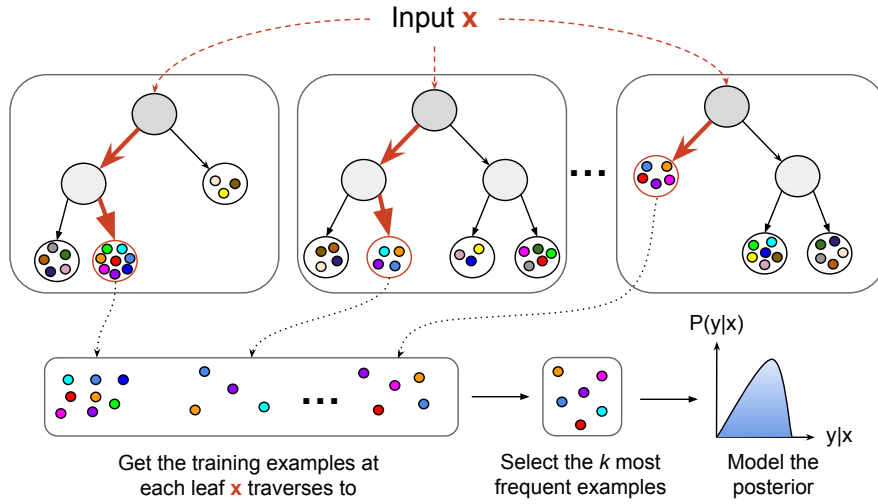


Figure 4: IBUG workflow. Given a GBRT model and an input instance  $x$ , IBUG collects the training examples at each leaf  $x$  traverses to, keeps the  $k$  most frequent examples, and then uses those examples to model the output distribution.

### A.1 Ethical Statement

In general, this work has no foreseeable negative societal impacts; however, users should carefully validate their models as imprecise uncertainty estimates may adversely affect certain domains (e.g., healthcare, weather).

## B Implementation and Experiment Details

We implement IBUG in Python, using Cython—a Python package allowing the development of C extensions—to store a unified representation of the model structure. IBUG currently supports all major modern gradient boosting frameworks including XGBoost [12], LightGBM [36], and CatBoost [52]. Experiments are run on an Intel(R) Xeon(R) CPU E5-2690 v4 @ 2.6GHz with 60GB of RAM @ 2.4GHz. We run our experiments on publicly available datasets. Links to all data sources as well as the code for IBUG and all experiments is currently available online at <https://github.com/jjbrophy47/ibug>.

**Metrics.** We use the continuous ranked probability score (CRPS) and negative log likelihood (NLL) to measure probabilistic performance. CRPS is a quadratic measure of discrepancy between the cumulative distribution function (CDF)  $F$  of forecast  $\hat{y}$  and the empirical CDF of the scalar observation  $y$ :  $\int (F(\hat{y}) - \mathbb{1}[\hat{y} \geq y])^2 d\hat{y}$  in which  $\mathbb{1}$  is the indicator function [29, 76]. To evaluate point performance, we use root mean squared error (RMSE):  $\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$ .

### B.1 Datasets

This section gives a detailed description for each dataset we use in our experiments.

- **Ames** [19] consists of 2,930 instances of housing prices in the Ames, Iowa area characterized by 80 attributes. The aim is to predict the sale price of a given house.
- **Bike** [23, 21] contains 17,379 measurements of the number of bikes rented per hour characterized by 16 attributes. The aim is to predict the number of bikes rented for a given hour.
- **California** [49] consists of 20,640 instances of median housing prices in various California districts characterized by 8 attributes. The aim is to predict the median housing price for the given district.
- **Communities** [55, 21] consists of 1,994 measurements of violent crime statistics based on crime, survey, and census data. The dataset is characterized by 100 attributes, and the aim is to predict the violent crime rate for a given population.
- **Concrete** [75, 21] consists of 1,030 instances of concrete characterized by 8 attributes. The aim is to predict the compressive strength of the concrete.
- **Energy** [70, 21] consists of 768 buildings in which each building is one of 12 different shapes and is characterized by 8 features. The aim is to predict the cooling load associated with the building.
- **Facebook** [63, 21] consists of 40,949 Facebook posts characterized by 53 attributes. The aim is to predict the number of comments for a given post.
- **Kin8nm** [73] consists of 8,192 instances of the forward kinematics of an 8 link robotic arm. The aim is to predict the forward kinematics of the robotic arm.
- **Life** [53] consists of 2,928 instances of life expectancy estimates for various countries during one year. Each instance is characterized by 20 attributes, and the aim is to predict the life expectancy of the country during a specific year.
- **MEPS** [15] consists of 16,656 instances of medical expenditure survey data. Each instance is characterized by 139 attributes, and the aim is to predict the insurance utilization for the given medical expenditure.
- **MSD** [7] consists of 515,345 songs characterized by 90 audio features constructed from each song. The aim is to predict what year the song was released based on the audio features.
- **Naval** [16, 21] consists of 11,934 instances extracted from a high-performing gas turbine simulation. Each instance is characterized by 16 features. The aim is to predict the gas turbine decay coefficient.
- **News** [21, 24] consists of 39,644 Mashable articles characterized by 60 features. The aim is to predict the number of shares for a given article.

- **Obesity** [66] contains 48,346 instances of obesity rates for different states and regions with differing socioeconomic backgrounds. Each instance is characterized by 32 attributes. The aim is to predict the obesity rate of the region.
- **Power** [21, 35, 71] contains 9,568 readings of a Combined Cycle Power Plant (CCPP) at full work load. Each reading is characterized by 4 features. The aim is to predict the net hourly electrical energy output.
- **Protein** [21] contains 45,730 tertiary-protein-structure instances characterized by 9 attributes. The aim is to predict the armstrong coefficient of the protein structure.
- **STAR** [21, 65] contains 2,161 student-teacher achievement scores characterized by 39 attributes. The aim is to predict the student-teacher achievement based on the given intervention.
- **Superconductor** [21, 31] contains 21,263 potential superconductors characterized by 81 attributes. The aim is to predict the critical temperature of the given superconductor.
- **Synthetic** [10, 27] is a non-linear synthetic regression dataset in which the inputs are independent and uniformly distributed on the interval  $[0, 1]$ ; the dataset contains 10,000 instances characterized by 100 attributes.
- **Wave** [21] consists of 287,999 positions and absorbed power outputs of wave energy converters (WECs) in four real wave scenarios off the southern coast of Australia (Sydney, Adelaide, Perth and Tasmania). The aim is to predict the total power output of a given WEC.
- **Wine** [17, 21] consists of 6,497 instances of Portuguese “Vinho Verde” red and white wine characterized by 11 features. The aim is to predict the quality of the wine from 0-10.
- **Yacht** [21] consists of 308 instances of yacht-sailing performance characterized by 6 attributes. The aim is to predict the residual resistance per unit weight of displacement.

For each dataset, we generate one-hot encodings for any categorical variable and leave all numeric and binary variables as is. Table 4 shows a summary of the datasets after preprocessing.

Table 4: Dataset summary after preprocessing.

Dataset	Source	$n$	$p$
Ames	[19]	2,930	358
Bike	[23, 21]	17,379	37
California	[49]	20,640	100
Communities	[21, 55]	1,994	100
Concrete	[75, 21]	1,030	8
Energy	[70, 21]	768	16
Facebook	[63, 21]	40,949	133
Kin8nm	[73]	8,192	8
Life	[53]	2,928	204
MEPS	[15]	15,656	139
MSD	[7]	515,345	90
Naval	[21, 16]	11,934	17
News	[21, 24]	39,644	58
Obesity	[66]	48,346	100
Power	[21, 35, 71]	9,568	4
Protein	[21]	45,730	9
STAR	[21, 65]	2,161	95
Superconductor	[21, 31]	21,263	82
Synthetic	[10, 27]	10,000	100
Wave	[21]	287,999	48
Wine	[17, 21]	6,497	11
Yacht	[21]	308	6

## B.2 Hyperparameters

Tables 5 and 6 show hyperparameter values selected most often for each dataset when optimizing CRPS and NLL, respectively. We tune nearest-neighbor hyperparameter  $k$  using values [3, 5, 7, 9, 11, 15, 31, 61, 91, 121, 151, 201, 301, 401, 501, 601, 701],  $\gamma$  and  $\delta$  using values [1e-8, 1e-7, 1e-6, 1e-5, 1e-4, 1e-3, 1e-2, 1e-1, 0, 1e0, 1e1, 1e2, 1e3] with multipliers [1.0, 2.5, 5.0], number of trees  $T$  using values [10, 25, 50, 100, 250, 500, 1000, 2000] (since NGBoost has no hyperparameters to tune besides  $T$ , we tune  $T$  on the validation set using early stopping [22]), learning rate  $\eta$  using values [0.01, 0.1], maximum number of leaves  $h$  using values [15, 31, 61, 91], minimum number of leaves  $n_{\ell_0}$  using values [1, 20], maximum depth  $d$  using values [2, 3, 5, 7, -1 (unlimited)], and  $\rho$  which selects the minimum variance computed from the validation set predictions. For the MSD and Wave datasets, we use a bagging fraction of 0.1 [22, 64].

Table 5: Hyperparameters selected most often over 10 folds for each dataset when optimizing CRPS.

