

Probabilistic Gradient Boosting Machines for Large-Scale Probabilistic Regression

Olivier Sprangers
AIRLab, University of Amsterdam
Amsterdam, The Netherlands
o.r.sprangers@uva.nl

Sebastian Schelter
University of Amsterdam
Amsterdam, The Netherlands
s.schelter@uva.nl

Maarten de Rijke
University of Amsterdam
Amsterdam, The Netherlands
derijke@uva.nl

ABSTRACT

Gradient Boosting Machines (GBMs) are hugely popular for solving tabular data problems. However, practitioners are not only interested in point predictions, but also in probabilistic predictions in order to quantify the uncertainty of the predictions. Creating such probabilistic predictions is difficult with existing GBM-based solutions: they either require training multiple models or they become too computationally expensive to be useful for large-scale settings.

We propose Probabilistic Gradient Boosting Machines (PGBMs), a method to create probabilistic predictions with a single ensemble of decision trees in a computationally efficient manner. PGBM approximates the leaf weights in a decision tree as a random variable, and approximates the mean and variance of each sample in a dataset via stochastic tree ensemble update equations. These learned moments allow us to subsequently sample from a specified distribution after training.

We empirically demonstrate the advantages of PGBM compared to existing state-of-the-art methods: (i) PGBM enables probabilistic estimates without compromising on point performance in a single model, (ii) PGBM learns probabilistic estimates via a single model only (and without requiring multi-parameter boosting), and thereby offers a speedup of up to several orders of magnitude over existing state-of-the-art methods on large datasets, and (iii) PGBM achieves accurate probabilistic estimates in tasks with complex differentiable loss functions, such as hierarchical time series problems, where we observed up to 10% improvement in point forecasting performance and up to 300% improvement in probabilistic forecasting performance.

ACM Reference Format:

Olivier Sprangers, Sebastian Schelter, and Maarten de Rijke. 2021. Probabilistic Gradient Boosting Machines for Large-Scale Probabilistic Regression. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '21)*, August 14–18, 2021, Virtual Event, Singapore. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3447548.3467278>

1 INTRODUCTION

Forecasting practitioners are increasingly interested in probabilistic forecasts instead of point forecasts, in order to obtain a notion of

uncertainty of the forecast [2]. Even though probabilistic forecasting techniques have been around for quite some time, applying these techniques in large-scale industrial settings often remains challenging. For example, retailers may require thousands of new forecasts for each product for each store. In this setting, traditional confidence interval techniques based on single model estimates are often computationally too expensive to execute on a daily basis [14]. Existing Gradient Boosting methods for probabilistic forecasting often require training multiple models (e.g., LightGBM [10] or xgboost [4] require a separate model for each quantile of the forecast), or require computing expensive second-order derivative statistics [NGBoost, 6]. Therefore, we aim to find a method that efficiently generates high-quality probabilistic forecasts with a single model using GBMs.

We address the challenge of large-scale probabilistic forecasting by proposing a novel, simple probabilistic forecasting method that leverages the popular Gradient Boosting paradigm to provide accurate probabilistic forecasts in large-scale data settings (Section 3). We demonstrate that our approach achieves state-of-the-art point performance as well as probabilistic performance in forecasting tasks while only training a single ensemble of Gradient Boosted Decision Trees (GBDT). Our proposed method, *Probabilistic Gradient Boosting Machine* (PGBM), consists of two steps: (1) We treat the leaf weights in each tree as random variables that we approximate during training via the sample mean and sample variance of the samples in each leaf, and (2) We obtain an accurate estimate of the conditional mean and variance of our target for each sample by sequentially adding these random variables for each new tree. After training, we obtain a learned mean and variance for each sample which can be used during prediction. Based on the learned mean and variance, we can specify a distribution from which to obtain our probabilistic forecast after training. Our PGBM is simple as its learning procedure is comparable to standard gradient boosting such as LightGBM [10] or xgboost [4] while it requires only few additional computation steps during training and prediction. However, contrary to these existing methods, our method only requires training a single ensemble of decision trees to obtain a model capable of providing probabilistic predictions.

We empirically demonstrate that PGBM offers state-of-the-art point and probabilistic regression performance on 11 datasets of various sizes (Section 4.1). Therefore, PGBM provides the advantages of existing state-of-the-art point prediction gradient boosting packages such as LightGBM or xgboost, as well as the advantage of a state-of-the-art probabilistic prediction package such as NGBoost. In addition, we show how to optimise PGBM's probabilistic estimates *after* training by varying a single hyperparameter and choosing different sets of posterior distributions. This offers the

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
KDD '21, August 14–18, 2021, Virtual Event, Singapore.

© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-8332-5/21/08...\$15.00
<https://doi.org/10.1145/3447548.3467278>

benefit of training a model only once, and optimizing for probabilistic performance thereafter. Neither existing standard gradient boosting packages or probabilistic packages such as NGBoost offer this. Furthermore, we demonstrate that our GPU-based implementation of PGBM trains an order of magnitude faster than existing state-of-the-art probabilistic gradient boosting methods on large datasets. Finally, we showcase the use of PGBM on the problem of probabilistic hierarchical time series forecasting, demonstrating that our implementation enables the optimization of complex differentiable loss functions without manually specifying an analytical gradient and hessian (in contrast to existing gradient boosting packages that rely on a priori specification of an analytical gradient and hessian).

In summary, the contributions of this paper are the following:

- We introduce PGBM, a gradient boosting framework for probabilistic regression problems (Section 3);
- We demonstrate state-of-the-art point performance and probabilistic performance of PGBM on a set of regression benchmarks (Section 4.1);
- We show that PGBM’s probabilistic performance can be optimized *after* training the model, which allows practitioners to choose different posterior distributions without needing to re-train the model (Section 4.1);
- Our implementation of PGBM trains up to several orders of magnitude faster on larger datasets than competing methods (Section 4.1), and our implementation allows the use of complex differentiable loss functions, where we observed up to 10% improvement in point forecasting performance and up to 300% improvement in probabilistic forecasting performance (Section 4.2).

2 BACKGROUND

Gradient boosting optimizes a loss function by iteratively adding a set of weak learners into an ensemble [7]. Each new weak learner is added sequentially, such that this new learner reduces the aggregate error from the existing ensemble of weak learners. Typically the weak learners are decision trees due to their strong empirical performance; and we also choose them as base learners in this work following the formalization of gradient boosting with decision trees due to Chen and Guestrin [4]. In gradient boosting, at each iteration k , we seek to construct a decision tree $f^{(k)}(\mathbf{x}_i)$ such that the update equation for our estimate for sample i reads:

$$\hat{y}_i^{(k)} = \hat{y}_i^{(k-1)} - \alpha f^{(k)}(\mathbf{x}_i), \quad (1)$$

in which α denotes the learning rate, typically chosen to be less than 1, such that only a tiny portion of each new base learner is added to the overall estimate at each iteration. We will derive a different set of update equations for PGBM in Section 3.3. To construct the decision tree $f^{(k)}$, we greedily split our training data based on its input features \mathbf{x} into left (I_L) and right (I_R) nodes by maximizing the following *gain*:

$$G = \frac{1}{2} \left[\frac{(\sum_{i \in I_L} g_i)^2}{\sum_{i \in I_L} h_i + \lambda} + \frac{(\sum_{i \in I_R} g_i)^2}{\sum_{i \in I_R} h_i + \lambda} - \frac{(\sum_{i \in I} g_i)^2}{\sum_{i \in I} h_i + \lambda} \right], \quad (2)$$

where λ is a regularization parameter, $I = I_L \cup I_R$, and g_i, h_i are the gradient and hessian, respectively, with respect to $\hat{y}_i^{(k-1)}$ of some differentiable loss function that we aim to minimize, for example

the mean-squared error loss in case of a regression problem.^{1,2} When constructing the decision tree, Eq. (2) is evaluated at each node to find the best possible split gain G^* among all features in the input \mathbf{x} , and typically a split is made if the gain exceeds a certain threshold. If no split is made, the node becomes a *leaf* and the corresponding leaf weight w_j follows from:

$$w_j = -\frac{\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda}, \quad (3)$$

in which $j \in \{0, 1, \dots, T\}$, with T the total number of leaves in our tree. We typically stop learning the tree if some pre-defined criterion is met, for example if no more splits with a positive split gain according to Eq. (2) can be made or the tree reaches a pre-defined fixed number of leaves. After training the tree, the output $f^{(k)}(\mathbf{x}_i)$ for a particular sample is then simply the leaf weight w_j , or:

$$\hat{y}_i^{(k)} = \hat{y}_i^{(k-1)} - \alpha w_j. \quad (4)$$

3 PROBABILISTIC GRADIENT BOOSTING MACHINES

We discuss our method *Probabilistic Gradient Boosting Machines* (PGBM). We introduce our problem setting, the core components of PGBM, and end with an analysis and discussion of PGBM.

