

---

# Formatting Instructions For NeurIPS 2023

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The abstract paragraph should be indented 1/2 inch (3 picas) on both the left- and  
2 right-hand margins. Use 10 point type, with a vertical spacing (leading) of 11 points.  
3 The word **Abstract** must be centered, bold, and in point size 12. Two line spaces  
4 precede the abstract. The abstract must be limited to one paragraph.

## 5 1 Submission of papers to NeurIPS 2023

## 6 2 Introduction

## 7 3 Diffusion Models

8 In this section we quickly review the basics of diffusion models. We focus on the stochastic differential  
9 equation formulation first presented by (author?) [5].

10 Let  $p(\mathbf{y})_{\text{data}}$  denote the data distribution. The goal of a diffusion model is to learn a mapping from a  
11 simple distribution  $p(\mathbf{z})$  to the data distribution  $p(\mathbf{y})_{\text{data}}$ .

12 This is achieved by reversing a diffusion process. In particular, we construct a stochastic differential  
13 equation  $\mathbf{y}(t)$  from  $t \in [0, T]$  such that  $\mathbf{y}(0) \sim p(\mathbf{y})_{\text{data}}$  and  $\mathbf{y}(1) \sim p(\mathbf{y}(T))$  is a simple distribution  
14 we can sample from and whose evolution is given by

$$d\mathbf{y}(t) = \mathbf{f}(\mathbf{y}(t), t)dt + \mathbf{g}(t)d\mathbf{w}(t), \quad (1)$$

15 where  $\mathbf{w}(t)$  is a standard Brownian motion and  $\mathbf{f} : \mathbb{R}^d \times [0, T] \rightarrow \mathbb{R}^d$  and  $\mathbf{g} : [0, T] \rightarrow \mathbb{R}$  are called  
16 the drift coefficient and the diffusion coefficient, respectively.

17 It is possible to reverse this SDE and sample from  $p(\mathbf{y})_{\text{data}}$  by first sampling  $\mathbf{y}(1) \sim p(\mathbf{y}(T))$  and  
18 then evolving the system backwards in time. This is done by solving the reverse SDE (Cite anderson  
19 1982)

$$d\mathbf{y}(t) = [f(\mathbf{y}(t), t) - g(t)^2 \nabla_{\mathbf{y}(t)} \log p(\mathbf{y}(t))]dt + g(t)d\bar{\mathbf{w}}(t). \quad (2)$$

20 where  $\bar{\mathbf{w}}(t)$  is a standard Brownian with reversed time. Thus because  $\mathbf{f}$  and  $g$  are known, and we  
21 construct the SDE so that  $p(\mathbf{y}(T))$  is simple, as long as we know the score  $\nabla_{\mathbf{y}(t)} \log p(\mathbf{y}(t))$  we can  
22 sample from  $p(\mathbf{y})_{\text{data}}$ .

## 23 Estimating the Score

24 An important result by (author?) [6] is that it is possible to estimate the score  $\nabla_{\mathbf{y}(t)} \log p(\mathbf{y}(t))$  by  
25 computing

$$\mathbf{s}^* = \operatorname{argmin}_{\mathbf{s} \in \mathcal{S}} \mathbb{E}_t \mathbb{E}_{p(\mathbf{y}(0))_{\text{data}}} \mathbb{E}_{p(\mathbf{y}(t)|\mathbf{y}(0))} \left[ \left\| \nabla_{\mathbf{y}(t)} \log p(\mathbf{y}(t)|\mathbf{y}(0)) - \mathbf{s}(\mathbf{y}(t), t) \right\|^2 \right]. \quad (3)$$

where  $\mathcal{S} = \{s : \mathbb{R}^d \times [0, T] \rightarrow \mathbb{R}^d\}$  is the set of all possible score functions indexed by time  $t$ , and  $\mathbb{E}_t$  denotes the expectation over uniformly sampled  $t \in [0, T]$ .

## Conditional Diffusion Models

Although it is most common to train diffusion models unconditionally as explained above, one can also train diffusion models conditionally on some input  $\mathbf{x}$ .