Dataset	NGBoost		PGBM					CBU				
	$T$	$\gamma/\delta$	$T$	$\eta$	$h$	$n_{\ell_0}$	$\gamma/\delta$	$T$	$\eta$	$d$	$n_{\ell_0}$	$\gamma/\delta$
Ames	2000	$\gamma:1e+00$	2000	0.1	15	1	$\gamma:2e+00$	2000	0.1	7	1	$\delta:5e+03$
Bike	2000	$\gamma:5e-01$	2000	0.01	61	1	$\delta:1e-08$	2000	0.1	2	1	$\delta:5e-02$
California	2000	$\delta:3e-02$	1000	0.1	31	20	$\delta:3e-02$	2000	0.1	-1	1	$\delta:1e-01$
Communities	223	$\delta:1e-02$	500	0.01	15	20	$\gamma:1e+01$	2000	0.01	7	1	$\delta:5e-02$
Concrete	2000	$\delta:1e+00$	2000	0.1	15	20	$\delta:1e-08$	2000	0.1	5	1	$\delta:2e+00$
Energy	2000	$\gamma:5e-01$	2000	0.1	15	1	$\gamma:5e-01$	2000	0.1	3	1	$\delta:1e-01$
Facebook	2000	$\gamma:1e+00$	2000	0.01	15	1	$\gamma:2e+00$	2000	0.1	5	1	$\gamma:1e+00$
Kin8nm	581	$\delta:3e-02$	2000	0.1	61	20	$\delta:5e-02$	2000	0.1	7	1	$\delta:5e-02$
Life	2000	$\gamma:1e+00$	2000	0.1	15	1	$\delta:2e-01$	2000	0.1	5	1	$\delta:1e+00$
MEPS	583	$\delta:1e-01$	50	0.1	15	1	$\delta:1e-08$	100	0.01	-1	1	$\gamma:1e+00$
MSD	2000	$\delta:1e-01$	2000	0.01	91	20	$\gamma:1e+01$	2000	0.1	7	1	$\delta:5e-01$
Naval	2000	$\delta:0e+00$	2000	0.1	61	20	$\gamma:3e-02$	2000	0.1	7	1	$\delta:3e-04$
News	2000	$\gamma:5e-01$	100	0.01	15	20	$\delta:2e+03$	100	0.01	2	1	$\gamma:5e-01$
Obesity	2000	$\delta:1e-01$	500	0.1	91	20	$\delta:1e-08$	2000	0.1	7	1	$\delta:5e-01$
Power	2000	$\delta:2e-01$	500	0.1	91	1	$\delta:5e-01$	2000	0.1	7	1	$\delta:1e+00$
Protein	2000	$\delta:1e-01$	2000	0.1	91	20	$\gamma:2e+00$	2000	0.1	7	1	$\delta:1e+00$
STAR	187	$\delta:2e+01$	1000	0.01	15	1	$\gamma:1e+01$	2000	0.01	-1	1	$\delta:5e+01$
Superconductor	162	$\gamma:5e-01$	1000	0.01	15	20	$\delta:5e-02$	2000	0.1	-1	1	$\delta:3e-02$
Synthetic	208	$\delta:5e-01$	500	0.01	15	20	$\delta:1e+01$	2000	0.01	3	1	$\delta:1e+00$
Wave	2000	$\delta:5e+03$	2000	0.1	15	1	$\delta:1e-08$	2000	0.1	-1	1	$\delta:2e+02$
Wine	309	$\delta:5e-02$	2000	0.01	91	20	$\delta:5e-01$	2000	0.1	7	1	$\delta:2e-01$
Yacht	2000	$\gamma:5e-01$	2000	0.1	15	1	$\delta:0e+00$	2000	0.1	3	1	$\gamma:2e+00$

Dataset	CatBoost				IBUG		
	$T$	$\eta$	$d$	$n_{\ell_0}$	$k$	$\rho$	$\gamma/\delta$
Ames	2000	0.1	-1	1	5	2206	$\delta:3e-04$
Bike	2000	0.1	2	1	3	0.471	$\gamma:2e-01$
California	2000	0.1	-1	1	7	2e-15	$\delta:1e-08$
Communities	2000	0.01	-1	1	15	0.017	$\delta:0e+00$
Concrete	2000	0.1	5	1	3	0.049	$\gamma:5e-01$
Energy	2000	0.1	5	1	3	0.035	$\gamma:1e-01$
Facebook	2000	0.1	-1	1	15	0.213	$\delta:1e-01$
Kin8nm	2000	0.1	7	1	3	0.003	$\gamma:5e-01$
Life	2000	0.1	5	1	3	0.047	$\gamma:5e-01$
MEPS	250	0.01	5	1	201	1.08	$\delta:1e-07$
MSD	2000	0.1	7	1	31	1.25	$\delta:1e-07$
Naval	2000	0.1	7	1	3	1e-15	$\gamma:5e-01$
News	1000	0.01	2	1	15	163	$\gamma:5e-01$
Obesity	2000	0.1	7	1	5	0.306	$\gamma:5e-01$
Power	2000	0.1	7	1	5	0.220	$\delta:1e-01$
Protein	2000	0.1	7	1	31	0.028	$\delta:1e-01$
STAR	250	0.01	5	1	121	192	$\delta:1e-08$
Superconductor	2000	0.1	5	1	3	5e-15	$\gamma:1e-01$
Synthetic	1000	0.01	7	1	401	9.34	$\delta:1e-08$
Wave	2000	0.1	-1	1	3	2e-10	$\gamma:2e-01$
Wine	2000	0.1	7	1	15	0.268	$\delta:2e-08$
Yacht	2000	0.1	2	1	3	0.196	$\gamma:1e-01$

Table 6: Hyperparameters selected most often over 10 folds for each dataset when optimizing NLL.

Dataset	NGBoost		PGBM					CBU				
	$T$	$\gamma/\delta$	$T$	$\eta$	$h$	$n_{\ell_0}$	$\gamma/\delta$	$T$	$\eta$	$d$	$n_{\ell_0}$	$\gamma/\delta$
Ames	373	$\delta:2e+03$	2000	0.1	15	1	$\gamma:2e+00$	2000	0.1	7	1	$\gamma:1e+01$
Bike	926	$\delta:0e+00$	2000	0.01	61	1	$\delta:1e-08$	2000	0.1	2	1	$\delta:0e+00$
California	2000	$\delta:5e-02$	1000	0.1	31	20	$\delta:1e-01$	2000	0.1	-1	1	$\delta:2e-01$
Communities	156	$\delta:1e-02$	500	0.01	15	20	$\gamma:1e+01$	2000	0.01	7	1	$\delta:1e-01$
Concrete	383	$\delta:1e+00$	2000	0.1	15	20	$\delta:1e+00$	2000	0.1	5	1	$\delta:2e+00$
Energy	422	$\delta:1e-02$	2000	0.1	15	1	$\delta:1e-08$	2000	0.1	3	1	$\delta:1e-01$
Facebook	549	$\delta:0e+00$	2000	0.01	15	1	$\gamma:5e+00$	2000	0.1	5	1	$\gamma:2e+00$
Kin8nm	975	$\delta:1e-02$	2000	0.1	61	20	$\delta:5e-02$	2000	0.1	7	1	$\delta:1e-01$
Life	366	$\delta:2e-01$	2000	0.1	15	1	$\delta:1e+00$	2000	0.1	5	1	$\delta:1e+00$
MEPS	188	$\delta:1e+00$	50	0.1	15	1	$\delta:1e-08$	100	0.01	-1	1	$\delta:1e+00$
MSD	2000	$\delta:3e-02$	2000	0.01	91	20	$\gamma:1e+01$	2000	0.1	7	1	$\delta:1e+00$
Naval	2000	$\delta:5e-05$	2000	0.1	61	20	$\gamma:5e-02$	2000	0.1	7	1	$\delta:3e-04$
News	38	$\delta:1e+03$	100	0.01	15	20	$\gamma:1e+01$	100	0.01	2	1	$\delta:2e+03$
Obesity	2000	$\delta:0e+00$	500	0.1	91	20	$\delta:1e-08$	2000	0.1	7	1	$\delta:5e-01$
Power	275	$\delta:2e-01$	500	0.1	91	1	$\delta:1e+00$	2000	0.1	7	1	$\delta:2e+00$
Protein	2000	$\delta:2e-01$	2000	0.1	91	20	$\delta:2e+00$	2000	0.1	7	1	$\delta:1e+00$
STAR	176	$\delta:1e+01$	1000	0.01	15	1	$\delta:2e+02$	2000	0.01	-1	1	$\delta:5e+01$
Superconductor	378	$\gamma:1e+00$	1000	0.01	15	20	$\gamma:2e+00$	2000	0.1	-1	1	$\delta:1e-01$
Synthetic	284	$\delta:5e-01$	500	0.01	15	20	$\delta:1e+01$	2000	0.01	3	1	$\delta:1e+00$
Wave	2000	$\gamma:1e+00$	2000	0.1	15	1	$\delta:5e+02$	2000	0.1	-1	1	$\delta:2e+02$
Wine	390	$\delta:5e-02$	2000	0.01	91	20	$\gamma:2e+01$	2000	0.1	7	1	$\delta:5e-01$
Yacht	356	$\delta:0e+00$	2000	0.1	15	1	$\delta:5e-02$	2000	0.1	3	1	$\delta:5e-01$

Dataset	CatBoost				IBUG		
	$T$	$\eta$	$d$	$n_{\ell_0}$	$k$	$\rho$	$\gamma/\delta$
Ames	2000	0.1	-1	1	11	4673	$\delta:1e-08$
Bike	2000	0.1	2	1	5	0.4	$\gamma:2e-01$
California	2000	0.1	-1	1	31	0.063	$\delta:0e+00$
Communities	2000	0.01	-1	1	61	0.026	$\delta:0e+00$
Concrete	2000	0.1	5	1	5	0.56	$\delta:1e-08$
Energy	2000	0.1	5	1	3	0.087	$\gamma:2e-01$
Facebook	2000	0.1	-1	1	301	0.175	$\delta:1e-01$
Kin8nm	2000	0.1	7	1	7	0.031	$\delta:0e+00$
Life	2000	0.1	5	1	7	0.22	$\delta:2e-08$
MEPS	250	0.01	5	1	301	1.76	$\delta:1e+00$
MSD	2000	0.1	7	1	61	1.75	$\delta:1e-07$
Naval	2000	0.1	7	1	5	4e-04	$\gamma:5e-01$
News	1000	0.01	2	1	301	994	$\delta:2e+03$
Obesity	2000	0.1	7	1	9	0.529	$\delta:1e-07$
Power	2000	0.1	7	1	15	0.861	$\delta:1e-07$
Protein	2000	0.1	7	1	121	0.218	$\delta:5e-08$
STAR	250	0.01	5	1	121	189	$\delta:1e-05$
Superconductor	2000	0.1	5	1	7	0.019	$\gamma:2e-01$
Synthetic	1000	0.01	7	1	401	9.39	$\delta:1e-08$
Wave	2000	0.1	-1	1	31	349	$\gamma:2e-01$
Wine	2000	0.1	7	1	61	0.297	$\delta:2e-08$
Yacht	2000	0.1	2	1	3	0.196	$\gamma:2e-01$

### B.3 Additional Metrics

In this section, we show results for point performance and probabilistic performance with additional metrics. Each table shows average results over the 10 random folds for each dataset, with standard errors in subscripted parentheses. We use the *Uncertainty Toolbox*<sup>3</sup> [14] to compute each metric. Lower is better for all metrics.

**Point performance and negative-log likelihood.** Tables 7 and 8 show point (RMSE) and probabilistic (NLL) performance of each method.

Table 7: Point (RMSE ↓) performance for each method on each dataset.

