
Boosted Generative Models

Aditya Grover¹ Stefano Ermon¹

Abstract

We propose a new approach for using unsupervised boosting to create an ensemble of generative models, where models are trained in sequence to correct earlier mistakes. Our meta-algorithmic framework can leverage any existing base learner that permits likelihood evaluation, including recent latent variable models. Further, our approach allows the ensemble to include discriminative models trained to distinguish real data from model-generated data. We show theoretical conditions under which incorporating a new model in the ensemble will improve the fit and empirically demonstrate the effectiveness of boosting on density estimation and sample generation on synthetic and benchmark real datasets.

1. Introduction

A variety of deep generative models have recently shown promising results in tasks spanning computer vision, speech recognition, and natural language processing (Oord et al., 2016; Kingma & Welling, 2014; Goodfellow et al., 2014). Despite significant progress, existing generative models cannot fit complex distributions with a sufficiently high degree of accuracy, limiting their applicability.

In this paper, we propose a technique for ensembling (imperfect) generative models to improve their overall performance. Our meta-algorithm is inspired by boosting, a technique used in supervised learning to combine weak classifiers (e.g., decision stumps or trees), which individually might not perform well on a given classification task, into a more powerful ensemble. The boosting algorithm will attempt to learn a classifier to correct for the mistakes made by reweighting the original dataset, and repeat this procedure recursively. Under some conditions on the weak classifiers’ effectiveness, this procedure can drive the (training) error to zero (Freund et al., 1999).

¹Stanford University, USA. Correspondence to: Aditya Grover <adityag@cs.stanford.edu>.

We show that a similar procedure can be applied to generative models. Given an initial generative model that provides an imperfect fit to the data distribution, we construct a second model to correct for the error, and repeat recursively. The second model is also a generative one, which is trained on a reweighted version of the original training set. Our meta-algorithm is general and can construct ensembles of any existing generative model that permits (approximate) likelihood evaluation such as fully-visible belief networks and variational autoencoders. Interestingly, our method can also leverage powerful discriminative models. Specifically, we train a binary classifier to distinguish true data samples from “fake” ones generated by the current model and provide a principled way to include this discriminator in the ensemble.

A prior attempt at boosting density estimation proposed a *sum-of-experts* formulation (Rosset & Segal, 2002). We show theoretically and empirically the limitations of additive approaches in learning complex distributions, and instead propose to use multiplicative boosting or equivalently a *product-of-experts* formulation. Building on the proposed formulation, this paper makes the following contributions:

1. We provide theoretical conditions under which incorporating a new model is guaranteed to improve the ensemble fit, eventually recovering the true distribution.
2. We design and analyze a flexible meta-algorithmic boosting framework for including both generative and discriminative models in the ensemble.
3. We provide an algorithm for boosting *normalizing flow* models that permits exact and efficient likelihood evaluation and sampling.
4. We empirically demonstrate the effectiveness of boosted generative models on density estimation and sample generation on mixture of Gaussians, MNIST, and CIFAR-10 datasets.

2. Unsupervised boosting

Supervised boosting provides an algorithmic formalization of the hypothesis that a sequence of weak learners can create a single strong learner (Schapire & Freund, 2012). Here, we propose a framework that extends boosting to

unsupervised settings for learning generative models. For ease of presentation, all distributions are w.r.t. any arbitrary $\mathbf{x} \in \mathbb{R}^d$, unless otherwise specified. We use upper-case symbols to denote probability distributions and assume they all admit absolutely continuous densities (denoted by the corresponding lower-case notation) on a reference measure $d\mathbf{x}$. Our analysis naturally extends to discrete distributions, which we skip for brevity. Please refer to Appendix A for the proofs of all results in the following two sections.

Formally, we consider the following maximum likelihood estimation (MLE) setting. Given some data points $X = \{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^m$ sampled i.i.d. from an unknown distribution P , we provide a model class \mathcal{Q} parameterizing the distributions that can be represented by the generative model and minimize the Kullback-Liebler (KL) divergence w.r.t. the true distribution,

$$\min_{Q \in \mathcal{Q}} D_{KL}(P||Q). \quad (1)$$

In practice, we only observe samples from P and hence, maximize the log-likelihood of the observed data X . Selecting the model class for maximum likelihood learning is non-trivial; the MLE w.r.t. a small class can be far from the true distribution, whereas a large class poses the risk of overfitting in the absence of sufficient data, or even underfitting due to difficulty in optimizing non-convex objectives that frequently arise due to the use of latent variable models, neural networks, etc.

The boosting intuition is to greedily increase model capacity by learning a sequence of weak intermediate models $\{h_t \in \mathcal{H}_t\}_{t=0}^T$ that can correct for mistakes made by previous models in the ensemble. Here, \mathcal{H}_t is a predefined model class (such as \mathcal{Q}) for h_t . We defer the algorithms pertaining to the learning of such intermediate models to the next section, and first discuss two mechanisms for deriving the final estimate q_T from the individual density estimates at each round, $\{h_t\}_{t=0}^T$.

2.1. Additive boosting

In additive boosting, the final density estimate is an arithmetic average of the intermediate models,

$$q_T = \sum_{t=0}^T \alpha_t \cdot h_t$$

where $0 \leq \alpha_t \leq 1$ denote the weights assigned to the intermediate models. The weights are re-normalized at every round to sum to 1 which gives us a valid probability density estimate. Starting with a base model h_0 , we can express the density estimate after a round of boosting recursively as,

$$q_t = (1 - \hat{\alpha}_t) \cdot q_{t-1} + \hat{\alpha}_t \cdot h_t$$

where $\hat{\alpha}_t$ denotes the normalized weight for h_t at round t . We now derive conditions on the intermediate models that guarantee “progress” in every round of boosting.

Theorem 1. *Let $\delta_{KL}^t(h_t, \hat{\alpha}_t) = D_{KL}(P||Q_{t-1}) - D_{KL}(P||Q_t)$ denote the reduction in KL-divergence at the t^{th} round of additive boosting. The following conditions hold,*

1. *Sufficient: If $\mathbb{E}_P \left[\log \frac{h_t}{q_{t-1}} \right] \geq 0$, then $\delta_{KL}^t(h_t, \hat{\alpha}_t) \geq 0$ for all $\hat{\alpha}_t \in [0, 1]$.*
2. *Necessary: If $\exists \hat{\alpha}_t \in (0, 1]$ such that $\delta_{KL}^t(h_t, \hat{\alpha}_t) \geq 0$, then $\mathbb{E}_P \left[\frac{h_t}{q_{t-1}} \right] \geq 1$.*

The sufficient and necessary conditions require that the expected log-likelihood and likelihood respectively of the current intermediate model, h_t are better-or-equal than those of the combined previous model, q_{t-1} under the true distribution when compared using density ratios. Learning such an intermediate model at every round is algorithmically difficult and furthermore, counterintuitive to the boosting paradigm which combines the strength of *weak* intermediate learners. Next, we consider an alternative formulation that eases the requirements for learning good intermediate models.