To do so we make the following modifications to the above formulation.

1. We construct one separate SDE per value of  $\mathbf{x}$ . Each SDE shares the same drift and diffusion coefficients but the initial distribution  $p_{\mathbf{x}}(\mathbf{y}(0))$  is given by  $p(\mathbf{y}|\mathbf{x})_{\text{data}}$ .
2. The reverse SDE is now given by

$$d\mathbf{y}(t) = [f(\mathbf{y}(t), t) - g(t)^2 \nabla_{\mathbf{y}(t)} \log p_{\mathbf{x}}(\mathbf{y}(t))]dt + g(t)d\bar{\mathbf{w}}(t). \quad (4)$$

where  $\nabla_{\mathbf{y}(t)} \log p_{\mathbf{x}}(\mathbf{y}(t))$  is the score of the conditional distribution  $p_{\mathbf{x}}(\mathbf{y}(t))$ . Importantly, because we choose the diffusion and drift coefficients so that at  $t = T$  the distribution is the same for all values of  $\mathbf{x}$ , we can still sample from the data distribution in the same way as before.

3. The final change is that the score function is now estimated by

$$\mathbf{s}^* = \operatorname{argmin}_{\mathbf{s} \in \mathcal{S}} \mathbb{E}_t \mathbb{E}_{p(\mathbf{y}(0), \mathbf{x})_{\text{data}}} \mathbb{E}_{p(\mathbf{y}(t)|\mathbf{y}(0))} \left[ \|\nabla_{\mathbf{y}(t)} \log p(\mathbf{y}(t)|\mathbf{y}(0)) - \mathbf{s}(\mathbf{y}(t), \mathbf{x}, t)\|^2 \right].$$

with the changes being that now  $\mathcal{S} = \{s : \mathbb{R}^d \times \mathbb{R}^m \times [0, T] \rightarrow \mathbb{R}^d\}$  is the set of all possible score functions but now allowing for the score to depend on the input  $\mathbf{x}$ , and the expectation is taken over the joint distribution  $p(\mathbf{y}(0), \mathbf{x})_{\text{data}}$ . We emphasize that the score of the conditional distribution  $p(\mathbf{y}(t)|\mathbf{y}(0))$  is still the same because once we condition on  $\mathbf{y}(0)$  the distribution is the same for all values of  $\mathbf{x}$ .

This formulation of conditional diffusion models is different than controllable generation as presented in [5]. There, a conditional diffusion model is constructed by noting that

$$\nabla_{\mathbf{y}(t)} p(\mathbf{y}(t)|\mathbf{x}) = \nabla_{\mathbf{y}(t)} p(\mathbf{y}(t)) + \nabla_{\mathbf{y}(t)} p(\mathbf{x}|\mathbf{y}(t))$$

and hence if we obtain the first term from an unconditional diffusion model, and the second term by differentiating through another trained model  $p(\mathbf{x}|\mathbf{y}(t))$ , we can obtain the score of the conditional distribution. In our case this is not feasible because in general the dimension of  $\mathbf{y}$  will be much smaller than the dimension of  $\mathbf{x}$ .

## 4 Gradient Boosted Trees

Gradient Boosted Trees (GBT) [2] are a popular non-parametric machine learning model for function approximation. The objective is to find a function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  that minimizes

$$L(F) = \mathbb{E}_{\mathbf{x}, y} [l(y, F(\mathbf{x}))], \quad (5)$$

where  $l : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  is a loss function, and the expectation is taken over the joint distribution of the input  $\mathbf{x}$  and the target  $y$ . It does this by imposing the requirement that  $F$  is as a scaled sum of  $M$  decision trees  $f_m : \mathbb{R}^d \rightarrow \mathbb{R}$ , i.e.