Dataset	NGBoost	PGBM	CBU	IBUG	IBUG+CBU
Ames	24580 <sub>(804)</sub>	<b>23541</b> <sub>(1225)</sub>	<b>22576</b> <sub>(924)</sub>	<b>22942</b> <sub>(1388)</sub>	<b>22391</b> <sub>(1119)</sub>
Bike	4.173 <sub>(0.076)</sub>	3.812 <sub>(0.225)</sub>	<b>2.850</b> <sub>(0.192)</sub>	<b>2.826</b> <sub>(0.200)</sub>	<b>2.708</b> <sub>(0.202)</sub>
California	0.503 <sub>(0.003)</sub>	0.445 <sub>(0.001)</sub>	0.449 <sub>(0.002)</sub>	<b>0.432</b> <sub>(0.001)</sub>	0.434 <sub>(0.002)</sub>
Communities	0.137 <sub>(0.004)</sub>	<b>0.135</b> <sub>(0.004)</sub>	<b>0.133</b> <sub>(0.004)</sub>	<b>0.133</b> <sub>(0.004)</sub>	<b>0.132</b> <sub>(0.004)</sub>
Concrete	5.485 <sub>(0.182)</sub>	3.840 <sub>(0.209)</sub>	<b>3.682</b> <sub>(0.202)</sub>	<b>3.629</b> <sub>(0.183)</sub>	<b>3.617</b> <sub>(0.188)</sub>
Energy	0.461 <sub>(0.030)</sub>	0.291 <sub>(0.022)</sub>	0.381 <sub>(0.023)</sub>	<b>0.264</b> <sub>(0.023)</sub>	0.303 <sub>(0.023)</sub>
Facebook	<b>20.8</b> <sub>(1.102)</sub>	<b>20.5</b> <sub>(0.867)</sub>	<b>20.1</b> <sub>(0.913)</sub>	<b>20.0</b> <sub>(0.903)</sub>	<b>19.9</b> <sub>(0.929)</sub>
Kin8nm	0.176 <sub>(0.001)</sub>	0.108 <sub>(0.001)</sub>	0.103 <sub>(0.001)</sub>	<b>0.086</b> <sub>(0.001)</sub>	0.091 <sub>(0.001)</sub>
Life	2.280 <sub>(0.032)</sub>	1.678 <sub>(0.059)</sub>	<b>1.637</b> <sub>(0.058)</sub>	<b>1.652</b> <sub>(0.055)</sub>	<b>1.610</b> <sub>(0.056)</sub>
MEPS	<b>23.7</b> <sub>(0.955)</sub>	<b>24.1</b> <sub>(0.760)</sub>	<b>23.5</b> <sub>(0.950)</sub>	<b>23.7</b> <sub>(0.932)</sub>	<b>23.6</b> <sub>(0.945)</sub>
MSD	9.121 <sub>(0.010)</sub>	8.804 <sub>(0.008)</sub>	8.743 <sub>(0.008)</sub>	8.747 <sub>(0.008)</sub>	<b>8.722</b> <sub>(0.008)</sub>
Naval	0.002 <sub>(0.000)</sub>	0.001 <sub>(0.000)</sub>	0.001 <sub>(0.000)</sub>	<b>0.000</b> <sub>(0.000)</sub>	<b>0.000</b> <sub>(0.000)</sub>
News	<b>11162</b> <sub>(1153)</sub>	<b>11047</b> <sub>(1106)</sub>	<b>11036</b> <sub>(1118)</sub>	<b>11036</b> <sub>(1116)</sub>	<b>11032</b> <sub>(1118)</sub>
Obesity	5.315 <sub>(0.022)</sub>	3.658 <sub>(0.033)</sub>	<b>3.572</b> <sub>(0.038)</sub>	<b>3.576</b> <sub>(0.037)</sub>	<b>3.567</b> <sub>(0.037)</sub>
Power	3.836 <sub>(0.045)</sub>	3.017 <sub>(0.056)</sub>	<b>2.924</b> <sub>(0.065)</sub>	<b>2.941</b> <sub>(0.059)</sub>	<b>2.912</b> <sub>(0.063)</sub>
Protein	4.525 <sub>(0.040)</sub>	<b>3.455</b> <sub>(0.021)</sub>	3.520 <sub>(0.019)</sub>	3.512 <sub>(0.017)</sub>	3.493 <sub>(0.018)</sub>
STAR	233 <sub>(2.388)</sub>	<b>229</b> <sub>(2.076)</sub>	<b>229</b> <sub>(1.850)</sub>	<b>228</b> <sub>(1.985)</sub>	<b>228</b> <sub>(1.857)</sub>
Superconductor	<b>0.170</b> <sub>(0.101)</sub>	0.425 <sub>(0.091)</sub>	0.463 <sub>(0.087)</sub>	0.427 <sub>(0.088)</sub>	0.419 <sub>(0.089)</sub>
Synthetic	<b>10.2</b> <sub>(0.068)</sub>	<b>10.1</b> <sub>(0.072)</sub>	<b>10.2</b> <sub>(0.072)</sub>	<b>10.1</b> <sub>(0.073)</sub>	<b>10.1</b> <sub>(0.073)</sub>
Wave	13537 <sub>(32.7)</sub>	7895 <sub>(86.0)</sub>	4803 <sub>(37.5)</sub>	4899 <sub>(55.0)</sub>	<b>4020</b> <sub>(33.5)</sub>
Wine	0.693 <sub>(0.010)</sub>	<b>0.603</b> <sub>(0.010)</sub>	0.626 <sub>(0.010)</sub>	<b>0.596</b> <sub>(0.012)</sub>	<b>0.598</b> <sub>(0.011)</sub>
Yacht	0.761 <sub>(0.106)</sub>	0.809 <sub>(0.103)</sub>	<b>0.677</b> <sub>(0.124)</sub>	<b>0.668</b> <sub>(0.125)</sub>	<b>0.645</b> <sub>(0.124)</sub>
IBUG W-TL	16-5-1	13-8-1	6-16-0	-	2-13-7
IBUG+CBU W-TL	18-3-1	12-9-1	16-6-0	7-13-2	-

Table 8: Probabilistic (NLL ↓) performance for each method on each dataset.

Dataset	NGBoost	PGBM	CBU	IBUG	IBUG+CBU
Ames	11.3 <sub>(0.018)</sub>	11.3 <sub>(0.029)</sub>	11.9 <sub>(0.140)</sub>	<b>11.2</b> <sub>(0.030)</sub>	11.5 <sub>(0.092)</sub>
Bike	1.942 <sub>(0.024)</sub>	1.929 <sub>(0.078)</sub>	<b>1.184</b> <sub>(0.034)</sub>	1.886 <sub>(0.056)</sub>	1.382 <sub>(0.042)</sub>
California	0.545 <sub>(0.007)</sub>	0.580 <sub>(0.005)</sub>	0.524 <sub>(0.004)</sub>	0.477 <sub>(0.010)</sub>	<b>0.437</b> <sub>(0.016)</sub>
Communities	<b>-0.697</b> <sub>(0.045)</sub>	<b>-0.666</b> <sub>(0.034)</sub>	<b>-0.614</b> <sub>(0.109)</sub>	<b>-0.639</b> <sub>(0.135)</sub>	<b>-0.665</b> <sub>(0.116)</sub>
Concrete	3.043 <sub>(0.030)</sub>	2.802 <sub>(0.083)</sub>	<b>2.766</b> <sub>(0.086)</sub>	2.980 <sub>(0.146)</sub>	<b>2.695</b> <sub>(0.060)</sub>
Energy	0.604 <sub>(0.192)</sub>	<b>0.322</b> <sub>(0.182)</sub>	<b>0.406</b> <sub>(0.116)</sub>	1.644 <sub>(0.514)</sub>	0.658 <sub>(0.165)</sub>
Facebook	<b>2.102</b> <sub>(0.026)</sub>	3.116 <sub>(0.077)</sub>	2.574 <sub>(0.191)</sub>	2.175 <sub>(0.067)</sub>	2.276 <sub>(0.140)</sub>
Kin8nm	-0.414 <sub>(0.007)</sub>	-0.774 <sub>(0.034)</sub>	-0.772 <sub>(0.008)</sub>	<b>-0.841</b> <sub>(0.008)</sub>	<b>-0.847</b> <sub>(0.010)</sub>
Life	2.163 <sub>(0.029)</sub>	1.943 <sub>(0.033)</sub>	1.932 <sub>(0.079)</sub>	1.858 <sub>(0.033)</sub>	<b>1.783</b> <sub>(0.041)</sub>
MEPS	<b>3.722</b> <sub>(0.050)</sub>	3.902 <sub>(0.049)</sub>	<b>3.699</b> <sub>(0.038)</sub>	3.793 <sub>(0.052)</sub>	<b>3.675</b> <sub>(0.041)</sub>
MSD	3.454 <sub>(0.002)</sub>	3.571 <sub>(0.002)</sub>	3.415 <sub>(0.001)</sub>	3.415 <sub>(0.002)</sub>	<b>3.393</b> <sub>(0.001)</sub>
Naval	-5.408 <sub>(0.007)</sub>	-5.064 <sub>(0.338)</sub>	-6.141 <sub>(0.013)</sub>	-6.208 <sub>(0.010)</sub>	<b>-6.284</b> <sub>(0.007)</sub>
News	10.9 <sub>(0.268)</sub>	<b>10.7</b> <sub>(0.339)</sub>	<b>10.6</b> <sub>(0.205)</sub>	<b>10.6</b> <sub>(0.208)</sub>	<b>10.6</b> <sub>(0.192)</sub>
Obesity	2.940 <sub>(0.003)</sub>	2.604 <sub>(0.015)</sub>	<b>2.439</b> <sub>(0.009)</sub>	2.646 <sub>(0.009)</sub>	2.515 <sub>(0.010)</sub>
Power	2.752 <sub>(0.032)</sub>	<b>2.518</b> <sub>(0.021)</sub>	2.538 <sub>(0.019)</sub>	2.575 <sub>(0.036)</sub>	<b>2.514</b> <sub>(0.017)</sub>
Protein	2.840 <sub>(0.014)</sub>	2.661 <sub>(0.005)</sub>	2.553 <sub>(0.009)</sub>	2.653 <sub>(0.054)</sub>	<b>2.516</b> <sub>(0.010)</sub>
STAR	6.869 <sub>(0.013)</sub>	6.866 <sub>(0.012)</sub>	<b>6.866</b> <sub>(0.014)</sub>	<b>6.853</b> <sub>(0.008)</sub>	<b>6.852</b> <sub>(0.009)</sub>
Superconductor	<b>12.2</b> <sub>(13.1)</sub>	<b>0.035</b> <sub>(0.095)</sub>	<b>-0.014</b> <sub>(0.078)</sub>	0.783 <sub>(0.181)</sub>	0.108 <sub>(0.036)</sub>
Synthetic	3.745 <sub>(0.007)</sub>	<b>3.742</b> <sub>(0.006)</sub>	<b>3.741</b> <sub>(0.008)</sub>	<b>3.738</b> <sub>(0.007)</sub>	<b>3.738</b> <sub>(0.007)</sub>
Wave	10.7 <sub>(0.002)</sub>	10.3 <sub>(0.021)</sub>	<b>9.675</b> <sub>(0.003)</sub>	10.5 <sub>(0.030)</sub>	9.760 <sub>(0.046)</sub>
Wine	1.025 <sub>(0.013)</sub>	0.952 <sub>(0.020)</sub>	1.025 <sub>(0.028)</sub>	<b>0.910</b> <sub>(0.016)</sub>	0.933 <sub>(0.012)</sub>
Yacht	0.905 <sub>(0.232)</sub>	<b>0.357</b> <sub>(0.162)</sub>	0.951 <sub>(0.252)</sub>	1.799 <sub>(1.307)</sub>	0.840 <sub>(0.310)</sub>
IBUG W-TL	12-10-0	7-11-4	5-11-6	-	2-8-12
IBUG+CBU W-TL	15-6-1	10-10-2	13-6-3	12-8-2	-

<sup>3</sup><https://uncertainty-toolbox.github.io/>

**Check and interval scores.** Tables 9 and 10 show results when measuring performance with two additional proper scoring rules [29], *check score* (a.k.a. “pinball loss”) and *interval score* (evaluation using a pair of quantiles with expected coverage). Under these additional metrics, IBUG+CBU still outperform all other approaches.