3.1 Probabilistic Forecasting

In this work, we are interested in the problem of probabilistic forecasting. More generally, we are interested in the problem of probabilistic regression, in which one aims to estimate a conditional probability distribution $P(y|\mathbf{x})$ of some target scalar variable y based on a set of inputs \mathbf{x} . In the case of probabilistic forecasting, \mathbf{x} commonly includes lagged target variables as well as additional covariates. We are interested in finding a model $f(\mathbf{x})$ that provides us with an estimate of the mean μ and variance σ^2 of a target distribution such that we can obtain a sample of an estimate \hat{y} by sampling from a specified distribution D after training:

$$(\mu_{\hat{y}}, \sigma_{\hat{y}}^2) = f(\mathbf{x}) \quad (5)$$

$$\hat{y} \sim D(\mu_{\hat{y}}, \sigma_{\hat{y}}^2). \quad (6)$$

We construct our model $f(\mathbf{x})$ using gradient boosting. In order to find the estimate for the mean and variance of our target distribution, we next derive formulas for stochastic leaf weights (Eq. (3)) and new update equations (Eq. (4)).

3.2 Stochastic Leaf Weights

By creating stochastic leaf weights, we are able to learn a mean and variance of each leaf weight in each tree, thus enabling us to learn a mean and variance for each sample in our dataset. We assume that the gradient and hessian of our loss function are random variables with a mean (μ_g, μ_h) and finite variance (σ_g^2, σ_h^2) and covariance σ_{gh}^2 , which we approximate separately in each tree for each instance set

¹Some researchers refer to this method as *Newton boosting* rather than *gradient boosting* [20], as it employs a second-order derivative.

²Unlike [4], we drop the regularization term γT . We have no need for this regularization parameter, as the number of leaves T is a hyperparameter that needs to be specified in our method. Note also that the parameter γ is in fact by default set to zero in the xgboost implementation of [4].

I_j using the sample mean, sample variance and sample covariance for the n_j samples in an instance set I_j :

$$\mu_g \approx \bar{g}_j = \frac{1}{n_j} \sum_{i \in I_j} g_i \quad (7)$$

$$\mu_h \approx \bar{h}_j = \frac{1}{n_j} \sum_{i \in I_j} h_i \quad (8)$$

$$\sigma_g^2 \approx \bar{\sigma}_{g_j}^2 = \frac{1}{n_j - 1} \sum_{i \in I_j} (g_i - \bar{g})^2 \quad (9)$$

$$\sigma_h^2 \approx \bar{\sigma}_{h_j}^2 = \frac{1}{n_j - 1} \sum_{i \in I_j} (h_i - \bar{h})^2 \quad (10)$$

$$\sigma_{gh}^2 \approx \bar{\sigma}_{gh_j}^2 = \frac{1}{n_j - 1} \sum_{i \in I_j} (g_i - \bar{g})(h_i - \bar{h}). \quad (11)$$

Note that the sample variance and covariance require the Bessel correction $n_j - 1$ in order to obtain an unbiased estimate of the true variance and covariance. Moreover, the Central Limit Theorem dictates that we obtain our true mean and variance if $n_j \rightarrow \infty$. However, for tiny datasets, n_j is typically a small number as each leaf may only contain a few samples and thereby using sample statistics might be inappropriate. In the Experiments section (Section 4), we will demonstrate that we are still able to provide accurate probabilistic estimates even in such cases. Next, we write Eq. (3) in terms of the sample mean of the gradient and hessian:

$$w_j = -\frac{\frac{1}{n_j} \sum_{i \in I_j} g_i}{\frac{1}{n_j} \sum_{i \in I_j} h_i + \frac{1}{n_j} \lambda} = -\frac{\bar{g}_j}{\bar{h}_j + \bar{\lambda}_j}. \quad (12)$$

Now, we can model the expectation and variance of the leaf weight w_j using the sample statistics as follows (dropping the subscript j for readability):

$$\mu_j = E \left[\frac{\bar{g}}{(\bar{h} + \bar{\lambda})} \right] \approx \frac{\bar{g}}{(\bar{h} + \bar{\lambda})} - \frac{\bar{\sigma}_{gh}^2}{(\bar{h} + \bar{\lambda})^2} + \frac{\bar{g}\bar{\sigma}_h^2}{(\bar{h} + \bar{\lambda})^3} \quad (13)$$

$$\sigma_j^2 = V \left[\frac{\bar{g}}{(\bar{h} + \bar{\lambda})} \right] \approx \frac{\bar{\sigma}_g^2}{(\bar{h} + \bar{\lambda})^2} + \frac{\bar{g}^2 \bar{\sigma}_h^2}{(\bar{h} + \bar{\lambda})^4} - 2 \frac{\bar{g}\bar{\sigma}_{gh}^2}{(\bar{h} + \bar{\lambda})^3}. \quad (14)$$

We refer the reader to the Supplemental Materials A for the derivation of Eq. (13)–(14). Note that for most common loss functions such as the mean-squared error, the two final terms of Eq. (13)–(14) are zero as the Hessian h has no variation. When training a decision tree $f^{(k)}$, we store the obtained expectation and variance of each leaf of each tree and use these results to obtain our final estimate for the mean and variance using the update equations described in the next subsection.

3.3 Update Equations

Apart from the stochastic leaf weights, we require new update equations (Eq. (4)) in order to update the estimate for our mean and variance when adding a new tree at each iteration. These new update equations allow us to aggregate the stochastic weights over all the trees. For these equations, we make use of the following rules for the mean μ and variance σ^2 of some random variables (A, B) and a constant c :

$$\mu_{(A-cB)} = \mu_A - c \cdot \mu_B$$

$$\begin{aligned} \sigma_{cB}^2 &= c^2 \sigma_B^2 \\ \sigma_{(A-cB)}^2 &= \sigma_A^2 + c^2 \sigma_B^2 - 2c \cdot \sigma_{(A,B)}^2 \\ \sigma_{(A,B)}^2 &= \rho_{(A,B)} \sigma_A \sigma_B, \end{aligned}$$

in which ρ denotes Pearson’s correlation coefficient between the variables (A, B) . Using these rules, we can modify Eq. (4) to arrive at the formulas for the expectation E and variance V of our estimate $\hat{y}_i^{(k)}$:

$$\mu_{\hat{y}_i^{(k)}} = E \left[\hat{y}_i^{(k)} \right] = \mu_{\hat{y}_i^{(k-1)}} - \alpha \cdot \mu_{j^{(k)}} \quad (15)$$

$$\sigma_{\hat{y}_i^{(k)}}^2 = V \left[\hat{y}_i^{(k)} \right] = \sigma_{\hat{y}_i^{(k-1)}}^2 + \alpha^2 \sigma_{j^{(k)}}^2 - 2\alpha \rho \sigma_{\hat{y}_i^{(k-1)}} \sigma_{j^{(k)}}. \quad (16)$$

where the hyperparameter ρ denotes the correlation coefficient between trees k and $k - 1$. We provide further discussion around ρ in Section 3.5. Finally, the learned expectation and variance can be used for creating probabilistic predictions of new samples after training our model by sampling from a distribution parameterized by these learned quantities:

$$\hat{y}_i^{(k)} \sim D \left(\mu_{\hat{y}_i^{(k)}}, \sigma_{\hat{y}_i^{(k)}}^2 \right), \quad (17)$$

We are now ready to fully present our method PGBM.

3.4 PGBM

Algorithm. We provide a succinct overview of the procedures for training and prediction with Probabilistic Gradient Boosting Machines (PGBM) in Algorithms 1& 2.