2.2. Multiplicative boosting

In multiplicative boosting, we factorize the final density estimate as a geometric average of $T + 1$ intermediate models $\{h_t\}_{t=0}^T$, each assigned an exponentiated weight α_t ,

$$q_T = \frac{\prod_{t=0}^T h_t^{\alpha_t}}{Z_T}$$

where the partition function $Z_T = \int \prod_{t=0}^T h_t^{\alpha_t} d\mathbf{x}$. Recursively, we have the density estimate as,

$$\tilde{q}_t = h_t^{\alpha_t} \cdot \tilde{q}_{t-1} \quad (2)$$

where \tilde{q}_t is the unnormalized estimate (at round t). The base model h_0 is learned using MLE. The conditions on the intermediate models for reducing KL-divergence at every round are stated below.

Theorem 2. *Let $\delta_{KL}^t(h_t, \alpha_t) = D_{KL}(P||Q_{t-1}) - D_{KL}(P||Q_t)$ denote the reduction in KL-divergence at the t^{th} round of multiplicative boosting. The following conditions hold,*

1. *Sufficient: If $\mathbb{E}_P[\log h_t] \geq \log \mathbb{E}_{Q_{t-1}}[h_t]$, then $\delta_{KL}^t(h_t, \alpha_t) \geq 0$ for all $\alpha_t \in [0, 1]$.*
2. *Necessary: If $\exists \alpha_t \in (0, 1]$ such that $\delta_{KL}^t(h_t, \alpha_t) \geq 0$, then $\mathbb{E}_P[\log h_t] \geq \mathbb{E}_{Q_{t-1}}[\log h_t]$.*

Algorithm 1 GenBGM($X = \{\mathbf{x}_i\}_{i=1}^m, T$ rounds)

Initialize $d_0(\mathbf{x}_i) = 1/m$ for all $i = 1, 2, \dots, m$.
 Obtain base generative model h_0 .
 Set (unnormalized) density estimate $\tilde{q}_0 = h_0$
for $t = 1, 2, \dots, T$ **do**
 - Choose β_t and update d_t using Eq. (4).
 - Train generative model h_t to maximize Eq. (3).
 - Choose α_t .
 - Set density estimate $\tilde{q}_t = \tilde{q}_{t-1} \cdot h_t^{\alpha_t}$.
end for
 Estimate $Z_T = \int \tilde{q}_T d\mathbf{x}$.
 return $q_T = \tilde{q}_T / Z_T$.

Algorithm 2 DiscBGM($X = \{\mathbf{x}_i\}_{i=1}^m, T$ rounds, f)

Initialize $d_0(\mathbf{x}_i) = 1/m$ for all $i = 1, 2, \dots, m$.
 Obtain base generative model h_0 .
 Set (unnormalized) density estimate $\tilde{q}_0 = h_0$
for $t = 1, \dots, T$ **do**
 - Generate negative samples from q_{t-1}
 - Optimize r_t to maximize RHS in Eq.(5).
 - Set $h_t = [f']^{-1}(r_t)$.
 - Choose α_t .
 - Set density estimate $\tilde{q}_t = h_t^{\alpha_t} \cdot \tilde{q}_{t-1}$.
end for
 Estimate $Z_T = \int \tilde{q}_T d\mathbf{x}$.
 return $q_T = \tilde{q}_T / Z_T$.

In contrast to additive boosting, the conditions above compare expectations under the true distribution with expectations under the *model distribution* in the previous round, Q_{t-1} . The equality in the conditions holds for $\alpha_t = 0$, which corresponds to the trivial case where the current intermediate model is ignored in Eq. (2). For other valid α_t , the non-degenerate version of the sufficient inequality guarantees progress towards the true data distribution. Note that the intermediate models increase the overall capacity of the ensemble at every round.

From the necessary condition, we see that a “good” intermediate model h_t necessarily assigns a better-or-equal log-likelihood under the true distribution as opposed to the model distribution, Q_{t-1} . This condition suggests two learning objectives for intermediate models which we discuss next.

3. Boosted generative models

In this section, we design and analyze meta-algorithms for multiplicative boosting of generative models. Given any base model which permits (approximate) likelihood evaluation, we provide a mechanism for boosting this model using an ensemble of generative and/or discriminative models.

3.1. Generative boosting

Supervised boosting algorithms such as AdaBoost (Freund & Schapire, 1995) typically involve a reweighting procedure for training weak learners. We can similarly train an ensemble of generative models for unsupervised boosting, where every subsequent model performs MLE w.r.t a reweighted data distribution D_t ,

$$\max_{h_t} \mathbb{E}_{D_t} [\log h_t] \quad (3)$$

$$\text{where } d_t \propto \left(\frac{p}{q_{t-1}} \right)^{\beta_t} \quad (4)$$

and $\beta_t \in [0, 1]$ is the reweighting coefficient (at round t). Note that these coefficients are in general different from the model weights α_t that appear in the density estimate in Eq. (2).

Proposition 1. *If we can maximize the objective in Eq. (3) optimally, then $\delta_{KL}^t(h_t, \alpha_t) \geq 0$ for any $\beta_t \in [0, 1]$ with the equality holding for $\beta_t = 0$.*

While the objective in Eq. (3) can be hard to optimize in practice, the target distribution becomes easier to approximate as we reduce the reweighting coefficient. For the extreme case of $\beta_t = 0$, the reweighted data distribution is simply uniform. There is no free lunch however, since a low β_t results in a slower reduction in KL-divergence leading to a computational-statistical trade-off.

The pseudocode for the corresponding boosting meta-algorithm, referred to as GenBGM, is given in Algorithm 1. In practice, we only observe samples from the true data distribution, and hence, approximate p based on the empirical data distribution which is defined to be uniform over the dataset X . At every subsequent round, GenBGM learns an intermediate model that maximizes the log-likelihood of data sampled from a reweighted data distribution.

3.2. Discriminative boosting

A base generative model can be boosted using a discriminative approach as well. Here, the intermediate model is specified as the density ratio obtained from a binary classifier. Consider the following setup: we observe an equal number of samples drawn i.i.d. from the true data distribution (w.l.o.g. assigned the label $y = +1$) and the model distribution in the previous round Q_{t-1} (label $y = -1$).

Definition 1. *Let $f : \mathbb{R}^+ \rightarrow \mathbb{R}$ be any convex, lower semi-continuous function satisfying $f(1) = 0$. The f -divergence between P and Q is defined as,*

$$D_f(P||Q) = \int q \cdot f\left(\frac{p}{q}\right) d\mathbf{x}$$

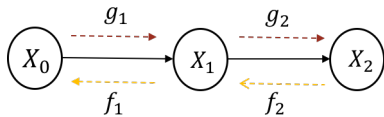


Figure 1: The normalizing flow model for $T+1 = 3$ layers.