Table 9: Probabilistic (check score a.k.a. “pinball loss”  $\downarrow$ ) performance.

Dataset	NGBoost	PGBM	CBU	IBUG	IBUG+CBU
Ames	19358 <sup>(276)</sup>	5487 <sup>(179)</sup>	5551 <sup>(167)</sup>	<b>5266</b> <sup>(185)</sup>	<b>5145</b> <sup>(186)</sup>
Bike	6.264 <sup>(0.482)</sup>	0.597 <sup>(0.020)</sup>	0.420 <sup>(0.018)</sup>	0.490 <sup>(0.024)</sup>	<b>0.386</b> <sup>(0.016)</sup>
California	8e+10 <sup>(8e+10)</sup>	0.112 <sup>(4e-04)</sup>	0.110 <sup>(4e-04)</sup>	0.107 <sup>(5e-04)</sup>	<b>0.104</b> <sup>(4e-04)</sup>
Communities	0.034 <sup>(0.001)</sup>	0.034 <sup>(1e-03)</sup>	0.034 <sup>(9e-04)</sup>	<b>0.033</b> <sup>(9e-04)</sup>	<b>0.033</b> <sup>(9e-04)</sup>
Concrete	1.722 <sup>(0.092)</sup>	0.972 <sup>(0.043)</sup>	<b>0.902</b> <sup>(0.039)</sup>	0.932 <sup>(0.049)</sup>	<b>0.878</b> <sup>(0.041)</sup>
Energy	0.262 <sup>(0.022)</sup>	<b>0.074</b> <sup>(0.003)</sup>	0.099 <sup>(0.005)</sup>	<b>0.072</b> <sup>(0.005)</sup>	0.079 <sup>(0.004)</sup>
Facebook	2.024 <sup>(0.049)</sup>	1.788 <sup>(0.047)</sup>	1.617 <sup>(0.030)</sup>	1.551 <sup>(0.033)</sup>	<b>1.502</b> <sup>(0.035)</sup>
Kin8nm	0.048 <sup>(3e-04)</sup>	0.031 <sup>(5e-04)</sup>	0.029 <sup>(3e-04)</sup>	<b>0.026</b> <sup>(3e-04)</sup>	<b>0.026</b> <sup>(3e-04)</sup>
Life	1.462 <sup>(0.739)</sup>	0.411 <sup>(0.014)</sup>	0.389 <sup>(0.012)</sup>	0.400 <sup>(0.011)</sup>	<b>0.368</b> <sup>(0.011)</sup>
MEPS	<b>2.779</b> <sup>(0.098)</sup>	3.246 <sup>(0.046)</sup>	3.050 <sup>(0.055)</sup>	3.100 <sup>(0.057)</sup>	3.033 <sup>(0.056)</sup>
MSD	2.283 <sup>(0.003)</sup>	2.310 <sup>(0.002)</sup>	2.203 <sup>(0.002)</sup>	2.226 <sup>(0.002)</sup>	<b>2.195</b> <sup>(0.002)</sup>
Naval	0.002 <sup>(3e-05)</sup>	2e-04 <sup>(2e-05)</sup>	2e-04 <sup>(2e-06)</sup>	1e-04 <sup>(1e-06)</sup>	<b>1e-04</b> <sup>(8e-07)</sup>
News	<b>1102</b> <sup>(23.7)</sup>	1188 <sup>(26.3)</sup>	1181 <sup>(26.2)</sup>	1280 <sup>(20.5)</sup>	1198 <sup>(26.0)</sup>
Obesity	1.620 <sup>(0.014)</sup>	0.939 <sup>(0.011)</sup>	<b>0.879</b> <sup>(0.009)</sup>	0.941 <sup>(0.010)</sup>	0.894 <sup>(0.009)</sup>
Power	1.063 <sup>(0.012)</sup>	0.773 <sup>(0.010)</sup>	<b>0.744</b> <sup>(0.011)</sup>	0.778 <sup>(0.010)</sup>	<b>0.743</b> <sup>(0.011)</sup>
Protein	2739 <sup>(2730)</sup>	0.920 <sup>(0.006)</sup>	0.902 <sup>(0.005)</sup>	0.900 <sup>(0.004)</sup>	<b>0.880</b> <sup>(0.004)</sup>
STAR	66.6 <sup>(0.803)</sup>	<b>65.9</b> <sup>(0.697)</sup>	<b>65.7</b> <sup>(0.647)</sup>	<b>65.4</b> <sup>(0.613)</sup>	<b>65.4</b> <sup>(0.605)</sup>
Superconductor	1.215 <sup>(0.014)</sup>	<b>0.064</b> <sup>(0.002)</sup>	0.076 <sup>(0.002)</sup>	0.077 <sup>(0.003)</sup>	<b>0.064</b> <sup>(0.002)</sup>
Synthetic	2.918 <sup>(0.021)</sup>	<b>2.897</b> <sup>(0.020)</sup>	<b>2.898</b> <sup>(0.020)</sup>	<b>2.894</b> <sup>(0.020)</sup>	<b>2.894</b> <sup>(0.020)</sup>
Wave	2.9e+05 <sup>(446)</sup>	1964 <sup>(37.3)</sup>	1186 <sup>(5.194)</sup>	1350 <sup>(8.028)</sup>	<b>1023</b> <sup>(4.813)</sup>
Wine	0.194 <sup>(0.002)</sup>	<b>0.163</b> <sup>(0.003)</sup>	0.170 <sup>(0.003)</sup>	<b>0.162</b> <sup>(0.003)</sup>	<b>0.162</b> <sup>(0.003)</sup>
Yacht	0.594 <sup>(0.080)</sup>	<b>0.147</b> <sup>(0.021)</sup>	<b>0.142</b> <sup>(0.024)</sup>	<b>0.139</b> <sup>(0.024)</sup>	<b>0.128</b> <sup>(0.023)</sup>
IBUG W-T-L	17-3-2	11-9-2	9-5-8	-	1-6-15
IBUG+CBU W-T-L	17-3-2	15-6-1	18-2-2	15-6-1	-

Table 10: Probabilistic (interval score  $\downarrow$ ) performance.

Dataset	NGBoost	PGBM	CBU	IBUG	IBUG+CBU
Ames	2.0e+05 <sup>(3492)</sup>	59165 <sup>(1952)</sup>	66337 <sup>(2499)</sup>	<b>57219</b> <sup>(1941)</sup>	<b>55551</b> <sup>(1994)</sup>
Bike	66.4 <sup>(6.425)</sup>	7.048 <sup>(0.411)</sup>	<b>4.270</b> <sup>(0.136)</sup>	6.775 <sup>(0.324)</sup>	<b>4.263</b> <sup>(0.191)</sup>
California	1e+12 <sup>(1e+12)</sup>	1.257 <sup>(0.008)</sup>	1.168 <sup>(0.008)</sup>	1.230 <sup>(0.020)</sup>	<b>1.119</b> <sup>(0.006)</sup>
Communities	0.361 <sup>(0.013)</sup>	0.366 <sup>(0.012)</sup>	0.352 <sup>(0.011)</sup>	<b>0.343</b> <sup>(0.013)</sup>	<b>0.339</b> <sup>(0.011)</sup>
Concrete	17.2 <sup>(0.917)</sup>	11.1 <sup>(0.698)</sup>	<b>10.3</b> <sup>(0.491)</sup>	12.1 <sup>(0.800)</sup>	<b>10.1</b> <sup>(0.523)</sup>
Energy	2.711 <sup>(0.198)</sup>	<b>0.814</b> <sup>(0.050)</sup>	0.998 <sup>(0.067)</sup>	0.912 <sup>(0.083)</sup>	<b>0.819</b> <sup>(0.064)</sup>
Facebook	28.4 <sup>(0.909)</sup>	26.8 <sup>(1.125)</sup>	21.6 <sup>(0.692)</sup>	<b>17.4</b> <sup>(0.476)</sup>	<b>17.1</b> <sup>(0.509)</sup>
Kin8nm	0.458 <sup>(0.003)</sup>	0.311 <sup>(0.007)</sup>	0.292 <sup>(0.005)</sup>	0.302 <sup>(0.009)</sup>	<b>0.262</b> <sup>(0.005)</sup>
Life	17.5 <sup>(9.900)</sup>	5.051 <sup>(0.239)</sup>	4.617 <sup>(0.198)</sup>	5.093 <sup>(0.264)</sup>	<b>4.332</b> <sup>(0.207)</sup>
MEPS	42.2 <sup>(1.973)</sup>	44.3 <sup>(1.294)</sup>	<b>37.7</b> <sup>(1.223)</sup>	<b>38.3</b> <sup>(1.375)</sup>	<b>37.2</b> <sup>(1.254)</sup>
MSD	24.5 <sup>(0.039)</sup>	24.8 <sup>(0.035)</sup>	22.3 <sup>(0.020)</sup>	22.4 <sup>(0.029)</sup>	<b>22.0</b> <sup>(0.025)</sup>
Naval	0.014 <sup>(3e-04)</sup>	0.003 <sup>(3e-04)</sup>	0.002 <sup>(2e-05)</sup>	0.001 <sup>(3e-05)</sup>	<b>0.001</b> <sup>(1e-05)</sup>
News	<b>16557</b> <sup>(519)</sup>	<b>16242</b> <sup>(556)</sup>	<b>16166</b> <sup>(580)</sup>	18694 <sup>(373)</sup>	<b>16426</b> <sup>(551)</sup>
Obesity	15.5 <sup>(0.153)</sup>	9.731 <sup>(0.125)</sup>	<b>8.747</b> <sup>(0.083)</sup>	10.7 <sup>(0.139)</sup>	9.162 <sup>(0.086)</sup>
Power	10.6 <sup>(0.136)</sup>	8.146 <sup>(0.122)</sup>	<b>7.837</b> <sup>(0.165)</sup>	8.512 <sup>(0.152)</sup>	<b>7.803</b> <sup>(0.156)</sup>
Protein	36689 <sup>(36570)</sup>	10.1 <sup>(0.149)</sup>	9.277 <sup>(0.062)</sup>	9.322 <sup>(0.045)</sup>	<b>8.853</b> <sup>(0.052)</sup>
STAR	642 <sup>(7.014)</sup>	637 <sup>(6.564)</sup>	<b>636</b> <sup>(6.131)</sup>	<b>630</b> <sup>(4.545)</sup>	<b>630</b> <sup>(4.968)</sup>
Superconductor	12.0 <sup>(0.133)</sup>	0.776 <sup>(0.023)</sup>	0.755 <sup>(0.030)</sup>	1.150 <sup>(0.060)</sup>	<b>0.692</b> <sup>(0.033)</sup>
Synthetic	28.4 <sup>(0.228)</sup>	<b>28.1</b> <sup>(0.188)</sup>	<b>28.1</b> <sup>(0.211)</sup>	<b>28.0</b> <sup>(0.197)</sup>	<b>28.0</b> <sup>(0.199)</sup>
Wave	3e+06 <sup>(3727)</sup>	20256 <sup>(323)</sup>	11748 <sup>(55.8)</sup>	16669 <sup>(117)</sup>	<b>10569</b> <sup>(47.4)</sup>
Wine	1.930 <sup>(0.023)</sup>	<b>1.723</b> <sup>(0.032)</sup>	1.793 <sup>(0.035)</sup>	<b>1.716</b> <sup>(0.030)</sup>	<b>1.692</b> <sup>(0.031)</sup>
Yacht	5.798 <sup>(0.808)</sup>	<b>1.621</b> <sup>(0.248)</sup>	<b>1.796</b> <sup>(0.419)</sup>	<b>1.955</b> <sup>(0.433)</sup>	<b>1.619</b> <sup>(0.406)</sup>
IBUG W-T-L	18-3-1	8-9-5	4-8-10	-	0-6-16
IBUG+CBU W-T-L	18-4-0	16-6-0	16-4-2	16-6-0	-

**Calibration error.** Table 11 shows the average MACE (mean absolute calibration error) and sharpness scores. Sharpness quantifies the average of the standard deviations and thus does not depend on the actual ground-truth label; therefore, MACE and sharpness are shown together, with better methods having both low calibration error and low sharpness scores.