Training (Algorithm 1). In PGBM, gradient boosting is performed comparable to LightGBM [10] or xgboost [4], and PGBM employs global equal density histogram binning to bin continuous features into discrete bins in order to reduce the computational effort required to find the optimal split decision (Line 1). At the start of training, we initialize the estimate \hat{y} , typically with the mean of the training set in a regression setting (Line 2). Then, gradient boosting is performed for a fixed number of iterations by first computing the gradient and hessian (Lines 4–5) of the training set and optionally choosing a subsample of the dataset (commonly referred to as bootstrapping, Line 6) on which to build the decision tree. The decision tree is then constructed up to a fixed number of leaves (Line 7) by first selecting the samples in the current node (Line 8), second by finding the best split for this node (Line 9), and third by splitting the current node or creating stochastic leaf weights if no split can be made (Line 10), for example, when the split does not result in a positive gain according to Eq. (2). After the tree construction has finished, predictions for the entire training set are generated (Line 12) and the overall estimate is updated (Line 13) and the process repeats for the next iteration.

Note that the learned variance is only used to create the probabilistic estimate in the prediction algorithm; it can also serve as a validation criterion during training (for example, by performing a prediction step on a validation set and deciding based on some probabilistic metric whether to continue training or not).

Prediction (Algorithm 2). During prediction, we initialize the estimate using the stored initial estimate of the training set (Line 1). Then, we make predictions on the dataset by iterating over all the

trees (Line 2) using our new update equations (Line 3). Finally, we obtain our probabilistic estimate by sampling from a distribution parameterized by our learned mean and variance (Line 5).

Algorithm 1 PGBM training algorithm

Input: Input dataset $\mathbf{X} \in \mathbb{R}^{n \times f}$ with n samples and f features, target output $\mathbf{y} \in \mathbb{R}^n$, differentiable loss function $l(\mathbf{y}, \hat{\mathbf{y}})$ and model hyperparameters.

Output:

- 1: Bin features such that for each feature $|x| \leq \max_bins$
 - 2: Set initial estimate $\hat{\mathbf{y}}$, e.g. to mean $\bar{\mathbf{y}}$ of target output
 - 3: **for** $k = 1$ **to** num_iterations **do**
 - 4: Compute gradient $\mathbf{g}^{(k)} = \nabla_{\hat{\mathbf{y}}^{(k)}} l(\mathbf{y}, \hat{\mathbf{y}})$
 - 5: Compute hessian $\mathbf{h}^{(k)} = \nabla_{\hat{\mathbf{y}}^{(k)}}^2 l(\mathbf{y}, \hat{\mathbf{y}})$
 - 6: Select subsample of input dataset as instance set I_1
 - 7: **for** $j = 1$ **to** \max_leaves **do**
 - 8: Select instance set I_j of \mathbf{X} , $\mathbf{g}^{(k)}$, $\mathbf{h}^{(k)}$
 - 9: Find best split for all (features, bins) (Eq. (2))
 - 10: Create split if split criteria are met else create stochastic leaf weight (Eq. (7)–(14))
 - 11: **end for**
 - 12: Predict \mathbf{X} to obtain $\mu_{\hat{\mathbf{y}}}^{(k)}$ (Eq. (15))
 - 13: Update estimate $\hat{\mathbf{y}} = \mu_{\hat{\mathbf{y}}}^{(k)}$
 - 14: **end for**
-

Algorithm 2 PGBM prediction algorithm

Input: Input dataset $\mathbf{X} \in \mathbb{R}^{n \times f}$ with n samples and f features, target distribution D and model hyperparameters.

Output:

- 1: Set initial estimate $\hat{\mathbf{y}}$ to mean $\bar{\mathbf{y}}$ of training dataset
 - 2: **for** $k = 1$ **to** num_iterations **do**
 - 3: Predict \mathbf{X} to obtain $(\mu_{\hat{\mathbf{y}}}^{(k)}, \sigma_{\hat{\mathbf{y}}}^2)^{(k)}$ (Eq. (15)–(16))
 - 4: **end for**
 - 5: Draw $n_samples$ $\hat{\mathbf{y}} \sim D\left(\mu_{\hat{\mathbf{y}}}^{(k)}, \sigma_{\hat{\mathbf{y}}}^2\right)$
-

Implementation. We implement PGBM in PyTorch [16] and offer it as a Python package.³ PyTorch offers (multi-)GPU acceleration by default, which allows us to scale PGBM to problems involving a large number of samples (we trained on datasets of over 10M samples) as we can distribute training across multiple GPUs. More importantly, our implementation allows the use of the automated differentiation engine of PyTorch, such that we can employ arbitrary complex differentiable loss functions without requiring an analytical gradient and hessian. This is in stark contrast to existing popular packages such as LightGBM [10] or xgboost [4], where custom loss functions require the manual derivation of an analytical gradient and hessian. We provide an example of this benefit in Section 4.2. Note that PGBM can be made compatible with existing gradient boosting packages relatively easily too, as it only requires storing one additional sample statistic (the variance), changing the

³<https://github.com/elephant/pgbm>

update equations according to Section 3.3 and choosing a distribution D after training to sample from. We intend to provide an implementation of PGBM within existing popular gradient boosting packages in the future.

Furthermore, we implemented a custom CUDA kernel that integrates with PyTorch to calculate the optimal splitting decision (Eq. (2)), the most compute intensive part of PGBM. Our kernel leverages the parallel processing power of modern CUDA-capable GPUs, by parallelizing the split decision across the sample and feature dimension using parallel reductions. We demonstrate the effectiveness of our GPU training in Section 4.1.

3.5 Analysis & Discussion

Computational complexity. Even though two parameters – a mean and variance – are learned in PGBM, the trees are constructed comparable to standard gradient boosting such as in [4, 10]. Therefore, the additional cost of our second parameter is negligible as only the sample statistics need to be calculated in the leaves. In contrast to NGBM [6], PGBM also does not require calculation of a natural gradient, which involves the inversion of many small matrices. PGBM’s runtime generally scales with the number of samples, the number of features and the number of bins used to bin the features, in accordance with existing GBM packages.

Higher-order moments and leaf sample quantiles. We only consider the first two moments of a distribution (i.e., the mean and variance) to derive our stochastic leaf weights, which limits the output distribution D to distributions parameterized using location and scale parameters (i.e., our learned mean and variance). This is a limitation compared to, e.g., NGBM [6]. We considered calculating higher order sample statistics such as the sample skewness (third moment) and sample kurtosis (fourth moment), however the disadvantage is that (i) there is no unbiased sample statistic for those measures, (ii) deriving approximations of the form of Eq. (13)–(14) becomes exceedingly complex, and (iii) higher-order sample statistics require more samples in order for the sample statistic to provide a reasonable estimate of the true statistic. Moreover, as we learn separate sample statistics for each leaf in each tree, we are still able to model complex distributions over the entire dataset using distributions parameterized only by location and scale parameters.

Finally, one could also store the sample quantile information of each leaf and draw samples according to the stored quantile information. This would remove the need for specifying a particular distribution. While this seems an attractive option, calculating sample quantile information for each leaf is computationally difficult as it requires an expensive sorting operation, and storing a sufficient number of sample quantiles to reap the full benefits of this method requires storing at least 2–3x the number of leaf weights. In short, there is no real need to use higher order moments or leaf sample quantiles to provide accurate probabilistic estimates as we demonstrate in Section 4.1.

Output sampling. PGBM allows one to sample from different output distributions *after* training, which allows practitioners to train their model by minimizing some point metric (e.g., RMSE) and after training try different distributions for optimizing the probabilistic forecast based on some validation metric. The key benefit is that

this allows PGBM to achieve state-of-the-art point forecasting performance *as well as* accurate probabilistic forecasting performance *using the same model*. We will demonstrate this in Section 4.1. Note that practitioners can of course also directly optimize the probabilistic forecast by using a loss function that optimizes the probabilistic forecast.

Split decisions and tree dependence. In PGBM, split decisions in the tree are not recomputed based on the stochasticity of the leaf weights, even though it could be argued that this would be appropriate when sampling from the trees. We intentionally avoid this as it is computationally expensive to recompute split decisions after training when sampling from the learned distribution. Secondly, by sequentially adding the mean and variance of each tree we omit the explicit covariance between tree k and trees $k - 2, k - 3, \dots$, and only model the covariance between subsequent trees. We implicitly model both these effects using a single constant correlation hyperparameter ρ (Eq. (16)), for which we provide a more detailed analysis in Section 4.1.