Notable examples include the Kullback-Liebler (KL) divergence, Hellinger distance and the Jensen-Shannon (JS) divergence among many others. The binary classifier in discriminative boosting maximizes a variational lower bound on any f -divergence at round t ,

$$D_f(P\|Q_{t-1}) \geq \sup_{r_t \in \mathcal{R}_t} (\mathbb{E}_P[r_t] - \mathbb{E}_{Q_{t-1}}[f^*(r_t)]). \quad (5)$$

where f^* denotes the Fenchel conjugate of f and $r_t : \mathbb{R}^d \rightarrow \text{dom}_{f^*}$ parameterizes the classifier. Under mild conditions on f (Nguyen et al., 2010), the lower bound in Eq. (5) is tight if $r_t^* = f' \left(\frac{p}{q_{t-1}} \right)$.

Hence, a solution to Eq. (5) can be used to estimate density ratios. The density ratios naturally fit into the multiplicative boosting framework and provide a justification for the use of objectives of the form Eq. (5) for learning intermediate models as formalized in the proposition below.

Proposition 2. *For any given f -divergence, let r_t^* denote the optimal solution to Eq. (5) in the t^{th} round of boosting. Then, the model density at the end of the boosting round matches the true density if we set $\alpha_t = 1$ and $h_t = [f']^{-1}(r_t^*)$ where $[f']^{-1}$ denotes the inverse of the derivative of f .*

The pseudocode for the corresponding meta-algorithm, DiscBGM is given in Algorithm 2. At every round, we train a binary classifier to optimize the objective in Eq. (5) for a chosen f -divergence. As a special case, the negative of the cross-entropy loss commonly used for binary classification is also a lower bound on an f -divergence. While Algorithm 2 is applicable for any f -divergence, we will focus on the cross-entropy based objective henceforth to streamline the discussion.

Corollary 1. *Consider the (negative) cross-entropy objective maximized by a binary classifier,*

$$\sup_{c_t \in \mathcal{C}_t} \mathbb{E}_P[\log c_t] + \mathbb{E}_{Q_{t-1}}[\log(1 - c_t)]. \quad (6)$$

If a binary classifier c_t trained to optimize Eq. (6) is Bayes optimal, then the model density at the end of the boosting round matches the true density if we set $\alpha_t = 1$ and

$$h_t = \frac{c_t}{1 - c_t}. \quad (7)$$

In practice, a classifier with limited capacity trained on a finite dataset will not generally be Bayes optimal. The above

corollary, however, suggests that a good classifier can provide a ‘direction of improvement’, in a similar spirit to the gradient boosting algorithm for supervised learning (Freund & Schapire, 1995). Additionally, if the intermediate model distribution h_t obtained using Eq. (7) satisfies the conditions in Theorem 2, it is guaranteed to improve the fit.

The weights $\alpha_t \in [0, 1]$ can be interpreted as our confidence in the classifier density estimate, akin to the step size used in gradient descent. While in practice we use heuristic strategies for assigning weights to the intermediate models, the greedy optimum value of these weights at every round is a critical point for δ_{KL}^t (defined in Theorem 2). For example, in the extreme case where c_t is uninformative, *i.e.*, $c_t \equiv 0.5$, then $\delta_{KL}^t(h_t, \alpha_t) = 0$ for all $\alpha_t \in [0, 1]$. If c_t is Bayes optimal, then δ_{KL}^t attains a maxima when $\alpha_t = 1$ (Corollary 1).

3.3. Hybrid boosting

Intermediate models need not be exclusively generators or discriminators; we can design a boosting ensemble with any combination of generators and discriminators as intermediate models. If an intermediate model is chosen to be a generator, we learn a generative model using MLE after appropriately reweighting the data points. If a discriminator is used to implicitly specify an intermediate model, we set up the corresponding binary classification problem.

3.4. Regularization

In practice, we want boosted generative models (BGM) to generalize to data points outside the training set X . Regularization in BGMs is imposed primarily in two ways. First, every intermediate model can be independently regularized by early stopping of training based on validation error, incorporating explicit terms in the learning objective, heuristics such as dropout, etc. Moreover, restricting the number of rounds of boosting is another effective mechanism for regularizing BGMs. Fewer rounds of boosting are required if the intermediate models are sufficiently expressive.

3.5. Boosting normalized flow models

The boosting meta-algorithms presented above are applicable to a wide range of base generators learned using MLE as well as discriminators maximizing a variational lower bound on a chosen f -divergence. At any given round, the model density needs to be specified only up to a normalization constant. While many applications of generative modeling such as feature learning can sidestep computing the partition function, if needed it can be estimated using techniques such as Annealed Importance Sampling (Neal, 2001). Similarly Markov chain based methods can be used

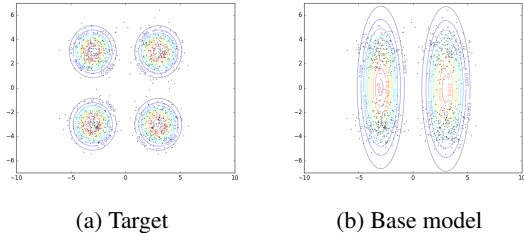


Figure 2: The mixture of Gaussians setup showing (a) true density and (b) base (misspecified) model.

to generate samples.

To offset these limitations, consider the normalizing flow model in Figure 1. The model consists of multiple layers of random variables $\{X_t\}_{t=0}^T$, with each layer expressing a distribution q_t over the same space as the data. The connections between the layers $t-1 \rightarrow t$ specify an invertible transformation $g_t = f_t^{-1}$ from one distribution to another using the change-of-variables formula,

$$q_t(\mathbf{x}) = q_{t-1}(f_t(\mathbf{x})) \left| \det \left(\frac{\partial f_t(\mathbf{x})}{\partial \mathbf{x}} \right) \right| \quad (8)$$

where $\frac{\partial f_t(\mathbf{x})}{\partial \mathbf{x}}$ denotes the Jacobian of f_t at \mathbf{x} . Normalizing flow models permit *exact* likelihood evaluation as long as the prior density of the first layer is tractable and computing the determinant of the Jacobian is inexpensive. As in (Dinh et al., 2014) we use affine location-scale transformations parameterized by neural networks that ensure that the Jacobian is an upper triangular matrix. Furthermore, efficient ancestral sampling is possible in flow models (Dinh et al., 2014).

In the boosting setting, the connections between any two layers of a normalizing flow model parametrize a weak learner that corrects for the shortcomings of the previous layers. The (recursive) change-of-variables formula in Eq. (8) is akin to a special form of multiplicative boosting with all model weights α_t set to unity and the intermediate model at any round, h_t corresponding to the absolute value of the determinant of the Jacobian of a parameterized invertible function f_t . We defer the theoretical analysis for multiplicative boosting in flow models to Appendix A.6. Importantly, flow models ensure that the resulting density estimate is normalized. Similar to Algorithm 1, we can learn a flow model with a fixed number of layers as a base generative model and greedily train additional layers to optimize the objective in Eq. (3), following an appropriate reweighting step as per Eq. (4). End-to-end training of very deep networks through backpropagation suffers from a decreasing signal-to-noise ratio due to vanishing or exploding gradients (Glorot & Bengio, 2010; Shalev-Shwartz et al., 2017). Greedy training offsets this limitation empirically,

Table 1: Average test NLL (with std error) for mixture of Gaussians.