We observe that NGBoost is particularly well-calibrated, but lacks sharpness, meaning the prediction intervals of NGBoost are generally too wide. PGBM tends to have very sharp prediction intervals, but high calibration error. In contrast, CBU tends to achieve both low calibration error and high sharpness in relation to the other methods. However, these results are with variance calibration (§3.2), which we note has a significant impact on the CBU approach. For example, the median improvement in MACE score (over datasets) for CBU when using variance calibration vs. without is greater than 3x.

Table 11: Probabilistic (MACE ↓ / sharpness ↓) performance. Standard errors are omitted for brevity.

Dataset	NGBoost	PGBM	CBU	IBUG	IBUG+CBU
Ames	0.082/74148	<b>0.040/18432</b>	0.073/18867	0.068/23186	0.063/19791
Bike	0.070/190	0.140/2.077	<b>0.045/2.136</b>	0.096/ <b>1.272</b>	0.051/1.595
California	<b>0.014/3e+13</b>	0.053/ <b>0.344</b>	0.021/0.367	0.089/0.382	0.037/0.364
Communities	0.039/0.129	0.067/ <b>0.120</b>	0.051/0.136	<b>0.035/0.133</b>	0.048/0.133
Concrete	0.056/6.889	0.068/3.002	0.096/3.177	0.115/ <b>2.503</b>	<b>0.054/2.708</b>
Energy	0.127/1.497	0.093/0.252	0.054/0.373	0.103/ <b>0.249</b>	<b>0.053/0.296</b>
Facebook	0.094/9.171	0.206/ <b>4.309</b>	0.072/7.332	<b>0.061/18.9</b>	0.091/12.5
Kin8nm	0.020/0.182	0.037/0.108	<b>0.020/0.096</b>	0.126/ <b>0.071</b>	0.045/0.081
Life	<b>0.039/111</b>	0.069/ <b>1.103</b>	0.079/1.189	0.115/1.401	0.069/1.216
MEPS	<b>0.030/6.680</b>	0.074/8.200	0.119/14.1	0.086/17.2	0.106/15.3
MSD	<b>0.007/7.749</b>	0.036/ <b>7.436</b>	0.012/8.137	0.039/9.088	0.031/8.519
Naval	<b>0.032/0.006</b>	0.279/1e-03	0.048/6e-04	0.059/ <b>5e-04</b>	0.086/5e-04
News	0.104/ <b>2170</b>	<b>0.085/3289</b>	0.101/2975	0.202/4803	0.109/3498
Obesity	0.012/5.996	0.065/3.451	<b>0.006/3.102</b>	0.095/2.957	0.043/ <b>2.956</b>
Power	0.020/3.761	0.026/2.558	<b>0.018/2.299</b>	0.030/3.328	0.019/2.729
Protein	0.029/2e+06	0.076/ <b>2.823</b>	0.037/3.144	<b>0.016/3.977</b>	0.046/3.498
STAR	0.025/248	0.031/250	0.030/ <b>242</b>	<b>0.023/245</b>	0.025/243
Superconductor	0.074/7.993	0.102/0.240	<b>0.028/0.322</b>	0.205/ <b>0.208</b>	0.041/0.240
Synthetic	<b>0.012/10.4</b>	0.023/10.9	0.019/ <b>10.4</b>	0.012/10.4	0.014/10.4
Wave	0.129/1e+06	0.018/6403	<b>0.007/4310</b>	0.089/6496	0.042/5127
Wine	<b>0.017/0.694</b>	0.070/ <b>0.540</b>	0.027/0.575	0.091/0.643	0.061/0.600
Yacht	0.115/4.057	0.174/0.690	0.098/0.508	0.124/ <b>0.371</b>	<b>0.078/0.412</b>



## B.4 Runtime

Tables 12 and 13 provide detailed runtime results for each method. Results are averaged over 10 folds, and standard deviations are shown in subscripted parentheses; lower is better. The last row in each table shows the Geometric mean over all datasets.

Table 12: Total train (including tuning) time (in seconds).

Dataset	NGBoost	PGBM	CBU	IBUG
Ames	<b>417</b> <sub>(587)</sub>	1.4e+05 <sub>(25170)</sub>	23181 <sub>(258)</sub>	22264 <sub>(796)</sub>
Bike	<b>195</b> <sub>(143)</sub>	58246 <sub>(8207)</sub>	23207 <sub>(239)</sub>	22417 <sub>(794)</sub>
California	<b>315</b> <sub>(90.8)</sub>	16958 <sub>(1173)</sub>	23141 <sub>(253)</sub>	22530 <sub>(649)</sub>
Communities	<b>38.0</b> <sub>(21.2)</sub>	24491 <sub>(4246)</sub>	23023 <sub>(260)</sub>	22429 <sub>(483)</sub>
Concrete	<b>57.1</b> <sub>(22.9)</sub>	4130 <sub>(3621)</sub>	22953 <sub>(265)</sub>	22402 <sub>(577)</sub>
Energy	<b>35.3</b> <sub>(33.0)</sub>	2706 <sub>(601)</sub>	22783 <sub>(278)</sub>	22423 <sub>(602)</sub>
Facebook	<b>731</b> <sub>(659)</sub>	3.5e+05 <sub>(58586)</sub>	23061 <sub>(310)</sub>	23145 <sub>(517)</sub>
Kin8nm	<b>77.8</b> <sub>(39.6)</sub>	10489 <sub>(2181)</sub>	23142 <sub>(296)</sub>	22694 <sub>(529)</sub>
Life	<b>105</b> <sub>(87.4)</sub>	83814 <sub>(25313)</sub>	23082 <sub>(273)</sub>	20531 <sub>(7050)</sub>
MEPS	<b>351</b> <sub>(477)</sub>	2.3e+05 <sub>(41039)</sub>	23139 <sub>(327)</sub>	20491 <sub>(7004)</sub>
MSD	<b>11720</b> <sub>(1022)</sub>	2.2e+05 <sub>(34478)</sub>	23972 <sub>(258)</sub>	52760 <sub>(16670)</sub>
Naval	<b>847</b> <sub>(1804)</sub>	38882 <sub>(11481)</sub>	23133 <sub>(210)</sub>	20607 <sub>(7059)</sub>
News	<b>2275</b> <sub>(138)</sub>	2.4e+05 <sub>(60642)</sub>	22960 <sub>(258)</sub>	22492 <sub>(448)</sub>
Obesity	<b>1569</b> <sub>(2208)</sub>	3.2e+05 <sub>(64847)</sub>	23169 <sub>(259)</sub>	21040 <sub>(7086)</sub>
Power	<b>107</b> <sub>(53.0)</sub>	12556 <sub>(1459)</sub>	23042 <sub>(338)</sub>	22445 <sub>(298)</sub>
Protein	<b>1430</b> <sub>(2298)</sub>	40132 <sub>(5779)</sub>	23043 <sub>(277)</sub>	22865 <sub>(344)</sub>
STAR	<b>17.0</b> <sub>(5.788)</sub>	20797 <sub>(3481)</sub>	22852 <sub>(263)</sub>	22074 <sub>(432)</sub>
Superconductor	<b>215</b> <sub>(29.3)</sub>	2.1e+05 <sub>(42859)</sub>	23291 <sub>(706)</sub>	22503 <sub>(426)</sub>
Synthetic	<b>439</b> <sub>(351)</sub>	1.0e+05 <sub>(15729)</sub>	23068 <sub>(333)</sub>	22394 <sub>(532)</sub>
Wave	<b>3487</b> <sub>(342)</sub>	1.0e+05 <sub>(16204)</sub>	23394 <sub>(173)</sub>	44282 <sub>(16994)</sub>
Wine	<b>33.1</b> <sub>(19.2)</sub>	15067 <sub>(2804)</sub>	20942 <sub>(7191)</sub>	22269 <sub>(451)</sub>
Yacht	<b>51.3</b> <sub>(41.8)</sub>	1965 <sub>(87.7)</sub>	22915 <sub>(239)</sub>	22184 <sub>(433)</sub>
Geo. mean	<b>265</b>	43604	23017	23726

Table 13: Average prediction time per text example (in milliseconds).