$$F(\mathbf{x}) = \sum_{m=1}^M \epsilon f_m(\mathbf{x}), \quad \epsilon \in (0, 1). \quad (6)$$

where  $\epsilon$  is a learning rate or shrinkage parameter. In the most basic form of the algorithm each tree is constructed to approximate gradient descent on the loss function  $L(F)$ . In particular, if we let  $F_i = \sum_{m=1}^i f_m$  denote the function after  $i$  iterations and then the  $i$ -th tree is constructed to approximately minimize the squared error

$$f_i = \operatorname{argmin}_f \mathbb{E}_{\mathbf{x}, y} \left( f(\mathbf{x}) - \frac{\partial l(y, \hat{y})}{\partial \hat{y}} \bigg|_{\hat{y}=F_i(\mathbf{x})} \right)^2. \quad (7)$$

61 using empirical risk minimization and a greedy algorithm to construct the tree.  
62 Various modifications to the basic algorithm have been proposed and implemented such as reg-  
63 ularization, special ways of optimizing the tree, support for categorical functions, higher order  
64 optimization[1, 3, 4]. In this paper we focus on the implementation of GBTs in the LightGBM library  
65 [3] with the understanding that the same principles apply to any other GBT implementation.

## 66 5 Treeffuser/Treeffusion Models

67 For  $\mathbf{x} \in \mathbb{R}^d$ ,  $\mathbf{y} \in \mathbb{R}^m$  the objective of probabilistic predictions is to produce an estimate of the full  
68 conditional distribution  $\mathbb{P}[\mathbf{y}|\mathbf{x}]$ . This objective is different than in standard regression where the goal  
69 is usually to predict  $\mathbb{E}[\mathbf{y}|\mathbf{x}]$ .

70 The most common approach to solve this problem is via parametric models. This procedure assumes  
71 that the distribution  $\mathbb{P}[\mathbf{y}|\mathbf{x}]$  can be well approximated by a parametric family of distributions

$$\mathbb{P}[\mathbf{y}|\mathbf{x}] = p[\mathbf{y}|\theta(\mathbf{x})],$$

72 where  $p$  is a well known distribution (e.g Gaussian) and  $\theta(\mathbf{x})$  is a function that maps  $\mathbf{x}$  to the  
73 parameters of the distribution  $p$  (e.g. the mean and covariance of a Gaussian). Optimization is then  
74 performed by finding the function  $\theta(\mathbf{x})$  that minimizes a proper-scoring rule such as the negative  
75 log-likelihood.

76 The advantage of this approach is that, because  $p$  is often a simple distribution, it is easy to sample  
77 from it, compute the log-likelihood, and evaluate its moments. However, if the assumption that  $p$  is a  
78 good approximation to  $\mathbb{P}[\mathbf{y}|\mathbf{x}]$  is not met, the model can be poorly calibrated, the predictions can be  
79 inaccurate and the uncertainty estimates can be unreliable. Moreover, it is often hard to know when  
80 these assumptions are met, and it often takes significant expertise to find an appropriate parametric  
81 family of distributions.

82 Non-parametric models, like diffusions have become the state-of-the-art for deep generative models.  
83 However, they are not as popular for probabilistic predictions and are often not considered for this  
84 task, specially when the probabilistic predictions are over tabular data. Coupled with the fact that  
85 diffusion models are often computationally expensive to train, and coding scratch requires significant  
86 expertise, it is not surprising that they have not become the standard for probabilistic predictions.

87 To close this gap we propose a new class of models that we call Treeffuser which adapts diffusion  
88 models to the task of probabilistic predictions using gradient boosted trees. The benefit of this  
89 approach are:

- 90 1. Fast training: Gradient boosted trees are fast to train and can be trained on large datasets.  
91 Training on CPUs can be done in a matter of seconds.
- 92 2. Easy to use: Because the model is non-parametric, it is not necessary to tune hyperparameters  
93 of the distribution  $p$ . This means that practioners don't need to have a deep understanding of  
94 the distribution  $p$  to use the model.
- 95 3. Large sample convergence guarantees: Due to the non-parametric nature of the model and  
96 the score-based formulation of the training, the model is guaranteed to converge to the true  
97 distribution as the number of samples goes to infinity.
- 98 4. Efficient support for probabilistic predictions over multi-dimensional outputs: The model  
99 can be used to predict the full conditional distribution  $\mathbb{P}[\mathbf{y}|\mathbf{x}]$  of multi-dimensional outputs  
100 in a manner that scales  $O(m)$  with the number of dimensions  $m$ . This is better than the  
101  $O(m^2)$  scaling of parametric models. that require the estimation of the covariance matrix.
- 102 5. Support for feautres important for tabular data: Native support for training with categorical  
103 variables, missing values, etc.