Model	NLL (in nats)
Base model	4.69 \pm 0.01
Add model	4.64 \pm 0.02
GenBGM	4.58 \pm 0.10
DiscBGM-NCE	4.42 \pm 0.01
DiscBGM-HD	4.35 \pm 0.01

Table 2: Test NLL for the CIFAR-10 dataset.

Model	NLL (in bits/dim)
Joint	3.58
$\beta_2 = 0$	3.56
$\beta_2 = 0.25$	3.55
$\beta_2 = 0.50$	3.55
$\beta_2 = 1.0$	3.57

and reweighting data points additionally lets us generalize better.

4. Empirical evaluation

We evaluated the performance of boosting generative models for the tasks of density estimation and sample generation on real and synthetic datasets.

4.1. Mixture of Gaussians

Experimental setup. The true data distribution is a equi-weighted mixture of four Gaussians centered symmetrically around the origin, each having an identity covariance matrix. The contours of the underlying density are shown in Figure 2 (a). We only observe 1,000 training samples drawn i.i.d. from the data distribution (shown as black dots in Figure 3), and the task is to learn this distribution. The test set contains 1,000 samples from the same distribution. We repeat the process 10 times.

As a base (misspecified) model, we fit a mixture of two Gaussians to the data; the contours for an example instance are shown in Figure 2 (b). We compare multiplicative and additive boosting, each run for $T = 2$ rounds. For additive boosting (Add), we use the algorithm proposed in (Rosset & Segal, 2002) setting $\hat{\alpha}_0$ to unity and doing a line search over $\hat{\alpha}_1, \hat{\alpha}_2 \in [0, 1]$. For the multiplicative boosting algorithms, all model weights, α 's to unity. The reweighting coefficients, β 's for GenBGM are all set to unity and the intermediate models are mixtures of two Gaussians as well. For DiscBGM, the classifiers are multi-layer perceptrons

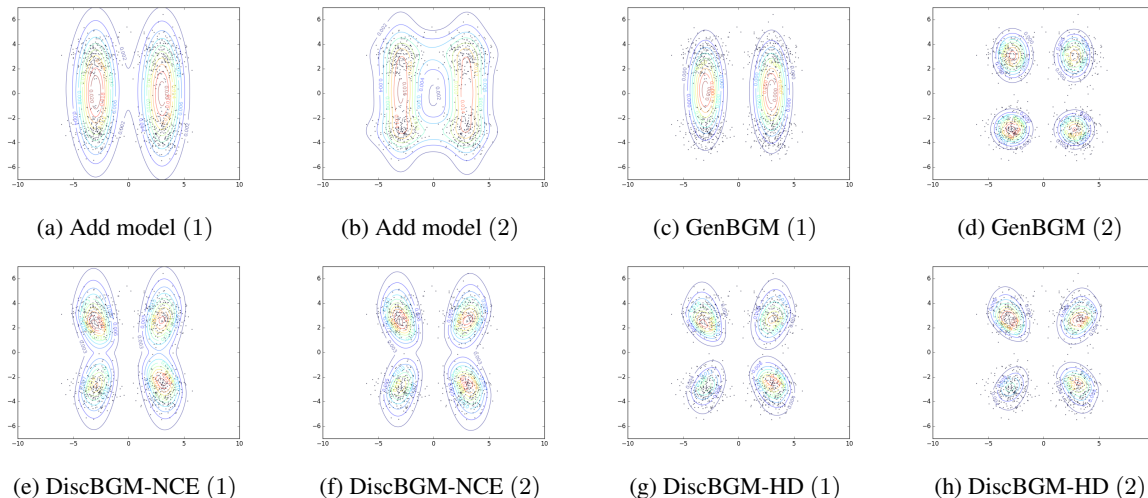


Figure 3: Multiplicative boosting algorithms such as GenBGM (c-d) and DiscBGM with negative cross-entropy (e-f) and Hellinger distance (g-h) outperform additive boosting (a-b) in correcting for model misspecification. Numbers in parenthesis indicate boosting round t .

with two hidden layers (each containing 100 units) maximizing f -divergences corresponding to the negative cross-entropy (NCE) and Hellinger distance (HD).

The test negative log-likelihood (NLL) estimates are listed in Table 1. Qualitatively, the contour plots for the estimated densities after every boosting round on a sample instance are shown in Figure 3. Multiplicative boosting algorithms outperform additive boosting in correcting for model misspecification. GenBGM initially leans towards maximizing coverage, whereas both versions of DiscBGM are relatively more conservative in assigning high densities to data points away from the modes. Refer to Appendix B.1 for results with additional heuristic model weighting (α 's) strategies.

4.2. MNIST

Experimental setup. Consider a baseline variational autoencoder (VAE) (Kingma & Welling, 2014) trained for sufficiently long on the binarized MNIST dataset (LeCun et al., 2010). Ancestral samples obtained by the baseline VAE model are shown in Figure 4 (a). We use the evidence lower bound (ELBO) as a proxy for approximately evaluating the marginal log-likelihood during learning.

The conventional approach to improving the performance of a latent variable model is to increase its representational capacity by adding hidden layers (Baseline VAE + depth) or increasing the number of hidden units in the existing layers (Baseline VAE + width). These lead to a marginal improvement in sample quality as seen in Figure 4 (b) and 4 (c). Table 4.1 lists the architecture for the baseline models. Appendix B.2 provides more details regarding the learning procedure and hyperparameters.

In contrast, boosting makes steady improvements in sample quality. We start with a VAE roughly half the size of Baseline VAE and generate samples after boosting it using GenBGM, DiscBGM, and GenDiscBGM (Figures 4 (d), 4 (e), and 4 (f) respectively). The discriminator used is a convolutional neural network (CNN) (LeCun & Bengio, 1995) trained to maximize the negative cross-entropy. The model weights, α 's and reweighting coefficient, β 's are set to unity. The boosted sequences generate sharper samples than all baselines in spite of having similar model capacity. The samples are generated using independent Monte Carlo Markov Chain (MCMC) sampling. See Appendix B.2 for further details.

Since boosting is particularly attractive for improving *weak* learners, it does not require the models in the ensemble to be trained until convergence. Tables 4.1 and 4.2 in Figure 4 compare the wall-clock time taken during training of baseline models and BGM sequences. In an attempt to equalize for training time, we train the VAEs used in the BGM sequences for only 70 epochs, such that the training time of the most expensive BGM sequence (GenDiscBGM) matches that of the weakest baseline (Baseline VAE). The convergence curves are shown in Appendix B.2 for reference. Hence, the boosting framework does not contribute towards any significant overhead during training.