Dataset	NGBoost	PGBM	CBU	IBUG
Ames	<b>5.583</b> <sub>(5.778)</sub>	9.505 <sub>(2.426)</sub>	<b>0.066</b> <sub>(0.010)</sub>	4.851 <sub>(2.766)</sub>
Bike	<b>0.514</b> <sub>(0.815)</sub>	<b>7.705</b> <sub>(8.198)</sub>	<b>0.010</b> <sub>(0.002)</sub>	61.6 <sub>(21.1)</sub>
California	0.243 <sub>(0.082)</sub>	<b>5.659</b> <sub>(9.562)</sub>	<b>0.004</b> <sub>(0.001)</sub>	23.4 <sub>(5.265)</sub>
Communities	0.393 <sub>(0.170)</sub>	11.5 <sub>(0.941)</sub>	<b>0.027</b> <sub>(0.010)</sub>	1.803 <sub>(1.118)</sub>
Concrete	2.154 <sub>(0.884)</sub>	<b>44.8</b> <sub>(57.5)</sub>	<b>0.043</b> <sub>(0.019)</sub>	1.876 <sub>(0.726)</sub>
Energy	1.830 <sub>(1.369)</sub>	32.9 <sub>(12.1)</sub>	<b>0.053</b> <sub>(0.027)</sub>	1.135 <sub>(0.300)</sub>
Facebook	0.533 <sub>(0.469)</sub>	<b>5.148</b> <sub>(7.166)</sub>	<b>0.024</b> <sub>(0.004)</sub>	105 <sub>(62.4)</sub>
Kin8nm	0.194 <sub>(0.107)</sub>	168 <sub>(89.7)</sub>	<b>0.008</b> <sub>(0.002)</sub>	4.713 <sub>(0.454)</sub>
Life	1.466 <sub>(1.171)</sub>	31.5 <sub>(28.1)</sub>	<b>0.064</b> <sub>(0.037)</sub>	6.376 <sub>(1.275)</sub>
MEPS	0.465 <sub>(0.121)</sub>	<b>9.510</b> <sub>(23.6)</sub>	<b>0.005</b> <sub>(0.002)</sub>	8.845 <sub>(7.866)</sub>
MSD	1.712 <sub>(0.347)</sub>	25.3 <sub>(1.933)</sub>	<b>0.003</b> <sub>(7e-04)</sub>	603 <sub>(97.4)</sub>
Naval	0.280 <sub>(0.129)</sub>	<b>187</b> <sub>(253)</sub>	<b>0.010</b> <sub>(0.007)</sub>	41.9 <sub>(22.3)</sub>
News	0.577 <sub>(0.100)</sub>	0.771 <sub>(0.403)</sub>	<b>0.002</b> <sub>(1e-03)</sub>	40.0 <sub>(38.0)</sub>
Obesity	0.988 <sub>(0.475)</sub>	10.1 <sub>(6.503)</sub>	<b>0.020</b> <sub>(0.003)</sub>	110 <sub>(9.535)</sub>
Power	0.252 <sub>(0.115)</sub>	6.904 <sub>(4.256)</sub>	<b>0.007</b> <sub>(0.002)</sub>	18.3 <sub>(17.5)</sub>
Protein	0.154 <sub>(0.080)</sub>	130 <sub>(73.7)</sub>	<b>0.004</b> <sub>(7e-04)</sub>	90.8 <sub>(25.7)</sub>
STAR	0.353 <sub>(0.121)</sub>	10.0 <sub>(1.191)</sub>	<b>0.038</b> <sub>(0.010)</sub>	0.937 <sub>(0.369)</sub>
Superconductor	0.096 <sub>(0.055)</sub>	<b>27.1</b> <sub>(79.5)</sub>	<b>0.005</b> <sub>(0.002)</sub>	52.6 <sub>(25.8)</sub>
Synthetic	<b>0.482</b> <sub>(0.541)</sub>	2.461 <sub>(0.524)</sub>	<b>0.009</b> <sub>(0.005)</sub>	<b>7.595</b> <sub>(10.1)</sub>
Wave	<b>1.018</b> <sub>(1.339)</sub>	<b>9.116</b> <sub>(17.0)</sub>	<b>0.003</b> <sub>(3e-04)</sub>	719 <sub>(106)</sub>
Wine	0.135 <sub>(0.061)</sub>	<b>280</b> <sub>(281)</sub>	<b>0.009</b> <sub>(0.002)</sub>	6.001 <sub>(1.696)</sub>
Yacht	4.192 <sub>(2.669)</sub>	71.5 <sub>(12.1)</sub>	<b>0.124</b> <sub>(0.081)</sub>	1.237 <sub>(0.334)</sub>
Geo. mean	0.585	18.5	<b>0.013</b>	15.9

## C Additional Experiments

In this section, we present additional experimental results.

### C.1 Probabilistic Performance Without Variance Calibration

Tables 14 and 15 show the probabilistic performance of each method *without* variance calibration. Even without variance calibration, IBUG+CBU generally outperforms competing methods. Standard errors are shown in subscripted parentheses.

Table 14: Probabilistic (CRPS  $\downarrow$ ) performance *without* variance calibration.

Dataset	NGBoost	PGBM	CBU	IBUG	IBUG+CBU
Ames	38279 <sub>(564)</sub>	12173 <sub>(484)</sub>	11948 <sub>(386)</sub>	<b>10442</b> <sub>(373)</sub>	<b>10208</b> <sub>(392)</sub>
Bike	13.9 <sub>(1.856)</sub>	1.274 <sub>(0.054)</sub>	<b>0.835</b> <sub>(0.035)</sub>	1.899 <sub>(0.224)</sub>	1.219 <sub>(0.105)</sub>
California	2e+11 <sub>(2e+11)</sub>	0.227 <sub>(0.004)</sub>	0.221 <sub>(0.001)</sub>	0.213 <sub>(1e-03)</sub>	<b>0.207</b> <sub>(9e-04)</sub>
Communities	0.068 <sub>(0.002)</sub>	0.077 <sub>(0.004)</sub>	0.070 <sub>(0.002)</sub>	<b>0.065</b> <sub>(0.002)</sub>	<b>0.065</b> <sub>(0.002)</sub>
Concrete	3.395 <sub>(0.181)</sub>	1.932 <sub>(0.088)</sub>	1.994 <sub>(0.095)</sub>	1.938 <sub>(0.079)</sub>	<b>1.780</b> <sub>(0.085)</sub>
Energy	0.539 <sub>(0.042)</sub>	<b>0.151</b> <sub>(0.007)</sub>	0.207 <sub>(0.010)</sub>	0.481 <sub>(0.041)</sub>	0.293 <sub>(0.022)</sub>
Facebook	4.022 <sub>(0.099)</sub>	3.860 <sub>(0.149)</sub>	3.214 <sub>(0.058)</sub>	3.072 <sub>(0.066)</sub>	<b>2.971</b> <sub>(0.072)</sub>
Kin8nm	0.095 <sub>(6e-04)</sub>	0.069 <sub>(0.003)</sub>	0.063 <sub>(8e-04)</sub>	0.052 <sub>(4e-04)</sub>	<b>0.052</b> <sub>(6e-04)</sub>
Life	2.897 <sub>(1.465)</sub>	0.836 <sub>(0.035)</sub>	0.852 <sub>(0.030)</sub>	0.794 <sub>(0.022)</sub>	<b>0.739</b> <sub>(0.024)</sub>
MEPS	<b>5.529</b> <sub>(0.196)</sub>	6.725 <sub>(0.126)</sub>	6.050 <sub>(0.109)</sub>	6.146 <sub>(0.113)</sub>	6.022 <sub>(0.114)</sub>
MSD	4.525 <sub>(0.005)</sub>	5.767 <sub>(0.006)</sub>	4.364 <sub>(0.004)</sub>	4.410 <sub>(0.005)</sub>	<b>4.342</b> <sub>(0.004)</sub>
Naval	0.003 <sub>(6e-05)</sub>	0.005 <sub>(0.002)</sub>	3e-04 <sub>(3e-06)</sub>	3e-04 <sub>(2e-06)</sub>	<b>3e-04</b> <sub>(2e-06)</sub>
News	<b>2476</b> <sub>(38.9)</sub>	2628 <sub>(94.9)</sub>	2712 <sub>(59.7)</sub>	2669 <sub>(43.9)</sub>	2593 <sub>(49.7)</sub>
Obesity	3.208 <sub>(0.028)</sub>	1.860 <sub>(0.022)</sub>	<b>1.754</b> <sub>(0.017)</sub>	1.882 <sub>(0.019)</sub>	1.772 <sub>(0.018)</sub>
Power	2.104 <sub>(0.024)</sub>	1.585 <sub>(0.057)</sub>	1.572 <sub>(0.024)</sub>	1.542 <sub>(0.020)</sub>	<b>1.488</b> <sub>(0.022)</sub>
Protein	5427 <sub>(5409)</sub>	1.932 <sub>(0.014)</sub>	1.822 <sub>(0.010)</sub>	1.785 <sub>(0.008)</sub>	<b>1.740</b> <sub>(0.009)</sub>
STAR	132 <sub>(1.697)</sub>	157 <sub>(6.908)</sub>	132 <sub>(1.540)</sub>	<b>129</b> <sub>(1.225)</sub>	<b>130</b> <sub>(1.327)</sub>
Superconductor	3.200 <sub>(0.031)</sub>	<b>0.134</b> <sub>(0.005)</sub>	0.151 <sub>(0.004)</sub>	0.303 <sub>(0.025)</sub>	0.201 <sub>(0.013)</sub>
Synthetic	5.778 <sub>(0.043)</sub>	6.946 <sub>(0.242)</sub>	<b>5.769</b> <sub>(0.049)</sub>	<b>5.731</b> <sub>(0.040)</sub>	<b>5.735</b> <sub>(0.042)</sub>
Wave	5.7e+05 <sub>(886)</sub>	4152 <sub>(247)</sub>	<b>2350</b> <sub>(10.3)</sub>	4905 <sub>(12.3)</sub>	3112 <sub>(8.952)</sub>
Wine	0.385 <sub>(0.005)</sub>	0.383 <sub>(0.015)</sub>	0.355 <sub>(0.007)</sub>	<b>0.322</b> <sub>(0.006)</sub>	<b>0.321</b> <sub>(0.006)</sub>
Yacht	1.187 <sub>(0.142)</sub>	<b>0.310</b> <sub>(0.056)</sub>	<b>0.291</b> <sub>(0.050)</sub>	0.644 <sub>(0.068)</sub>	0.394 <sub>(0.042)</sub>
IBUG W-T-L	16-4-2	12-4-6	10-4-8	-	0-5-17
IBUG+CBU W-T-L	17-3-2	15-4-3	14-2-6	17-5-0	-

Table 15: Probabilistic (NLL  $\downarrow$ ) performance *without* variance calibration.