Hessian distribution. The distribution of the Hessian h should have a support of $[0, \infty)$ to avoid division by zero in Eq. (13)–(14), or equivalently, we require the sum of the Hessians (plus regularization constant λ) of all samples in an instance set I_j of a leaf to be positive. For common convex loss functions such as the mean-squared error this is not an issue, however for non-convex loss functions this might pose a problem in rare cases where all Hessians in an instance set add up to zero. In those cases, numerical issues can usually be avoided by requiring a decent (e.g., > 10) minimum number of samples in each leaf in a tree – this can be set as a hyperparameter in PGBM.

4 EXPERIMENTS

In this section, we first demonstrate how PGBM can provide accurate point and probabilistic predictions on a set of common regression benchmark datasets from the UCI Machine Learning Repository (Section 4.1). We show how PGBM allows practitioners to optimize their probabilistic estimate after training, thereby removing the need to retrain a model under different posterior distribution assumptions. Next, we demonstrate the efficiency of our implementation of PGBM compared to existing gradient boosting methods. Finally, we demonstrate PGBM on the problem of forecasting for hierarchical time series, which requires optimizing a complex loss function for which deriving an analytical gradient is too complex (Section 4.2).

To facilitate reproducibility of the results in this paper, our experiments are run on open data, and our experimentation code is available online.⁴

4.1 UCI regression benchmarks

Task. We perform probabilistic regression on a set of regression datasets. Our goal is to obtain the lowest probabilistic prediction score as well as the lowest point performance score.

Evaluation. We evaluate the probabilistic performance of each method using the Continuously Ranked Probability Score (CRPS),

which is a measure of discrepancy between the empirical cumulative distribution function of an observation y and the cumulative distribution F of a forecast \hat{y} [23]:

$$C = \int [F(\hat{y}) - \mathbb{1}(\hat{y} \geq y)]^2 d\hat{y}, \quad (18)$$

in which $\mathbb{1}$ denotes the indicator function. We compute the empirical CRPS based on 1,000 sample predictions generated by the trained models on the test set. We evaluate point performance using Root Mean Squared Error (RMSE):

$$RMSE = \sqrt{\frac{1}{n} \sum_i^n (y_i - \hat{y}_i)^2}. \quad (19)$$

We present these metrics relative to the median of PGBM over all the folds tested for a dataset, and we refer the reader to Table 4 of Supplemental Materials B for further details.

Protocol. We follow the same protocol as Duan et al. [6], and create 20 random folds for each dataset except for `msd` for which we only create one. For each of these folds, we keep 10% of the samples as test set. The remaining 90% is first split into an 80/20 validation/training set to find the number of training iterations that results in the lowest validation score. After validation, the full 90% training set is trained using the number of iterations found in the validation step. As output distribution for the probabilistic prediction we use a Normal distribution, similar to Duan et al. [6].

Baseline models. For probabilistic performance, we compare against NGBoost [6], which has recently been shown to outperform other comparable methods on the current set of benchmarks. We use the same settings for NGBoost as in [6]. For point performance, we also compare to LightGBM [10], one of the most popular and best-performing gradient boosting packages available. We configure LightGBM to have the same settings as PGBM.

PGBM. For all datasets, we use the same hyperparameters for PGBM, except that we use a bagging fraction of 0.1 for MSD in correspondence with Duan et al. [6]. Our training objective in PGBM is to minimize the mean-squared error (MSE). We refer the reader to Table 5 of Supplemental Materials B for an overview of key hyperparameters of each method.

Results. We provide the results of our first experiment in Figure 1 and observe the following:

- On probabilistic performance, PGBM outperforms NGBoost on average by approximately 15% as demonstrated by the relatively lower CRPS across all but one dataset (`msd`, where the difference is tiny). This is remarkable, as the training objective in PGBM is to minimize the MSE, rather than optimize the probabilistic forecast as does NGBoost.
- PGBM outperforms NGBoost on all datasets on point performance, and on average by almost 20%, which is in line with expectation as we explicitly set out to optimize the MSE as training objective in PGBM in this experiment. However, as becomes clear from this result, PGBM does not have to sacrifice point performance in order to provide state-of-the-art probabilistic estimates nonetheless. Compared to LightGBM, PGBM performs slightly worse on average (approx. 3%) on point performance. We suspect this is due to implementation specifics.

⁴Repository at <https://github.com/elephaint/pgbm>

The main takeaway from this experiment is that even though we only optimized for a point metric (MSE) in PGBM, we were still able to achieve similar probabilistic performance compared to a method that explicitly optimizes for probabilistic performance.

Analysis: correlation hyperparameter. We perform a brief analysis of the correlation hyperparameter ρ (Eq. (16)). This hyperparameter controls the dependence between variance estimates of subsequent trees in PGBM, and is critical for probabilistic performance. Figure 2a shows the CRPS evaluated on the validation set at different settings for ρ for each dataset for a single fold. We normalized the CRPS scores on the lowest CRPS for each dataset. Across all datasets, the CRPS seems to follow a parabolic shape and consequently there seems to be an optimal choice for ρ across different datasets: a value of 0.02–0.07 typically seems appropriate. Empirically, we found that an initial value of $\rho = \frac{\log_{10} n}{100}$, where n denotes the size of the training set generally works well and therefore we used that in our experiments as default value. Intuitively, the positiveness of the correlation between subsequent trees seems logical: if the leaf weight of a given tree shifts more positively (negatively) as a result of stochasticity, the residual on which the next tree will be constructed shifts in the same direction. Consequently, the leaf weights of this next tree will also shift in the same direction, thus exhibiting a positive correlation with the previous tree’s leaf weights. Furthermore, larger datasets tend to cluster together in behavior, as can be seen from the curves for the protein, msd, kin8nm, power and naval datasets. It seems that for larger datasets, typically larger settings for ρ are appropriate. Our hypothesis is that for trees containing many samples per leaf, the correlation between subsequent trees is higher, as more samples in the tree’s leaves will generally imply that the model has not yet (over)fit to the training set and there is likely more information left in the residuals compared to the situation where there are few samples per leaf. This would explain the behavior observed in Figure 2a as we train each model with a maximum of 8 leaves per tree, resulting in more samples per leaf for larger datasets. We test this hypothesis by training the relatively larger protein and msd datasets using different settings for the maximum number of leaves, for which we show the results in Figure 2b. Confirming our hypothesis, we indeed observe the optimal correlation parameter decreasing when we train PGBM using a higher number of maximum leaves per tree (i.e., the minimum of the parabola shifts to the left).

Analysis: posterior distribution. One of the key benefits of PGBM is that it allows us to optimize the probabilistic estimate *after* training, as the choice of distribution D in Eq. (6) is independent from the training process. This offers the benefit of training to optimize a certain point metric such as RMSE, and choosing the distribution that best fits the learned mean and variance after training by validating a set of distribution choices on a validation set. To demonstrate this benefit, we repeated the experiments from our first experiment for a single fold. For each dataset, we evaluated the learned model on CRPS on the validation set for a set of common distributions and a range of tree correlations $\rho = \{0.00, 0.01, \dots, 0.09\}$. The optimal choice of distribution and tree correlation on the validation set was subsequently used for calculating the CRPS on the test set. We report the results in Table 1, where ‘Base case’ refers to the base case scenario from our first experiment, where we chose a Normal

distribution and a tree correlation hyperparameter of $\rho = \frac{\log_{10} n}{100}$ across all datasets, and ‘Optimal’ refers to the result on the test set when choosing the distribution and tree correlation according to the lowest CRPS on the validation set. We see that for most datasets, the minimum CRPS on the validation set is similar across choices of distribution, which implies that it is more beneficial to optimize the tree correlation rather than the choice of output distribution for these datasets. On the test set, we see improved scores compared to the base case on all but the smallest dataset, thereby showcasing the benefit of optimizing the probabilistic forecast after training. We would advice practitioners to start with a generic distribution such as the normal or Student’s $t(3)$, and optimize the probabilistic estimate after training by testing for different tree correlations and distribution choices.