4.3. CIFAR-10

Experimental setup. We boost a normalized flow model trained on the CIFAR-10 dataset (Krizhevsky & Hinton, 2009). The base model is a multi-scale architecture similar to the one proposed by (Dinh et al., 2016), marginally

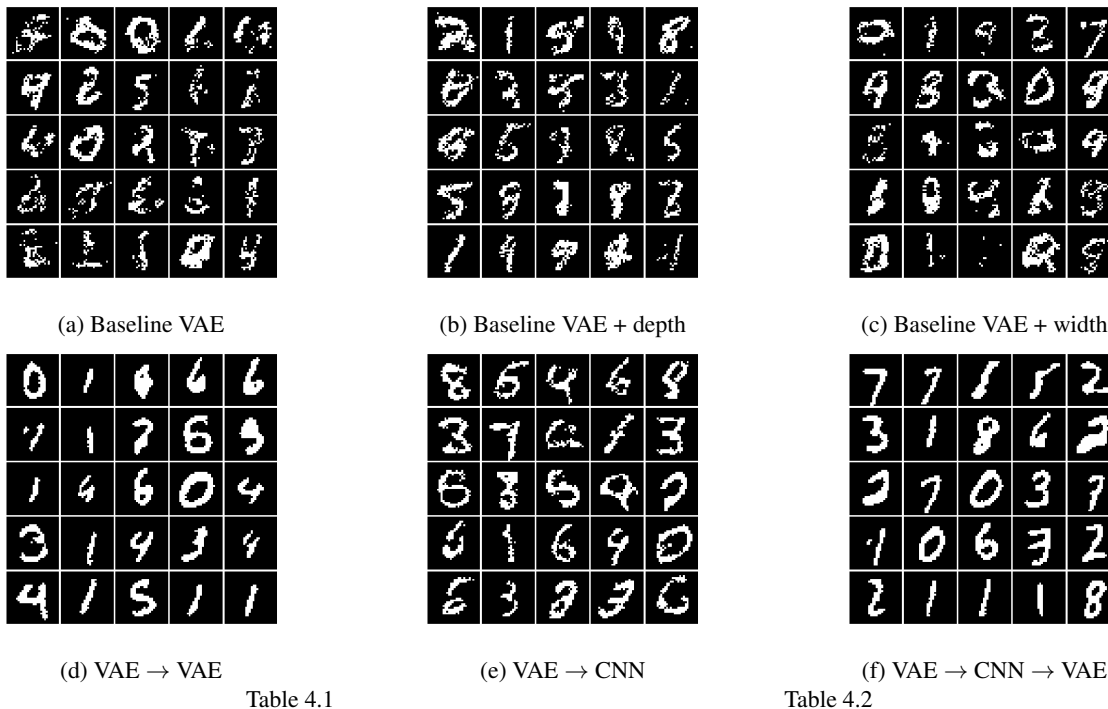


Figure 4: The boosted generative models (d-f) demonstrate how ensembles of weak learners can generate sharper samples, compared to naively increasing model capacity (a-c). Note that we show actual samples of binarized digits and not mean values for the pixels. The VAEs used in (d-f) have a 784-100-50 layered architecture.

downsized due to resource constraints. The architecture details and learning procedure are deferred to Appendix B.3.

The prior density, q_0 is an isotropic unit norm Gaussian. The base model implicitly specifies an intermediate learner h_1 . A second intermediate learner h_2 is obtained by adding another scale to the model. Learning the parameters for the final architecture (h_1 and h_2) can be done in two ways. The “Joint” baseline jointly learns the parameters based on MLE. Alternatively, we learn the model greedily (*i.e.*, h_1 followed by h_2) using the boosting procedure described in Section 3.5. The reweighting coefficient β_1 is set to zero (*i.e.*, all points have same weights) since the prior density is not expected to be very informative. For subsequent learning of h_2 , we do a line search for different values of β_2 .

The test NLLs are listed in Table 2. Following common practice, the results are expressed in bits and normalized by the total number of dimensions ($32 \times 32 \times 3$). We observe that boosted models are significantly better than the

jointly trained baseline (up to 90 bits over the full space). Hence, boosting indirectly provides a better optimization mechanism offsetting the low signal-to-noise ratio problem commonly observed in training very deep networks (Glorot & Bengio, 2010; Shalev-Shwartz et al., 2017). The samples generated by the models do not show any significant difference and are deferred to Appendix B.3.

5. Related work and discussion

Boosting, a meta-algorithmic framework initially proposed in the context of supervised learning arose in response to a seminal question posed by (Kearns & Valiant, 1994): can a set of weak learners create a strong learner? Boosting has offered interesting theoretical insights into the fundamental limits of supervised learning and led to the development of practical algorithms that work well in practice (Schapire, 1990; Freund et al., 1999; Friedman, 2002; Caruana & Niculescu-Mizil, 2006). In the supervised set-

ting, Lebanon & Lafferty (2002) have shown theoretical connections between the log-likelihood loss similar to the one used in our framework and the exponential loss minimized by AdaBoost (Freund et al., 1999). Our work provides a foundational framework for unsupervised boosting with connections to prior work discussed below.

Sum- vs. product-of-experts. (Rosset & Segal, 2002) proposed an algorithm for density estimation using Bayesian networks similar to gradient boosting. Theorem 1 highlights the key limitation of additive formulations for learning complex multi-modal distributions as requiring a more powerful learner at every round. These models are normalized and easy to sample, but are generally outperformed by multiplicative formulations in correcting for model misspecification in practice (Section 4.1). Recent preprints concurrent with this work use additive boosting in the context of improving approximate posteriors for variational inference (Miller et al., 2016; Guo et al., 2016) and generative adversarial networks (Tolstikhin et al., 2017). The latter shows improvements for a GAN-specific problem of mode dropping on a toy 2D problem.

The product-of-experts formulation, which we adopt in the current work, was initially proposed for feature learning in energy based models such as Boltzmann machines. For example, the hidden units in a restricted Boltzmann machine can be interpreted as weak learners performing MLE. If the number of weak learners is fixed, they can be efficiently updated in parallel but there is a risk of learning redundant features (Hinton, 1999; 2002). Weak learners can also be added incrementally based on the learner’s ability to distinguish observed data and model-generated data (Welling et al., 2002). (Tu, 2007) generalized the latter to boost arbitrary probabilistic models; their algorithm is a special case of DiscBGM with all α ’s set to 1 and the discriminator maximizing the negative cross-entropy. (Gutmann & Hirayama, 2011) employ a similar procedure for learning probabilistic models by minimizing Bregman divergences. DiscBGM additionally accounts for imperfections in learning classifiers through flexible model weights and can be adapted to maximize any f -divergence.

Unsupervised-as-supervised learning. The use of density ratios learned by a binary classifier for estimation was first proposed by (Friedman et al., 2001) and has been subsequently applied elsewhere, notably for parameter estimation using noise-contrastive estimation (Gutmann & Hyvärinen, 2010) and sample generation in generative adversarial networks (GAN) (Goodfellow et al., 2014).