Dataset	NGBoost	PGBM	CBU	IBUG	IBUG+CBU
Ames	11.3 <sub>(0.023)</sub>	23.7 <sub>(6.813)</sub>	1676 <sub>(1093)</sub>	<b>11.2</b> <sub>(0.031)</sub>	<b>11.3</b> <sub>(0.070)</sub>
Bike	1.942 <sub>(0.024)</sub>	11.1 <sub>(4.073)</sub>	<b>1.264</b> <sub>(0.080)</sub>	2.958 <sub>(0.106)</sub>	2.386 <sub>(0.091)</sub>
California	0.551 <sub>(0.010)</sub>	7.821 <sub>(7.026)</sub>	2.261 <sub>(0.341)</sub>	0.484 <sub>(0.009)</sub>	<b>0.375</b> <sub>(0.009)</sub>
Communities	<b>-7e-01</b> <sub>(0.056)</sub>	20.7 <sub>(7.235)</sub>	2.438 <sub>(1.168)</sub>	<b>-6e-01</b> <sub>(0.136)</sub>	<b>-4e-01</b> <sub>(0.269)</sub>
Concrete	3.062 <sub>(0.031)</sub>	3.102 <sub>(0.277)</sub>	684 <sub>(358)</sub>	<b>2.848</b> <sub>(0.055)</sub>	<b>2.822</b> <sub>(0.093)</sub>
Energy	<b>0.670</b> <sub>(0.250)</sub>	<b>0.481</b> <sub>(0.341)</sub>	6.129 <sub>(3.597)</sub>	1.461 <sub>(0.113)</sub>	1.048 <sub>(0.162)</sub>
Facebook	2.099 <sub>(0.026)</sub>	14.7 <sub>(5.523)</sub>	5.147 <sub>(1.834)</sub>	2.195 <sub>(0.070)</sub>	<b>2.044</b> <sub>(0.045)</sub>
Kin8nm	-4e-01 <sub>(0.007)</sub>	35.0 <sub>(23.1)</sub>	59.5 <sub>(24.2)</sub>	-8e-01 <sub>(0.009)</sub>	<b>-9e-01</b> <sub>(0.024)</sub>
Life	2.188 <sub>(0.044)</sub>	23.5 <sub>(20.6)</sub>	71.9 <sub>(38.8)</sub>	<b>1.889</b> <sub>(0.038)</sub>	<b>1.885</b> <sub>(0.100)</sub>
MEPS	<b>3.732</b> <sub>(0.056)</sub>	11.3 <sub>(3.246)</sub>	<b>3.722</b> <sub>(0.044)</sub>	3.820 <sub>(0.064)</sub>	<b>3.678</b> <sub>(0.054)</sub>
MSD	3.454 <sub>(0.002)</sub>	65.6 <sub>(0.162)</sub>	3.450 <sub>(0.004)</sub>	3.415 <sub>(0.002)</sub>	<b>3.383</b> <sub>(0.002)</sub>
Naval	-5e+00 <sub>(0.007)</sub>	-4e+00 <sub>(0.357)</sub>	-5e+00 <sub>(0.057)</sub>	-6e+00 <sub>(0.007)</sub>	<b>-6e+00</b> <sub>(0.006)</sub>
News	<b>10.9</b> <sub>(0.335)</sub>	130 <sub>(49.5)</sub>	<b>10.8</b> <sub>(0.368)</sub>	<b>11.0</b> <sub>(0.415)</sub>	<b>10.7</b> <sub>(0.307)</sub>
Obesity	2.940 <sub>(0.003)</sub>	2.603 <sub>(0.015)</sub>	<b>2.488</b> <sub>(0.009)</sub>	2.646 <sub>(0.009)</sub>	<b>2.493</b> <sub>(0.009)</sub>
Power	2.769 <sub>(0.042)</sub>	11.0 <sub>(8.270)</sub>	5.304 <sub>(0.672)</sub>	<b>2.575</b> <sub>(0.036)</sub>	<b>2.569</b> <sub>(0.057)</sub>
Protein	2.841 <sub>(0.015)</sub>	5.299 <sub>(0.268)</sub>	3.291 <sub>(0.042)</sub>	2.747 <sub>(0.123)</sub>	<b>2.531</b> <sub>(0.028)</sub>
STAR	6.872 <sub>(0.015)</sub>	23.2 <sub>(5.203)</sub>	6.989 <sub>(0.040)</sub>	<b>6.853</b> <sub>(0.008)</sub>	<b>6.857</b> <sub>(0.012)</sub>
Superconductor	<b>12.1</b> <sub>(13.4)</sub>	10.5 <sub>(4.631)</sub>	<b>-6e-03</b> <sub>(0.093)</sub>	1.151 <sub>(0.111)</sub>	0.602 <sub>(0.094)</sub>
Synthetic	3.746 <sub>(0.008)</sub>	27.3 <sub>(6.288)</sub>	3.782 <sub>(0.032)</sub>	<b>3.738</b> <sub>(0.007)</sub>	<b>3.744</b> <sub>(0.010)</sub>
Wave	10.7 <sub>(0.002)</sub>	22.2 <sub>(7.985)</sub>	<b>9.679</b> <sub>(0.004)</sub>	10.9 <sub>(0.004)</sub>	10.4 <sub>(0.004)</sub>
Wine	1.029 <sub>(0.014)</sub>	109 <sub>(24.9)</sub>	578 <sub>(428)</sub>	<b>0.910</b> <sub>(0.016)</sub>	0.968 <sub>(0.030)</sub>
Yacht	<b>0.904</b> <sub>(0.232)</sub>	7.227 <sub>(4.096)</sub>	4.770 <sub>(2.330)</sub>	1.502 <sub>(0.308)</sub>	<b>1.204</b> <sub>(0.519)</sub>
IBUG W-T-L	11-7-4	10-10-2	8-9-5	-	1-10-11
IBUG+CBU W-T-L	13-7-2	11-10-1	8-11-3	11-10-1	-

## C.2 Comparison to $k$ -Nearest Neighbors

In this section, we compare IBUG to  $k$ -nearest neighbors, in which similarity is defined by Euclidean distance. For the nearest-neighbors approach, we tune two different  $k$  values, one for estimating the conditional mean, and one for estimating the variance. We also apply standard scaling to the data before training, and denote this method  $k$ NN in our results. Table 16 shows that IBUG is consistently better than  $k$ NN in terms of probabilistic performance. However, we note that point predictions from GBRTs is typically better than  $k$ NNs, thus we also compare IBUG to a variant of  $k$ NN that uses CatBoost as a base model to estimate the conditional mean.

**Euclidean Distance vs. Affinity.** To test which similarity measure (Euclidean distance or affinity) is more effective, we use the output from CatBoost to model the conditional mean, we then use  $k$ NN or IBUG to find their respective  $k$ -nearest training examples to estimate the variance; we denote these methods  $k$ NN-CB<sup>4</sup> and IBUG-CB. For  $k$ NN-CB, we also reduce the dimensionality of the data by only using the most important features identified by the CatBoost model;<sup>5</sup> this helps  $k$ NN-CB combat the curse of dimensionality when computing similarity. Results of this comparison are in Table 16, in which we observe IBUG-CB is always statistically the same or better than  $k$ NN-CB. These results suggest affinity is a more effective similarity measure than Euclidean distance for uncertainty estimation in GBRTs.

Table 16: Probabilistic (CRPS) performance comparison of IBUG against two different nearest-neighbor models.  $k$ NN estimates the conditional mean and variance using two different  $k$  values; and  $k$ NN-CB estimates the variance in the same way as  $k$ NN, but uses the scalar output from the CatBoost model to estimate the conditional mean. Overall, these results suggest affinity is a better measure of similarity than Euclidean distance for uncertainty estimation in GBRTs.

Dataset	$k$ NN	$k$ NN-CB	IBUG-CB
Bike	<b>0.932</b> <sub>(0.029)</sub>	<b>0.978</b> <sub>(0.049)</sub>	<b>0.974</b> <sub>(0.048)</sub>
California	<b>0.579</b> <sub>(0.001)</sub>	<b>0.219</b> <sub>(1e-03)</sub>	<b>0.213</b> <sub>(9e-04)</sub>
Communities	<b>0.072</b> <sub>(0.002)</sub>	<b>0.065</b> <sub>(0.002)</sub>	<b>0.065</b> <sub>(0.002)</sub>
Concrete	4.645 <sub>(0.140)</sub>	<b>1.872</b> <sub>(0.085)</sub>	<b>1.849</b> <sub>(0.098)</sub>
Energy	0.875 <sub>(0.016)</sub>	<b>0.153</b> <sub>(0.010)</sub>	<b>0.143</b> <sub>(0.009)</sub>
Facebook	5.613 <sub>(0.065)</sub>	3.275 <sub>(0.068)</sub>	<b>3.073</b> <sub>(0.066)</sub>
Kin8nm	0.067 <sub>(7e-04)</sub>	<b>0.051</b> <sub>(5e-04)</sub>	<b>0.051</b> <sub>(6e-04)</sub>
Life	4.738 <sub>(0.078)</sub>	<b>0.785</b> <sub>(0.024)</sub>	<b>0.794</b> <sub>(0.023)</sub>
MEPS	7.283 <sub>(0.220)</sub>	<b>6.181</b> <sub>(0.107)</sub>	<b>6.150</b> <sub>(0.114)</sub>
MSD	5.312 <sub>(0.006)</sub>	4.446 <sub>(0.004)</sub>	<b>4.410</b> <sub>(0.005)</sub>
Naval	8e-04 <sub>(2e-05)</sub>	3e-04 <sub>(2e-06)</sub>	<b>2e-04</b> <sub>(2e-06)</sub>
News	2654 <sub>(52.0)</sub>	<b>2597</b> <sub>(52.3)</sub>	<b>2545</b> <sub>(41.0)</sub>
Obesity	5.526 <sub>(0.013)</sub>	<b>1.900</b> <sub>(0.043)</sub>	<b>1.866</b> <sub>(0.021)</sub>
Power	2.074 <sub>(0.022)</sub>	<b>1.553</b> <sub>(0.020)</sub>	<b>1.542</b> <sub>(0.020)</sub>
Protein	3.241 <sub>(0.010)</sub>	<b>1.787</b> <sub>(0.008)</sub>	<b>1.784</b> <sub>(0.008)</sub>
STAR	140 <sub>(1.553)</sub>	<b>129</b> <sub>(1.204)</sub>	<b>130</b> <sub>(1.214)</sub>
Superconductor	3.445 <sub>(0.041)</sub>	<b>0.156</b> <sub>(0.006)</sub>	<b>0.153</b> <sub>(0.006)</sub>
Synthetic	6.136 <sub>(0.047)</sub>	<b>5.735</b> <sub>(0.039)</sub>	<b>5.731</b> <sub>(0.040)</sub>
Wave	11987 <sub>(36.6)</sub>	2700 <sub>(17.0)</sub>	<b>2679</b> <sub>(16.0)</sub>
Wine	0.445 <sub>(0.004)</sub>	<b>0.322</b> <sub>(0.006)</sub>	<b>0.322</b> <sub>(0.006)</sub>
Yacht	3.354 <sub>(0.408)</sub>	<b>0.275</b> <sub>(0.048)</sub>	<b>0.276</b> <sub>(0.048)</sub>
IBUG-CB W-T-L	21-1-0	10-12-0	-

<sup>4</sup>Again, we apply standard scaling to the data before training  $k$ NN-CB.

<sup>5</sup>We tune the number of important features to use for  $k$ NN-CB using values [5, 10, 20].

### C.3 Comparison to Bayesian Additive Regression Trees

BART (Bayesian Additive Regression Trees) takes a Bayesian approach to uncertainty estimation in trees [13]. Although well-grounded theoretically, BART requires expensive sampling techniques such as MCMC (Markov Chain Monte Carlo) to provide approximate solutions.

In this section, we compare IBUG to BART using a popular open-source implementation.<sup>6</sup> However, due to BART’s computational complexity, we tune the number of trees for both IBUG and BART using values [10, 50, 100, 200], set the number of chains for BART to 5, and run our comparison using a subset of the datasets in our empirical evaluation consisting of 11 relatively small datasets.

Tables 17 and 18 show IBUG consistently outperforms BART in terms of both probabilistic and point performance.

Table 17: Probabilistic (CRPS ↓) performance comparison between IBUG and BART.