Analysis: training time. Our implementation in PyTorch allows us to use GPU training by default, which allows us to significantly speed up training for larger datasets. We demonstrate this benefit in Table 2 where we compare training times for datasets of different size against a baseline of PGBM (we refer to Table 6 in the Supplemental Materials for the absolute timings). For this experiment, we also included the higgs dataset, which is a 10M sample UCI dataset commonly used to benchmark gradient boosting packages. For PGBM, we show results for training on GPU-only and CPU-only. NGBoost does not offer GPU training and runs on top of the default scikit-learn decision tree regressor. We ran our experiments on a 6-core machine with a nVidia RTX 2080Ti GPU. As can be seen, PGBM is up to several orders of magnitude faster than NGBoost as the dataset size increases. This demonstrates that PGBM and our implementation allow practitioners to solve probabilistic regression problems for datasets much larger than NGBoost.

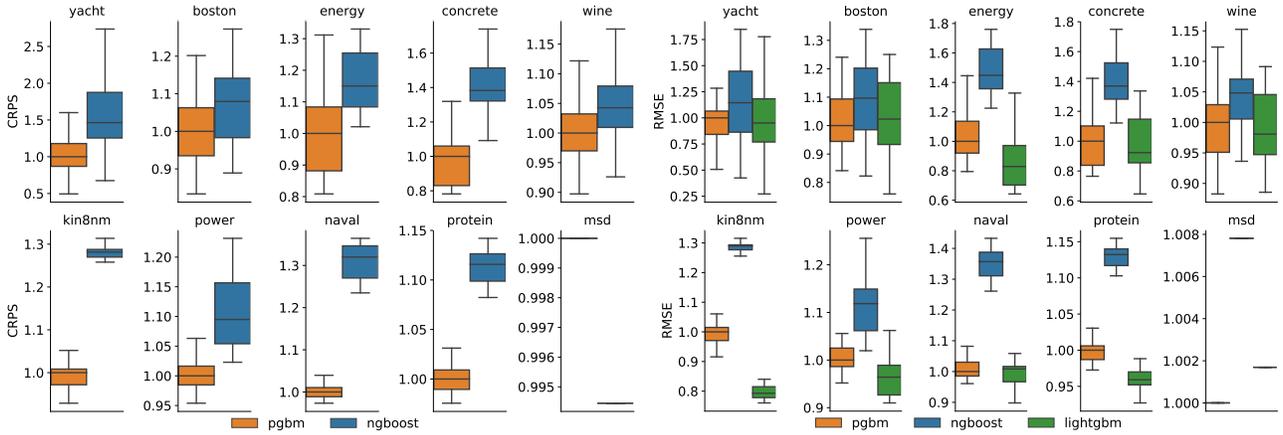
In general, for smaller datasets, training on cpu offers the best timings for the two methods. We include the relative timings to LightGBM for reference in Table 2, which shows that PGBM even becomes competitive to LightGBM for the largest dataset (higgs). However, the timings for LightGBM represent the timings to train a *single* LightGBM model. If one is interested in obtaining a probabilistic forecast, the timings would be multiplied by the number of quantiles required. Hence, for a fine-grained probability distribution, the timings would be 5–10x higher for LightGBM, again demonstrating the effectiveness of our implementation for probabilistic forecasting.

4.2 Hierarchical time series

Task. So far, our experimental results were obtained using the mean-squared error as objective function for which an analytical gradient and hessian can be easily derived. In this experiment, we apply PGBM to the problem of hierarchical time series forecasting, which is a problem where our loss function is rather complex, so that it becomes very tedious to manually calculate an analytical gradient and hessian for it:

$$L = \sum_j^N w_j (y_j - \hat{y}_j)^2, \quad (20)$$

where w_j is the weight of the j -th time series, and N is the number of time series. In hierarchical time series, we aggregate time series across multiple levels. For example, in the case of two time series



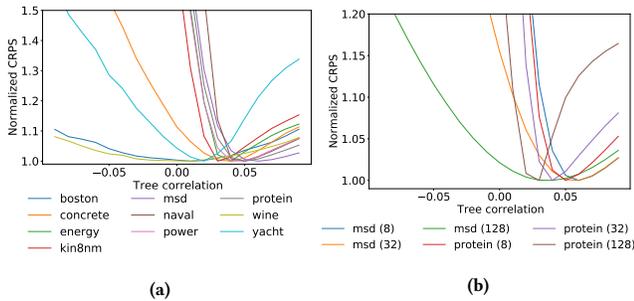
(a) Probabilistic performance (CRPS)

(b) Point performance (RMSE)

Figure 1: Results for probabilistic (CRPS, left) and point (RMSE, right) performance for each dataset for each method, with the smallest dataset yacht in the top left corner and the largest dataset msd in the lower right corner. Lower is better, and results have been indexed against the median test score of PGBM. PGBM outperforms NGBoost both in probabilistic and point performance.

Table 1: Results for probabilistic (CRPS) performance for each dataset when varying the posterior distribution and tree correlations on the validation set (left) and subsequently using the optimal choice of distribution and tree correlation on the test set (right). Best results on validation and test set are in bold. We report the minimum CRPS on the validation set across all the tree correlations tested per distribution.

Dataset	Validation set									Test set	
	Normal	Student's t(3)	Logistic	Laplace	LogNormal	Gumbel	Weibull	Poisson	NegativeBinomial	Base case	Optimal
yacht	0.37	0.37	0.37	0.37	0.44	0.38	0.37	0.75	9.7	0.18	0.19
boston	1.41	1.4	1.39	1.39	1.4	1.42	1.39	1.58	12.27	2.08	2.06
energy	0.62	0.61	0.62	0.61	0.62	0.62	0.63	1.25	16.15	0.64	0.64
concrete	2.21	2.21	2.21	2.19	2.23	2.26	2.22	2.38	3.59	2.69	2.64
wine	0.32	0.33	0.32	0.32	0.32	0.32	0.32	0.59	0.64	0.35	0.34
kin8nm	0.08	0.08	0.08	0.08	1.04	0.08	0.08	0.26	0.26	0.08	0.08
power	1.86	1.86	1.86	1.86	1.86	1.87	1.88	5.22	454.34	1.81	1.80
naval	0.00	0.00	0.00	0.00	0.02	0.00	0.00	0.22	0.22	0.00	0.00
protein	2.18	2.18	2.18	2.18	45.3	2.17	2.54	2.15	2.15	2.14	2.12
msd	4.74	4.73	4.74	4.73	4.74	4.88	4.69	11.16	1982.19	4.78	4.74



(a)

(b)

Figure 2: Normalized CRPS on the validation set for different settings of tree correlation hyperparameter ρ , (a) for all datasets and (b) for protein and msd when trained using a maximum number of leaves per tree of $\{8, 32, 128\}$.

and two levels, $N = 3$ and our loss for each series reads $L_1 = w_1(y_1 - \hat{y}_1)^2$, $L_2 = w_2(y_2 - \hat{y}_2)^2$, $L_3 = w_3((y_1 + y_2) - (\hat{y}_1 + \hat{y}_2))^2$ with w_1, w_2, w_3 weights of each series, for example 0.25, 0.25, 0.5. Hence, the gradient and hessian of L with respect to the first estimate \hat{y}_1 becomes $\frac{\partial L}{\partial \hat{y}_1} = -w_1(2y_1 - 2\hat{y}_1) - w_3(2y_1 + 2y_2 - 2\hat{y}_1 + 2\hat{y}_2)$

Table 2: Training time for 2,000 iterations for 5 datasets of different size as a fraction of PGBM-gpu training time. Bold indicates probabilistic forecasting method with the lowest training time.

Dataset	Probabilistic forecast			Point forecast
	PGBM-gpu	PGBM-cpu	NGBoost	LightGBM
wine (n=1,599)	1.00	0.41	0.20	0.01
naval (n=12k)	1.00	1.62	1.05	0.02
protein (n=46k)	1.00	3.09	3.39	0.02
msd (n=515k)	1.00	26.85	48.83	0.09
higgs (n=10,5M)	1.00	101.76	147.78	0.43

and $\frac{\partial^2 L}{\partial \hat{y}_1^2} = 2w_1 + 2w_3$. It becomes clear that deriving this result analytically becomes increasingly complex when the number of levels and the number of time series increases, which necessitates the use of autodifferentiation packages such as PyTorch if we are interested in optimizing this objective.