Generative Adversarial Networks. GANs consist of a pair of generator-discriminator networks. The discriminator maximizes the negative cross-entropy as in Eq. (6) and the generator minimizes the same objective. GANs gener-

ate good samples, but are unstable to train in practice, and the training objective is not guaranteed to converge (Goodfellow, 2014). Borrowing terminology from (Mohamed & Lakshminarayanan, 2016), GANs can be seen as *implicit* probabilistic models that are directly optimizing for generating visually appealing samples as opposed to *prescribed* probabilistic models such as BGMs that explicitly characterize the log-likelihood and are typically stable to train. While the original GAN formulation was proposed for a cross-entropy based objective, (Nowozin et al., 2016) extend it to arbitrary f -divergences based on the same variational divergence minimization framework (Nguyen et al., 2010) we use for learning binary classifiers in boosted generative models.

6. Conclusion

We presented a general-purpose framework for boosting generative models by explicit factorization of the model likelihood as a product of simpler intermediate model densities. These intermediate models are learned greedily using discriminative or generative approaches, gradually increasing the overall model’s capacity. We further designed an algorithm for boosting normalizing flow models which can perform exact and efficient likelihood evaluation and sampling. We demonstrated the effectiveness of boosted generative models by designing several ensemble models which improve upon baseline generative models for the tasks of density estimation and sample generation without incurring any significant computational overhead.

In the future, we will explore the design of boosted models such as normalized flow models that permit efficient learning and inference. Further examination of the effect of various model weighting strategies and choice of f -divergences on learning is another interesting direction for future work.

References

- Caruana, Rich and Niculescu-Mizil, Alexandru. An empirical comparison of supervised learning algorithms. In *ICML*, 2006.
- Dinh, Laurent, Krueger, David, and Bengio, Yoshua. NICE: Non-linear independent components estimation. *arXiv preprint arXiv:1410.8516*, 2014.
- Dinh, Laurent, Sohl-Dickstein, Jascha, and Bengio, Samy. Density estimation using real NVP. *arXiv preprint arXiv:1605.08803*, 2016.
- Freund, Yoav and Schapire, Robert E. A decision-theoretic generalization of on-line learning and an application to boosting. In *European conference on computational learning theory*, 1995.
- Freund, Yoav, Schapire, Robert, and Abe, N. A short introduction to boosting. *Journal-Japanese Society For Artificial Intelligence*, 14(771-780):1612, 1999.
- Friedman, Jerome, Hastie, Trevor, and Tibshirani, Robert. *The elements of statistical learning*, volume 1. Springer series in statistics Springer, Berlin, 2001.
- Friedman, Jerome H. Stochastic gradient boosting. *Computational Statistics & Data Analysis*, 38(4):367–378, 2002.
- Glorot, Xavier and Bengio, Yoshua. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, 2010.
- Goodfellow, Ian, Pouget-Abadie, Jean, Mirza, Mehdi, Xu, Bing, Warde-Farley, David, Ozair, Sherjil, Courville, Aaron, and Bengio, Yoshua. Generative adversarial nets. In *NIPS*, 2014.
- Goodfellow, Ian J. On distinguishability criteria for estimating generative models. *arXiv preprint arXiv:1412.6515*, 2014.
- Guo, Fangjian, Wang, Xiangyu, Fan, Kai, Broderick, Tamara, and Dunson, David B. Boosting variational inference. *arXiv preprint arXiv:1611.05559*, 2016.
- Gutmann, Michael and Hirayama, Jun-ichiro. Bregman divergence as general framework to estimate unnormalized statistical models. In *UAI*, 2011.
- Gutmann, Michael and Hyvärinen, Aapo. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *AISTATS*, 2010.
- Hinton, Geoffrey E. Products of experts. In *ICANN*, 1999.
- Hinton, Geoffrey E. Training products of experts by minimizing contrastive divergence. *Neural computation*, 2002.
- Ioffe, Sergey and Szegedy, Christian. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- Kearns, Michael and Valiant, Leslie. Cryptographic limitations on learning boolean formulae and finite automata. *JACM*, 41(1):67–95, 1994.
- Kingma, Diederik and Ba, Jimmy. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- Kingma, Diederik P and Welling, Max. Auto-encoding variational bayes. In *ICLR*, 2014.
- Krizhevsky, Alex and Hinton, Geoffrey. Learning multiple layers of features from tiny images. 2009.
- Lebanon, Guy and Lafferty, John. Boosting and maximum likelihood for exponential models. In *NIPS*, 2002.
- LeCun, Yann and Bengio, Yoshua. Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks*, 3361(10):1995, 1995.
- LeCun, Yann, Cortes, Corinna, and Burges, Christopher JC. MNIST handwritten digit database. <http://yann.lecun.com/exdb/mnist>, 2010.
- Miller, Andrew C, Foti, Nicholas, and Adams, Ryan P. Variational boosting: Iteratively refining posterior approximations. *arXiv preprint arXiv:1611.06585*, 2016.
- Mohamed, Shakir and Lakshminarayanan, Balaji. Learning in implicit generative models. *arXiv preprint arXiv:1610.03483*, 2016.
- Neal, Radford M. Annealed importance sampling. *Statistics and Computing*, 11(2):125–139, 2001.
- Nguyen, XuanLong, Wainwright, Martin J, and Jordan, Michael I. Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*, 56(11):5847–5861, 2010.
- Nowozin, Sebastian, Cseke, Botond, and Tomioka, Ryota. f-GAN: Training generative neural samplers using variational divergence minimization. In *NIPS*, 2016.
- Oord, Aaron van den, Kalchbrenner, Nal, and Kavukcuoglu, Koray. Pixel recurrent neural networks. In *ICML*, 2016.
- Rosset, Saharon and Segal, Eran. Boosting density estimation. In *NIPS*, 2002.

- Salimans, Tim and Kingma, Diederik P. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In *NIPS*, 2016.
- Schapire, Robert E. The strength of weak learnability. *Machine learning*, 5(2):197–227, 1990.
- Schapire, Robert E and Freund, Yoav. *Boosting: Foundations and algorithms*. MIT press, 2012.
- Shalev-Shwartz, Shai, Shamir, Ohad, and Shammah, Shaked. Failures of deep learning. *arXiv preprint arXiv:1703.07950*, 2017.
- Tolstikhin, Ilya, Gelly, Sylvain, Bousquet, Olivier, Simon-Gabriel, Carl-Johann, and Schölkopf, Bernhard. Adagan: Boosting generative models. *arXiv preprint arXiv:1701.02386*, 2017.
- Tu, Zhuowen. Learning generative models via discriminative approaches. In *CVPR*, 2007.
- Uribe, Benigno, Murray, Iain, and Larochelle, Hugo. RNADE: The real-valued neural autoregressive density-estimator. In *NIPS*, 2013.
- Welling, Max, Zemel, Richard S, and Hinton, Geoffrey E. Self supervised boosting. In *NIPS*, 2002.