Dataset	BART	IBUG
Bike	4.521 <sub>(0.119)</sub>	<b>0.974</b> <sub>(0.048)</sub>
California	0.285 <sub>(0.001)</sub>	<b>0.213</b> <sub>(9e-04)</sub>
Communities	0.072 <sub>(0.002)</sub>	<b>0.065</b> <sub>(0.002)</sub>
Concrete	3.067 <sub>(0.073)</sub>	<b>1.849</b> <sub>(0.098)</sub>
Energy	0.402 <sub>(0.023)</sub>	<b>0.143</b> <sub>(0.009)</sub>
Kin8nm	0.107 <sub>(8e-04)</sub>	<b>0.051</b> <sub>(6e-04)</sub>
Naval	1e-03 <sub>(1e-05)</sub>	<b>2e-04</b> <sub>(2e-06)</sub>
Power	2.225 <sub>(0.018)</sub>	<b>1.542</b> <sub>(0.020)</sub>
STAR	134 <sub>(1.493)</sub>	<b>130</b> <sub>(1.214)</sub>
Wine	0.394 <sub>(0.005)</sub>	<b>0.322</b> <sub>(0.006)</sub>
Yacht	0.849 <sub>(0.039)</sub>	<b>0.276</b> <sub>(0.048)</sub>
IBUG W-T-L	11-0-0	-

Table 18: Point (RMSE ↓) performance comparison between IBUG and BART.

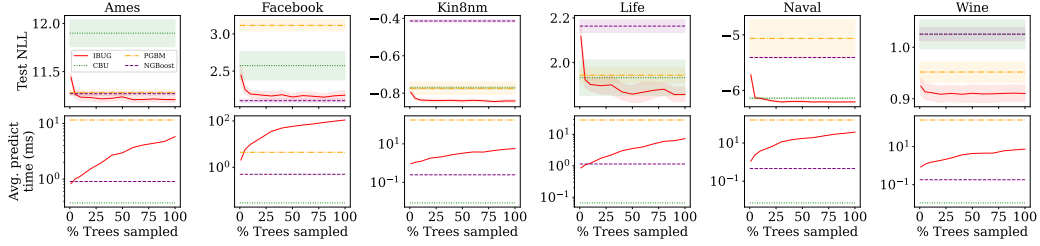
Dataset	BART	IBUG
Bike	8.396 <sub>(0.273)</sub>	<b>2.826</b> <sub>(0.200)</sub>
California	0.547 <sub>(0.003)</sub>	<b>0.432</b> <sub>(0.001)</sub>
Communities	0.137 <sub>(0.004)</sub>	<b>0.133</b> <sub>(0.004)</sub>
Concrete	5.507 <sub>(0.161)</sub>	<b>3.629</b> <sub>(0.183)</sub>
Energy	0.685 <sub>(0.039)</sub>	<b>0.264</b> <sub>(0.023)</sub>
Kin8nm	0.186 <sub>(0.001)</sub>	<b>0.086</b> <sub>(8e-04)</sub>
Naval	0.002 <sub>(2e-05)</sub>	<b>5e-04</b> <sub>(5e-06)</sub>
Power	4.057 <sub>(0.049)</sub>	<b>2.941</b> <sub>(0.059)</sub>
STAR	234 <sub>(2.479)</sub>	<b>228</b> <sub>(1.985)</sub>
Wine	0.708 <sub>(0.008)</sub>	<b>0.596</b> <sub>(0.012)</sub>
Yacht	1.624 <sub>(0.121)</sub>	<b>0.668</b> <sub>(0.125)</sub>
IBUG W-T-L	11-0-0	-

<sup>6</sup><https://github.com/JakeColtman/bartpy>

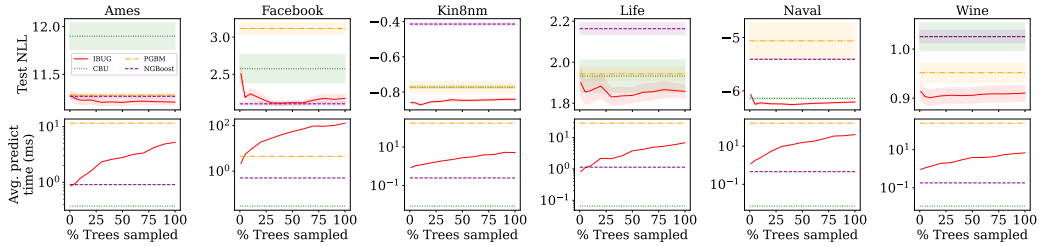
### C.4 Different Tree-Sampling Strategies

Figure 5 shows the probabilistic (NLL) performance of IBUG as the number of trees sampled ( $\tau$ ) increases using three different sampling strategies: *uniformly at random*, *first-to-last*, and *last-to-first*.

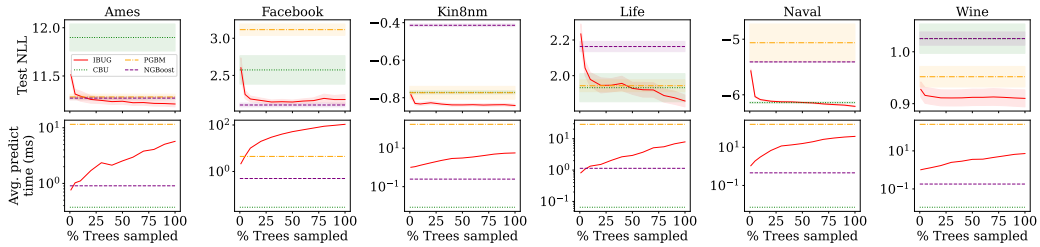
We observe that sampling trees *last-to-first* often requires sampling all trees in order to achieve the lowest NLL on the test set. When sampling *uniformly at random*, NLL tends to plateau starting around 10%. In contrast, sampling trees *first-to-last* on the Kin8nm, Naval, and Wine datasets requires 5% of the trees or less to result in the same or better NLL than when sampling all trees; these results provide some evidence that trees early in training contribute most, and suggest that sampling trees first-to-last may be most effective at obtaining the best probabilistic performance while sampling the fewest number of trees.



(a) Sampling trees *uniformly at random*.



(b) Sampling trees *first-to-last*.



(c) Sampling trees *last-to-first*.

Figure 5: Probabilistic (NLL) performance (lower is better) and average prediction time (in milliseconds) per test example (lower is better) as a function of  $\tau$  for different sampling techniques. *Top*: sample trees uniformly at random, *middle*: sample trees first-to-last (in terms of boosting iteration), *bottom*: sample trees last-to-first. All methods result in similar prediction times; however, *first-to-last* sampling typically provides the best NLL with the fewest number of trees sampled.

### C.5 Leaf Density

Figure 6 shows the average percentage of train instances visited per tree as a function of the total number of training instances for each dataset. We note that for some datasets, CatBoost, LightGBM, and XGBoost induce regression trees with very dense leaves where over half the training instances belong to those leaves. Figure 7 shows average leaf density for each tree in the GBRT.

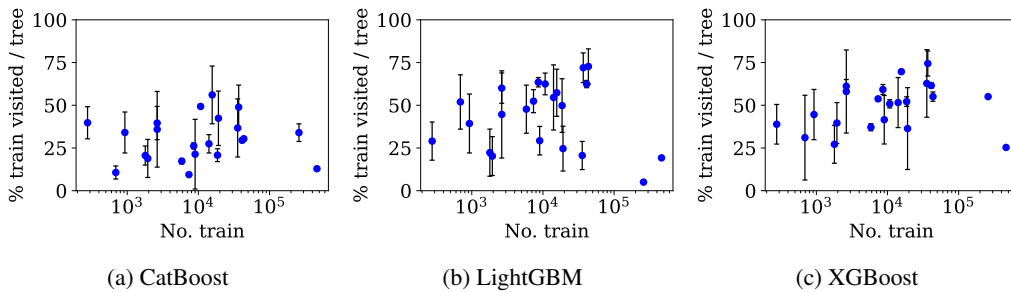


Figure 6: Average percentage of training instances visited per tree while computing affinity vectors on the test set test for each dataset. Results are averaged over all test instances, and error bars represent standard deviation; lower is better. In general, the number of training instances visited per tree is highly dependent on the dataset; and for some datasets, is also highly dependent on the test example (points with large standard deviations).

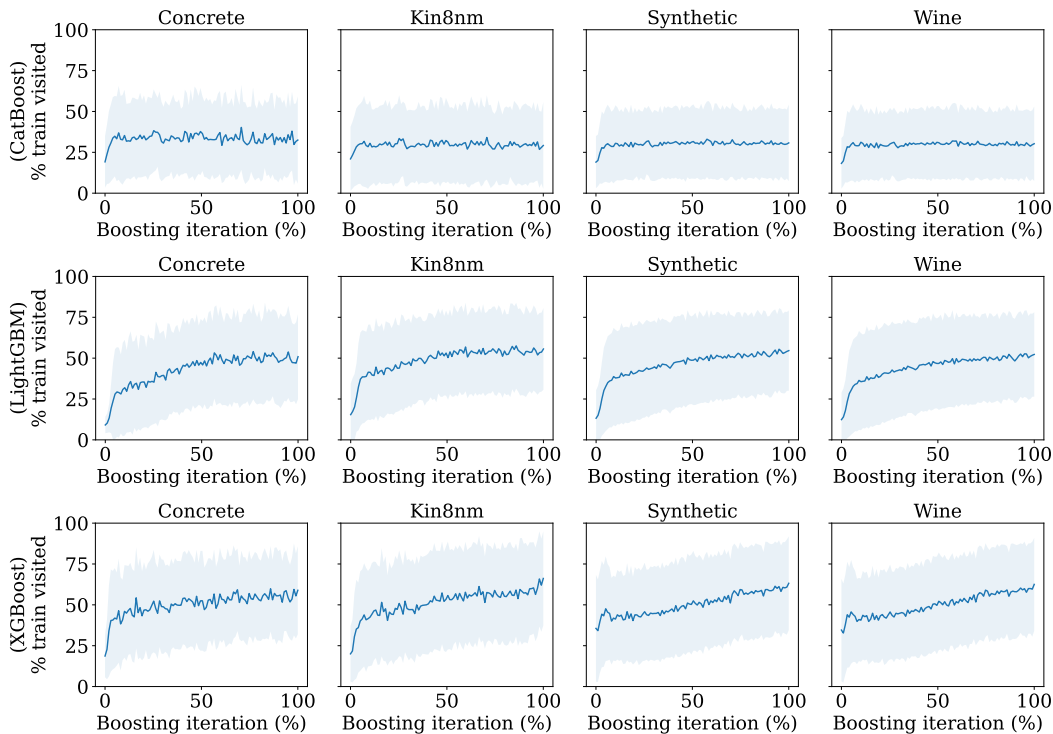


Figure 7: Average percentage of training instances visited at each iteration while computing affinity vectors on the test set for the Concrete, Kin8nm, Synthetic, and Wine datasets. Results are averaged over all test instances, and error bars represent standard deviation; lower is better. For LightGBM and XGBoost, weak learners later in training tend to pool a larger proportion of training instances into fewer leaves; in contrast, CatBoost has less dense leaves and training instances are more equally distributed among the leaves in each tree.