Dataset. We use a subset of the dataset from the M5 forecasting competition [13], in which participants were asked to provide hierarchical forecasts for Walmart retail store products. We use a single store and create forecasts for a single day. For each store, we are

Table 3: Point (RMSE) and probabilistic (CRPS) forecasting performance for the M5 dataset across aggregations when using MSE or weighted MSE (only for PGBM) as training objective. Lower is better, best results are indicated in bold.

objective	RMSE				CRPS		
	PGBM		NGBoost	LightGBM	PGBM		NGBoost
	MSE	wMSE	MSE	MSE	MSE	wMSE	MSE
Individual ¹	2.00	2.00	2.01	2.00	0.77	0.93	0.78
Category ²	67.9	67.6	77.4	70.0	72.6	40.4	82.0
Department ³	108	101	129	111	160	59	184
Total ⁴	213	190	276	225	472	136	560

1. n=85,372 2. n=196 3. n=84 4. n=28.

interested not only in accurately forecasting individual item sales, but also in optimizing the aggregate sales per day, aggregate sales per day per category and aggregate sales per day per department. Hence, we obtain four levels for our weighted loss function: (i) individual items, (ii) category aggregates per day, (iii) department aggregates per day, and (iv) total daily aggregates. For more details on the data and preprocessing we refer to Supplemental Materials B.

Protocol. We compare against a baseline of LightGBM, NGBoost and PGBM trained with the regular mean-squared error objective. All models are trained using the same hyperparameters. We validate on the last 28 days of the training dataset and pick the number of iterations resulting in the lowest item RMSE on this validation set. After validating, we retrain on the entire dataset for the number of optimal iterations and test on a held out test set of 28 days (with the first day starting the day after the last day in the validation set). We use the Normal distribution with a tree correlation of $\rho = \frac{\log_{10} n}{100}$ to generate our probabilistic forecasts for PGBM.

Results. We evaluate our model on RMSE and CRPS for each aggregation and the results are displayed in Table 3. We observe that using the weighted MSE that incorporates our four levels of aggregation results in a similar point forecast score for individual items, but in a much better forecast for the aggregations – differences up to 10% compared to the second-best point performance of PGBM are observed. Secondly, we see that the gain using the weighted MSE increases at hierarchically higher aggregation levels such as ‘total by date’. This is important, as this implies that we are able to generate item-level forecasts that are more consistent with higher-level aggregates. Finally, we observe that item-level CRPS is worse in the weighted MSE setting compared to the regular MSE setting, whereas our probabilistic estimate for higher aggregations improves up to 300% when using the weighted MSE. This can be expected: in the MSE setting, each individual item forecast ‘does not consider’ aggregates in the category or department, whereas in the weighted MSE setting, item forecasts are optimized to also consider the impact on the overall aggregations. All in all, this experiment demonstrates the benefit of our implementation: we can optimize over more complex loss functions, thereby enabling probabilistic forecasts of more complex problems such as hierarchical time series problems.

5 RELATED WORK

Traditional forecasting methods such as ARIMA [3] allow for probabilistic forecasts through specification of confidence intervals [9]. However, creating a confidence interval in these methods often

requires assuming normality of the distribution of the target variable or its residuals. Generalized Additive Models for Shape, Scale and Location (GAMLSS) [17] is a framework that allows for a more flexible choice of distribution of the target variable in probabilistic regression problems. A disadvantage is that the model needs to be pre-specified, limiting flexibility of this method. Prophet [22] is a more recent example of generalized additive models applied to the probabilistic forecasting setting. However, Prophet has been shown to underperform other contemporary probabilistic forecasting methods [1, 19] and to have difficulties scaling to very large datasets [19]. Other Bayesian methods exhibiting similar issues include Bayesian Additive Regression Trees (BART) [5], which requires computationally expensive sampling techniques to obtain probabilistic predictions.

Gradient Boosting Machines (GBMs) [7] are widely used for regression problems such as forecasting [4]. Popular GBM implementations such as LightGBM [10] or xgboost [4] have won various machine learning competitions [4]. The winning solution of the accuracy track of the recent M5 forecasting competition was based on a set of LightGBM models, and 4 out of the top-5 solutions used LightGBM models in their solutions [13]. However, GBMs are not naturally equipped to provide probabilistic forecasts as these models return point predictions, requiring multiple models when a practitioner desires a range of predictions. For example, the uncertainty track of the M5 forecasting competition required contestants to provide a set of quantiles for a hierarchical time series problem. The winning solution was based on 126 (!) separate LightGBM models, one for each requested quantile and time series aggregation level [14]. To address these limitations, NGBoost [6] allows for probabilistic regression with a single GBM by using the natural gradient to boost multiple parameters of a pre-specified distribution. Compared to NGBoost, our method PGBM does not require a natural gradient; it can achieve better or on-par predictive uncertainty estimates, without sacrificing performance on the point forecast of the same model as does NGBoost. Taieb et al. [21] also propose boosted additive models for probabilistic forecasting. Our work is different in that we use GBMs with stochastic leaf weights to estimate the conditional mean and variance simultaneously for each estimator. Gouk et al. [8] present a method for incrementally constructing decision trees using stochastic gradient information. Our method is different in that (i) we focus on the general case of probabilistic regression instead of incremental online learning of trees and (ii) we obtain our stochastic estimates by approximating the ratio of the gradient and hessian.

Outside of the GBM context, decision trees have also been used for probabilistic regression problems in Quantile Regression Forests (QRF) [15]. However, this method requires storing all observations when computing leaf weights of a decision tree, which makes this method less suitable for large datasets. In addition, GBMs commonly outperform random forests on regression tasks, making the former a better choice when performance is a key consideration.

Contemporary large-scale probabilistic forecasting techniques often leverage the power of neural networks, such as DeepAR [18] or Transformer-based models [11, 12]. However, in practice GBMs still seem the technique of choice—only one out of the top-5 solutions in the M5 uncertainty forecasting competition used a neural network method as its primary probabilistic forecasting tool.

In summary, we contribute the following on top of the related work discussed above: (i) PGBM is a single-parameter boosting method that achieves state-of-the-art point and probabilistic estimates using a single model, (ii) PGBM allows choosing an output distribution after training, which means the probabilistic forecast can be optimized after training, (iii) our implementation allows training of larger datasets up to several orders of magnitude faster than the existing state-of-the-art, and (iv) our implementation allows using complex differentiable loss functions, which removes the need to calculate an analytical gradient, thereby opening up a wider set of problems that can be effectively addressed.

6 CONCLUSION

In this work we introduced Probabilistic Gradient Boosting Machines (PGBM), a method for probabilistic regression using gradient boosting machines. PGBM creates probabilistic estimates by using stochastic tree leaf weights based on sample statistics. By combining the learned weights for each subsequent tree, PGBM learns a mean and variance for samples in a dataset which can be used to sample from an arbitrary distribution of choice. We demonstrated that PGBM provides state-of-the-art probabilistic regression results across a range of small to large datasets. Benefits of PGBM compared to existing work are that (i) PGBM is a single-parameter boosting method that optimizes a point regression but achieves state-of-the-art probabilistic estimates using the same model, (ii) PGBM enables the choice of an output distribution after training, which means practitioners can optimize the choice of distribution without requiring retraining of the model, (iii) our implementation allows training of larger datasets up to several orders of magnitude faster than the existing state-of-the-art, and (iv) our implementation in PyTorch allows using complex differentiable loss functions which removes the need to calculate an analytical gradient as is common in existing gradient boosting packages. We demonstrated the latter benefit for the problem of hierarchical time series forecasting, where we observed up to 10% improvement in point performance and up to 300% improvement in probabilistic forecasting performance.

Limitations of our work are that PGBM only learns the mean and variance in a tree, which limits the choice of output distribution. However, we observed no negative performance effects in our experiments thereof.

In the future, we intend to further work on the theoretical error bounds of PGBM. Under mild assumptions, sample statistics in each leaf of each tree appropriately represent the true statistics of the samples in each leaf provided a sufficient number of samples. However, we have yet to determine appropriate theoretical error bounds on the final estimated statistics when considering the simplifications we make during decision tree learning, such as the greedy approximate split finding, using a limited number of tree leaves, our approximation to the stochastic leaf weights, keeping decision points constant and treating the correlation between subsequent trees as a constant across samples and trees. Regarding the latter, we also expect that the probabilistic estimate can be further improved by using a better approximation to the tree correlations instead of our choice of keeping it fixed across trees and samples.