Appendices

A. Proofs of theoretical results

A.1. Theorem 1

The reduction in KL-divergence can be simplified as,

$$\begin{aligned}\delta_{KL}^t(h_t, \hat{\alpha}_t) &= \mathbb{E}_P \left[\log \frac{p}{q_{t-1}} \right] - \mathbb{E}_P \left[\log \frac{p}{q_t} \right] \\ &= \mathbb{E}_P \left[\log \frac{q_t}{q_{t-1}} \right] \\ &= \mathbb{E}_P \left[\log \left[(1 - \hat{\alpha}_t) + \hat{\alpha}_t \frac{h_t}{q_{t-1}} \right] \right].\end{aligned}$$

We first derive the sufficient condition by lower bounding $\delta_{KL}^t(h_t, \hat{\alpha}_t)$,

$$\begin{aligned}\delta_{KL}^t(h_t, \hat{\alpha}_t) &= \mathbb{E}_P \left[\log \left[(1 - \hat{\alpha}_t) + \hat{\alpha}_t \frac{h_t}{q_{t-1}} \right] \right] \\ &\geq \mathbb{E}_P \left[(1 - \hat{\alpha}_t) \log 1 + \hat{\alpha}_t \log \frac{h_t}{q_{t-1}} \right] \quad (\text{Arithmetic Mean} \geq \text{Geometric Mean}) \\ &= \hat{\alpha}_t \mathbb{E}_P \left[\log \frac{h_t}{q_{t-1}} \right] \quad (\text{Linearity of expectation})\end{aligned}$$

If the lower bound is non-negative, then so is $\delta_{KL}^t(h_t, \hat{\alpha}_t)$. Hence,

$$\mathbb{E}_P \left[\log \frac{h_t}{q_{t-1}} \right] \geq 0$$

which is the stated sufficient condition.

For the necessary condition to hold, we know that

$$\begin{aligned}0 &\leq \delta_{KL}^t(h_t, \hat{\alpha}_t) \\ &= \mathbb{E}_P \left[\log \left[(1 - \hat{\alpha}_t) + \hat{\alpha}_t \frac{h_t}{q_{t-1}} \right] \right] \\ &\leq \log \mathbb{E}_P \left[(1 - \hat{\alpha}_t) + \hat{\alpha}_t \frac{h_t}{q_{t-1}} \right] \quad (\text{Jensen's inequality}) \\ &= \log \left[(1 - \hat{\alpha}_t) + \hat{\alpha}_t \mathbb{E}_P \left[\frac{h_t}{q_{t-1}} \right] \right] \quad (\text{Linearity of expectation})\end{aligned}$$

Taking exponential on both sides,

$$\begin{aligned}(1 - \hat{\alpha}_t) + \hat{\alpha}_t \mathbb{E}_P \left[\frac{h_t}{q_{t-1}} \right] &\geq 1 \\ \mathbb{E}_P \left[\frac{h_t}{q_{t-1}} \right] &\geq 1\end{aligned}$$

which is the stated necessary condition.

A.2. Theorem 2

Proof. We first derive the sufficient condition,

$$\begin{aligned}
 \delta_{KL}^t(h_t, \alpha_t) &= \int p \log q_t \, d\mathbf{x} - \int p \log q_{t-1} \, d\mathbf{x} \\
 &= \int p \log \frac{h_t^{\alpha_t} \cdot q_{t-1}}{Z_t} - \int p \log q_{t-1} \quad (\text{using Eq. (2)}) \\
 &= \alpha_t \cdot \mathbb{E}_P[\log h_t] - \log \mathbb{E}_{Q_{t-1}}[h_t^{\alpha_t}] \\
 &\geq \alpha_t \cdot \mathbb{E}_P[\log h_t] - \log \mathbb{E}_{Q_{t-1}}[h_t]^{\alpha_t} \quad (\text{Jensen's inequality}) \\
 &= \alpha_t \cdot [\mathbb{E}_P[\log h_t] - \log \mathbb{E}_{Q_{t-1}}[h_t]] \\
 &\geq 0 \quad (\text{by assumption}).
 \end{aligned} \tag{9}$$

Note that if $\alpha_t = 1$, the sufficient condition is also necessary. For the necessary condition,

$$\begin{aligned}
 0 \leq \delta_{KL}^t(h_t, \alpha_t) &= \alpha_t \cdot \mathbb{E}_P[\log h_t] - \log \mathbb{E}_{Q_{t-1}}[h_t^{\alpha_t}] \\
 &\leq \alpha_t \cdot \mathbb{E}_P[\log h_t] - \mathbb{E}_{Q_{t-1}}[\log h_t^{\alpha_t}] \quad (\text{Jensen's inequality}) \\
 &= \alpha_t \cdot [\mathbb{E}_P[\log h_t] - \mathbb{E}_{Q_{t-1}}[\log h_t]] \quad (\text{Linearity of expectation}) \\
 &\leq \mathbb{E}_P[\log h_t] - \mathbb{E}_{Q_{t-1}}[\log h_t] \quad (\text{since } \alpha_t > 0).
 \end{aligned}$$

□

A.3. Theorem 1

Proof. By assumption, we can optimize Eq. (3) to get,

$$h_t \propto \left(\frac{p}{q_{t-1}} \right)^{\beta_t}.$$

Substituting for h_t in the multiplicative boosting formulation in Eq. (2),

$$\begin{aligned}
 q_t &\propto \frac{q_{t-1} \cdot h_t}{Z_{q_t}} \\
 &\propto q_{t-1} \cdot \left(\frac{p}{q_{t-1}} \right)^{\beta_t} \\
 &= \frac{p^{\beta_t} \cdot q_{t-1}^{1-\beta_t}}{Z_{q_t}}
 \end{aligned}$$

where the partition function $Z_{q_t} = \int p^{\beta_t} \cdot q_{t-1}^{1-\beta_t}$.

In order to prove the inequality, we first obtain a lower bound on the log-partition function, Z_{q_t} . For any given point, we have,

$$p^{\beta_t} \cdot q_{t-1}^{1-\beta_t} \leq \beta_t p + (1 - \beta_t) q_{t-1}$$

using the Arithmetic Mean \geq Geometric Mean inequality given that densities are non-negative. Integrating over all points in the domain,

$$\log Z_q \leq \log [\beta Z_p + (1 - \beta) Z_{q_{t-1}}] = 0 \tag{10}$$

where we have used the fact that p and q_{t-1} are normalized densities.