### 104 5.1 Treeffuser Model

105 Treeffuser is an algorithm that combines conditional diffusion with gradient boosted trees. The model  
106 is constructed by approximating the score function  $\nabla_{\mathbf{y}(t)} \log p(\mathbf{y}(t)|\mathbf{x})$  with  $m$  gradient boosted trees  
107  $s_i : \mathbb{R}^d \times \mathbb{R}^M \times [0, T] \rightarrow \mathbb{R}$  that depend on the input  $\mathbf{x}, t, \mathbf{y}(t)$ .

$$\mathbf{s}(\mathbf{y}(t), \mathbf{x}, t) = (s_1(\mathbf{y}(t), \mathbf{x}, t), \dots, s_m(\mathbf{y}(t), \mathbf{x}, t)).$$

108 **Training:** Each gradient boosted tree  $s_i$  is trained to minimize the loss function

$$L_i(s) = \mathbb{E}_t \mathbb{E}_{\mathbf{y}(0), \mathbf{x}} \mathbb{E}_{\mathbf{y}(t) | \mathbf{y}(0)} \left[ \left( \frac{\partial \log p(\mathbf{y}(t) | \mathbf{x})}{\partial y_i(t)} - s(\mathbf{y}(t), \mathbf{x}, t) \right)^2 \right]. \quad (8)$$

109 Add the term that we divide by using empirical risk minimization. Due to [6] we know that if the  
110 gradient boosted tree adequately minimizes this objective it is guaranteed to converge to the true score  
111 function. Algorithm ?? describes the training procedure for the Treeffuser model in more detail.

112 **Sampling:** Maybe add something about sampling from the model.

113 **Log-likelihood:** Describe the computation of the log-likelihood.

114 **Other quantities:** As sampling from the model is very easy it is also easy to compute other  
115 quantities via monte carlo integration and in principle it is straightforward to compute any quantity  
116 of the form  $\mathbb{E}_{\mathbf{y} \sim p(\mathbf{y} | \mathbf{x})} f(\mathbf{y})$ . This includes the mean, variance, quantiles, but also more complicated  
117 quantities where a flexible model is needed. We demonstrate this in the experiments using the  
118 estimator for causal effects estimation.

## 119 References

120 References follow the acknowledgments in the camera-ready paper. Use unnumbered first-level  
121 heading for the references. Any choice of citation style is acceptable as long as you are consistent. It  
122 is permissible to reduce the font size to small (9 point) when listing the references. Note that the  
123 Reference section does not count towards the page limit.(author?) [5]

## 124 References

- 125 [1] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of*  
126 *the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*,  
127 KDD '16. ACM, August 2016.
- 128 [2] Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of*  
129 *Statistics*, 29(5):1189–1232, 2001.
- 130 [3] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and  
131 Tie-Yan Liu. Lightgbm: A highly efficient gradient boosting decision tree. In I. Guyon, U. Von  
132 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances*  
133 *in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 134 [4] Liudmila Prokhorenkova, Gleb Gusev, Aleksandr Vorobev, Anna Veronika Dorogush, and Andrey  
135 Gulin. Catboost: unbiased boosting with categorical features, 2019.
- 136 [5] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and  
137 Ben Poole. Score-based generative modeling through stochastic differential equations, 2021.
- 138 [6] Pascal Vincent. A connection between score matching and denoising autoencoders. Technical  
139 Report 1358, Département d'Informatique et de Recherche Opérationnelle, Université de Mon-  
140 tréal, CP 6128, Succ. Centre-Ville, Montréal (QC) H3C 3J7, Canada, December 2010. THIS IS  
141 A PREPRINT VERSION OF A NOTE THAT HAS BEEN ACCEPTED FOR PUBLICATION  
142 IN NEURAL COMPUTATION.