REFERENCES

- [1] Alexander Alexandrov, Konstantinos Benidis, Michael Bohlke-Schneider, Valentin Flunkert, Jan Gasthaus, Tim Januschowski, Danielle C. Maddix, Syama Rangapuram, David Salinas, Jasper Schulz, Lorenzo Stella, Ali Caner Türkmen, and Yuyang Wang. 2020. GluonTS: Probabilistic and Neural Time Series Modeling in Python. *Journal of Machine Learning Research* 21, 116 (2020), 1–6.
- [2] Joos-Hendrik Böse, Valentin Flunkert, Jan Gasthaus, Tim Januschowski, Dustin Lange, David Salinas, Sebastian Schelter, Matthias Seeger, and Yuyang Wang. 2017. Probabilistic Demand Forecasting at Scale. *Proc. VLDB Endow.* 10, 12 (Aug. 2017), 1694–1705. <https://doi.org/10.14778/3137765.3137775>
- [3] G. E. P. Box and David A. Pierce. 1970. Distribution of Residual Autocorrelations in Autoregressive-Integrated Moving Average Time Series Models. *J. Amer. Statist. Assoc.* 65, 332 (1970), 1509–1526. <https://doi.org/10.2307/2284333>
- [4] Tianqi Chen and Carlos Guestrin. 2016. XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, San Francisco California USA, 785–794. <https://doi.org/10.1145/2939672.2939785>
- [5] Hugh A. Chipman, Edward I. George, and Robert E. McCulloch. 2010. BART: Bayesian Additive Regression Trees. *Ann. Appl. Stat.* 4, 1 (March 2010), 266–298. <https://doi.org/10.1214/09-AOAS285>
- [6] Tony Duan, Anand Avati, Daisy Ding, Khanh K. Thai, Sanjay Basu, Andrew Ng, and Alejandro Schuler. 2020. NGBoost: Natural Gradient Boosting for Probabilistic Prediction. *ICML (2020)*.
- [7] Jerome H. Friedman. 2001. Greedy Function Approximation: A Gradient Boosting Machine. *The Annals of Statistics* 29, 5 (2001), 1189–1232.
- [8] Henry Gouk, Bernhard Pfahringer, and Eibe Frank. 2019. Stochastic Gradient Trees. In *Asian Conference on Machine Learning*. PMLR, 1094–1109.
- [9] Rob J Hyndman. 2018. *Forecasting: Principles and Practice*.
- [10] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. 2017. LightGBM: A Highly Efficient Gradient Boosting Decision Tree. In *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc., 3146–3154.
- [11] Shiyang Li, Xiaoyong Jin, Yao Xuan, Xiyuo Zhou, Wenhui Chen, Yu-Xiang Wang, and Xifeng Yan. 2019. Enhancing the Locality and Breaking the Memory Bottleneck of Transformer on Time Series Forecasting. In *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc., 5244–5254.
- [12] Bryan Lim, Sercan O. Arik, Nicolas Loeff, and Tomas Pfister. 2019. Temporal Fusion Transformers for Interpretable Multi-Horizon Time Series Forecasting. *arXiv:1912.09363 [cs, stat]* (Dec. 2019). [arXiv:1912.09363 \[cs, stat\]](https://arxiv.org/abs/1912.09363)
- [13] Spyros Makridakis, Evangelos Spiliotis, and Vassilis Assimakopoulos. 2020. *The M5 Accuracy Competition: Results, Findings and Conclusions*.
- [14] Spyros Makridakis, Evangelos Spiliotis, Vassilis Assimakopoulos, Zhi Chen, Anil Gaba, Ilija Tsetlin, and Robert Winkler. 2020. *The M5 Uncertainty Competition: Results, Findings and Conclusions*.
- [15] Nicolai Meinshausen. 2006. Quantile Regression Forests. *Journal of Machine Learning Research* 7, 35 (2006), 983–999.
- [16] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc., 8024–8035.
- [17] R. A. Rigby and D. M. Stasinopoulos. 2005. Generalized Additive Models for Location, Scale and Shape. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 54, 3 (2005), 507–554. <https://doi.org/10.1111/j.1467-9876.2005.00510.x>
- [18] David Salinas, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. 2019. DeepAR: Probabilistic Forecasting with Autoregressive Recurrent Networks. *International Journal of Forecasting* (Oct. 2019). <https://doi.org/10.1016/j.ijforecast.2019.07.001>
- [19] Rajat Sen, Hsiang-Fu Yu, and Inderjit S Dhillon. 2019. Think Globally, Act Locally: A Deep Neural Network Approach to High-Dimensional Time Series Forecasting. In *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc., 4838–4847.
- [20] Fabio Sigrüst. 2020. Gradient and Newton Boosting for Classification and Regression. *Expert Systems with Applications* (Oct. 2020). <https://doi.org/10.1016/j.eswa.2020.114080>
- [21] Souhaib Ben Taieb, Raphael Huser, Rob J Hyndman, and Marc G Genton. 2015. Probabilistic Time Series Forecasting with Boosted Additive Models: An Application to Smart Meter Data. (2015), 30.
- [22] Sean J. Taylor and Benjamin Letham. 2018. Forecasting at Scale. *The American Statistician* 72, 1 (Jan. 2018), 37–45. <https://doi.org/10.1080/00031305.2017.1380080>
- [23] Michaël Zamo and Philippe Naveau. 2018. Estimation of the Continuous Ranked Probability Score with Limited Information and Applications to Ensemble Weather Forecasts. *Math Geosci* 50, 2 (Feb. 2018), 209–234. <https://doi.org/10.1007/s11004-017-9709-7>

A DERIVATION OF STOCHASTIC LEAF WEIGHTS

A.1 Expectation

We approximate the mean in each leaf by using a second-order Taylor approximation of the expectation of a function of the two random variables (g, h) around the point $\mathbf{a} = (\bar{g}, \bar{h})$:

$$\begin{aligned} E[f(g, h)] &= E[f(\mathbf{a}) + f'_g(\mathbf{a})(g - \bar{g}) + f'_h(\mathbf{a})(h - \bar{h}) \\ &\quad + \frac{1}{2}f''_{gg}(\mathbf{a})(g - \bar{g})^2 + f''_{gh}(\mathbf{a})(g - \bar{g})(h - \bar{h}) \\ &\quad + \frac{1}{2}f''_{hh}(\mathbf{a})(h - \bar{h})^2 + H] \end{aligned} \quad (21)$$

with H denoting the higher-order terms that we drop for our estimate. Using the laws of expectations we then obtain:

$$\begin{aligned} E[f(g, h)] &\approx E[f(\mathbf{a})] + E[f'_g(\mathbf{a})(g - \bar{g})] + E[f'_h(\mathbf{a})(h - \bar{h})] \\ &\quad + E[\frac{1}{2}f''_{gg}(\mathbf{a})(g - \bar{g})^2] + E[f''_{gh}(\mathbf{a})(g - \bar{g})(h - \bar{h})] \\ &\quad + E[\frac{1}{2}f''_{hh}(\mathbf{a})(h - \bar{h})^2] \end{aligned} \quad (22)$$

$$\begin{aligned} &= E[f(\mathbf{a})] + f'_g(\mathbf{a}) \underbrace{E[(g - \bar{g})]}_0 + f'_h(\mathbf{a}) \underbrace{E[(h - \bar{h})]}_0 \\ &\quad + \frac{1}{2}f''_{gg}(\mathbf{a}) \underbrace{E[(g - \bar{g})^2]}_{\sigma_g^2} + f''_{gh}(\mathbf{a}) \underbrace{E[(g - \bar{g})(h - \bar{h})]}_{\sigma_{gh}^2} \\ &\quad + \frac{1}{2}f''_{hh}(\mathbf{a}) \underbrace{E[(h - \bar{h})^2]}_{\sigma_h^2} \end{aligned} \quad (23)$$

$$\begin{aligned} &= E[f(\mathbf{a})] + \frac{1}{2}f''_{gg}(\mathbf{a})\sigma_g^2 + f''_{gh}(\mathbf{a})\sigma_{gh}^2 \\ &\quad + \frac{1}{2}f''_{hh}(\mathbf{a})\sigma_h^2, \end{aligned} \quad (24)$$

with σ_{gh}^2 denoting the covariance of the gradient and the hessian. For a function $f(g, h) = \frac{g}{h}$, we have:

$$\begin{aligned} f'_g &= h^{-1} \\ f'_h &= -gh^{-2} \\ f''_{gg} &= 0 \\ f''_{gh} &= -h^{-2} \\ f''_{hh} &= 2gh^{-3}. \end{aligned}$$

Substituting and using $\mathbf{a} = (\bar{g}, \bar{h})$, $f(\mathbf{a}) = \frac{\bar{g}}{\bar{h}}$:

$$E[f(g, h)] \approx E[f(\mathbf{a})] - h^{-2}(\mathbf{a})\sigma_{gh}^2 + gh^{-3}(\mathbf{a})\sigma_h^2 \quad (25)$$

$$= \frac{\bar{g}}{\bar{h}} - \frac{\sigma_{gh}^2}{\bar{h}^2} + \frac{\bar{g}\sigma_h^2}{\bar{h}^3}. \quad (26)$$

Finally, we can include the regularization constant $\bar{\lambda}$ to arrive at the final estimate of the expectation for the leaf weight w_j . This constant only affects the mean of the random variable h , therefore

we can safely add it to the terms containing \bar{h} :

$$E\left[\frac{\bar{g}}{(\bar{h} + \bar{\lambda})}\right] \approx \frac{\bar{g}}{(\bar{h} + \bar{\lambda})} - \frac{\sigma_{gh}^2}{(\bar{h} + \bar{\lambda})^2} + \frac{\bar{g}\sigma_h^2}{(\bar{h} + \bar{\lambda})^3}. \quad (27)$$

Note that we can obtain the first-order Taylor approximation of the mean by dropping the last two terms of Eq. (27):

$$E\left[\frac{\bar{g}}{(\bar{h} + \bar{\lambda})}\right] \approx \frac{\bar{g}}{(\bar{h} + \bar{\lambda})}. \quad (28)$$

A.2 Variance

For the variance, we start with the definition of variance for a function $f(g, h)$:

$$V[f(g, h)] = E[(f(g, h) - E[f(g, h)])^2]. \quad (29)$$

We perform a first-order Taylor expansion of $f(g, h)$ around the point $\mathbf{a} = (\bar{g}, \bar{h})$ and we substitute the first-order approximation of the mean:

$$\begin{aligned} V[f(g, h)] &\approx E\left[\left(f(\mathbf{a}) + f'_g(\mathbf{a})(g - \bar{g}) + f'_h(\mathbf{a})(h - \bar{h}) - E[f(\mathbf{a})]\right)^2\right] \end{aligned} \quad (30)$$

$$= E[(f'_g(\mathbf{a})(g - \bar{g}) + f'_h(\mathbf{a})(h - \bar{h}))^2] \quad (31)$$

$$= E[f_g'^2(\mathbf{a})(g - \bar{g})^2 + f_h'^2(\mathbf{a})(h - \bar{h})^2 + 2f'_g(\mathbf{a})(g - \bar{g})f'_h(\mathbf{a})(h - \bar{h})] \quad (32)$$

$$= f_g'^2(\mathbf{a})E[(g - \bar{g})^2] + f_h'^2(\mathbf{a})E[(h - \bar{h})^2] + 2f'_g(\mathbf{a})f'_h(\mathbf{a})E[(g - \bar{g})(h - \bar{h})] \quad (33)$$

$$= \bar{h}^{-2}\sigma_g^2 + \bar{g}^2\bar{h}^{-4}\sigma_h^2 - 2\bar{g}\bar{h}^{-3}\sigma_{gh}^2 \quad (34)$$

$$= \frac{\sigma_g^2}{\bar{h}^2} + \frac{\bar{g}^2\sigma_h^2}{\bar{h}^4} - 2\frac{\bar{g}\sigma_{gh}^2}{\bar{h}^3}. \quad (35)$$

Finally, including the regularization constant $\bar{\lambda}$ we obtain:

$$V\left[\frac{\bar{g}}{(\bar{h} + \bar{\lambda})}\right] \approx \frac{\sigma_g^2}{(\bar{h} + \bar{\lambda})^2} + \frac{\bar{g}^2\sigma_h^2}{(\bar{h} + \bar{\lambda})^4} - 2\frac{\bar{g}\sigma_{gh}^2}{(\bar{h} + \bar{\lambda})^3}. \quad (36)$$

B REPRODUCIBILITY

We report absolute scores and dataset statistics for the UCI benchmark in Table 4. An overview of the key hyperparameters for each method for both experiments is given in Table 5, and absolute timings for the timings of Table 2 in Table 6. An overview of the M5 dataset is given in Table 7. We refer to our code at <https://github.com/elephant/pgbm> for further details, such as the features of the M5 dataset, which mainly comprise lagged target variables, time indicators (e.g., day-of-week), event indicators (e.g., holidays) and item indicators.

Table 4: Results for probabilistic (CRPS) and point (RMSE) performance for each dataset. We report mean metrics over all folds per method and indicate the standard deviation in brackets. Lower is better.

Dataset	folds	samples	features	CRPS		RMSE		
				PGBM	NGBoost	PGBM	NGBoost	LightGBM
yacht	20	308	6	0.22 (0.070)	0.32 (0.104)	0.63 (0.213)	0.75 (0.297)	0.64 (0.281)
boston	20	506	13	1.61 (0.201)	1.73 (0.236)	3.05 (0.507)	3.31 (0.661)	3.11 (0.675)
energy	20	768	8	0.21 (0.034)	0.25 (0.022)	0.35 (0.062)	0.49 (0.055)	0.29 (0.075)
concrete	20	1,030	8	2.06 (0.335)	2.95 (0.326)	3.97 (0.759)	5.50 (0.642)	3.80 (0.762)
wine	20	1,599	11	0.33 (0.034)	0.34 (0.024)	0.60 (0.054)	0.62 (0.043)	0.60 (0.050)
kin8nm	20	8,192	8	0.07 (0.002)	0.10 (0.002)	0.13 (0.005)	0.17 (0.003)	0.11 (0.003)
power	20	9,568	4	1.81 (0.053)	2.01 (0.120)	3.35 (0.153)	3.70 (0.222)	3.20 (0.140)
naval	20	11,934	14	0.00 (0.000)	0.00 (0.000)	0.00 (0.000)	0.00 (0.000)	0.00 (0.000)
protein	20	45,730	9	2.19 (0.030)	2.44 (0.038)	3.98 (0.056)	4.50 (0.059)	3.82 (0.058)
msd	1	515,345	90	4.78	4.75	9.09	9.16	9.11
higgs	1	10,500,000	28	0.253	0.238	0.418	0.419	0.414

Table 5: Key hyperparameters for the UCI benchmark and hierarchical time series experiment.

	UCI benchmark			Hierarchical time series		
	PGBM	NGBoost	LightGBM	PGBM	LightGBM	NGBoost
min_split_gain	0	0	0	0	0	0
min_data_in_leaf	1	1	1	1	1	1
max_bin	64	n.a.	64	1024	1024	n.a.
max_leaves	16	n.a.	16	64	64	64
max_depth	-1	3	-1	-1	-1	-1
learning_rate	0.1	0.01	0.1	0.1	0.1	0.1
n_estimators	2000	2000	2000	1000	1000	1000
feature_fraction	1.0	1.0	1.0	0.7	0.7	0.7
bagging_fraction	1.0	1.0	1.0	0.7	0.7	0.7
seed	1	1	1	1	1	1
lambda	1.0	n.a.	1.0	1.0	1.0	n.a.
early_stopping_rounds	n.a.	n.a.	n.a.	20	20	20

Table 6: Average time in seconds for running 2,000 iterations for each dataset on the UCI benchmark datasets. For msd and higgs, a bagging fraction of 0.1 was used.

Dataset	Probabilistic forecast			Point forecast
	PGBM-gpu	PGBM-cpu	NGBoost	LightGBM
wine (n=1,599)	100	41	20	1
naval (n=12k)	103	167	108	2
protein (n=46k)	115	355	389	2
msd (n=515k)	136	3,645	6,628	12
higgs (n=10,5M)	316	32,200	46,744	135

Table 7: M5 dataset description.

M5		
time series	#	3,049
time series description	item product sales	
target	\mathbb{R}^+	
train samples	#	2,415,359
validation samples	#	85,372
test samples	#	85,372
time step	t	day
features	#	48