Now, consider the following quantity 1,

$$\begin{aligned}
 D_{KL}(P||Q_t) &= \mathbb{E}_P \left[\log \frac{p}{q_t} \right] \\
 &= \mathbb{E}_P \left[\log \frac{p}{\frac{p^{\beta_t} \cdot q_{t-1}^{1-\beta_t}}{Z_{q_t}}} \right] \\
 &= (1 - \beta_t) \mathbb{E}_P \left[\log \frac{p}{q_{t-1}} \right] + \log Z_{q_t} \\
 &\leq (1 - \beta_t) \mathbb{E}_P \left[\log \frac{p}{q_{t-1}} \right] \\
 &\leq \mathbb{E}_P \left[\log \frac{p}{q_{t-1}} \right] \quad (\text{since } \beta_t \geq 0) \\
 &= D_{KL}(P||Q_{t-1})
 \end{aligned}$$

where the second last inequality is due to Eq. (10). This finishes the proof. \square

A.4. Proposition 2

Proof. By the f -optimality assumption,

$$r_t = f' \left(\frac{p}{q_{t-1}} \right).$$

Hence, $h_t = \frac{p}{q_{t-1}}$. From Eq. (2),

$$q_t = q_{t-1} \cdot h_t^{\alpha_t} = p$$

finishing the proof. \square

A.5. Corollary 1

Proof. Let u_t denote the joint distribution over (\mathbf{x}, y) at round t . We will prove a slightly more general result where we have m positive training examples sampled from p and the k negative training examples sampled from q_{t-1} . Hence,

$$p = u(\mathbf{x}|y = +1) \quad u(y = +1) = \frac{m}{m+k} \quad (11)$$

$$q_{t-1} = u(\mathbf{x}|y = -1) \quad u(y = -1) = \frac{k}{m+k}. \quad (12)$$

The Bayes optimal density c_t can be expressed as,

$$\begin{aligned}
 c_t &= u(y = +1 | \mathbf{x}) \\
 &= u(\mathbf{x} | y = +1)u(y = +1)/u(\mathbf{x}).
 \end{aligned} \quad (13)$$

Similarly,

$$1 - c_t = u(\mathbf{x} | y = -1)u(y = -1)/u(\mathbf{x}). \quad (14)$$

From Eqs. (11- 14), we have,

$$h_t = \gamma \cdot \frac{c_t}{1 - c_t} = \frac{p}{q_{t-1}}.$$

where $\gamma = \frac{k}{m}$. Finally from Eq. (2),

$$q_t = q_{t-1} \cdot h_t^{\alpha_t} = p$$

finishing the proof. \square

A.6. Additional results

In this section, we present theoretical results regarding the proposed boosting framework.

First, we derive the necessary and sufficient conditions that guarantee reduction in KL-divergence for multiplicative boosting updates as specified by Eq. (8).

Corollary 2. [to Theorem 2] Let $\delta_{KL}^t(f_t) = D_{KL}(P||Q_{t-1}) - D_{KL}(P||Q_t)$ denote the reduction in KL-divergence at the t^{th} round of multiplicative boosting in normalizing flow models for an invertible transformation f_t . Then, $\delta_{KL}^t(f_t) \geq 0$ iff $\mathbb{E}_p \left[\log q_{t-1}(f_t(\mathbf{x})) \left| \det \frac{\partial f_t(\mathbf{x})}{\partial \mathbf{x}} \right| \right] \geq \mathbb{E}_p \left[\log q_t(g_t(\mathbf{x})) \left| \det \frac{\partial g_t(\mathbf{x})}{\partial \mathbf{x}} \right| \right]$ where $g_t = f_t^{-1}$.

Proof. We first note a change-of-variables equation analogous to Eq. (8), but in the reverse direction.

$$q_{t-1}(\mathbf{x}) = q_t(g_t(\mathbf{x})) \left| \det \frac{\partial g_t(\mathbf{x})}{\partial \mathbf{x}} \right| \quad (15)$$

where $\frac{\partial g_t(\mathbf{x})}{\partial \mathbf{x}}$ denotes the Jacobian of g_t at \mathbf{x} . The proof for both the forward and backward implications follow from the equivalence below,

$$\begin{aligned} \delta_{KL}^t(f_t) &= \mathbb{E}_p [\log q_t(\mathbf{x})] - \mathbb{E}_p [\log q_{t-1}(\mathbf{x})] \\ &= \mathbb{E}_p \left[\log \left(q_{t-1}(f_t(\mathbf{x})) \left| \det \frac{\partial f_t(\mathbf{x})}{\partial \mathbf{x}} \right| \right) \right] - \mathbb{E}_p [\log q_{t-1}(\mathbf{x})] \\ &= \mathbb{E}_p \left[\log \left(q_{t-1}(f_t(\mathbf{x})) \left| \det \frac{\partial f_t(\mathbf{x})}{\partial \mathbf{x}} \right| \right) \right] - \mathbb{E}_p \left[\log \left(q_t(g_t(\mathbf{x})) \left| \det \frac{\partial g_t(\mathbf{x})}{\partial \mathbf{x}} \right| \right) \right]. \end{aligned}$$

□

The second result in this section derives the optimal model weight, α_t for an adversarial Bayes optimal classifier. Define an adversarial Bayes optimal classifier c'_t as one that assigns the density $c'_t = 1 - c_t$ where c_t is the Bayes optimal classifier.

Corollary 3. [to Corollary 1] For an adversarial Bayes optimal classifier c'_t , δ_{KL}^t attains a maxima of zero when $\alpha_t = 0$.

Proof. For an adversarial Bayes optimal classifier,

$$c'_t = u(\mathbf{x} | y = -1)u(y = -1)/u(\mathbf{x}) \quad (16)$$

$$1 - c'_t = u(\mathbf{x} | y = +1)u(y = +1)/u(\mathbf{x}). \quad (17)$$

From Eqs. (11,12, 16,17),

$$h_t = \gamma \cdot \frac{c'_t}{1 - c'_t} = \frac{q_{t-1}}{p}.$$

Substituting the above intermediate model in Eq. (9),

$$\begin{aligned} \delta_{KL}^t(h_t, \alpha_t) &= \alpha_t \cdot \mathbb{E}_P \left[\log \frac{q_{t-1}}{p} \right] - \log \mathbb{E}_{Q_{t-1}} \left[\frac{q_{t-1}}{p} \right]^{\alpha_t} \\ &\leq \alpha_t \cdot \mathbb{E}_P \left[\log \frac{q_{t-1}}{p} \right] - \mathbb{E}_{Q_{t-1}} \left[\alpha_t \cdot \log \frac{q_{t-1}}{p} \right] \quad (\text{Jensen's inequality}) \\ &= \alpha_t \cdot \left[\mathbb{E}_P \left[\log \frac{q_{t-1}}{p} \right] - \mathbb{E}_{Q_{t-1}} \left[\log \frac{q_{t-1}}{p} \right] \right] \quad (\text{Linearity of expectation}) \\ &= -\alpha_t [D_{KL}(P || Q_{t-1}) + D_{KL}(Q_{t-1} || P)] \\ &\leq 0 \quad (D_{KL} \text{ is non-negative}). \end{aligned}$$

By inspection, the equality holds when $\alpha_t = 0$ finishing the proof.

□

B. Implementation details and additional experimental results

B.1. Mixture of Gaussians

Table 3: Heuristic strategies for setting intermediate model weights.

Weighting strategy	Per-round weights, α_t
Unity	1
Uniform	$1/(T + 1)$
Decay	$1/2^t$

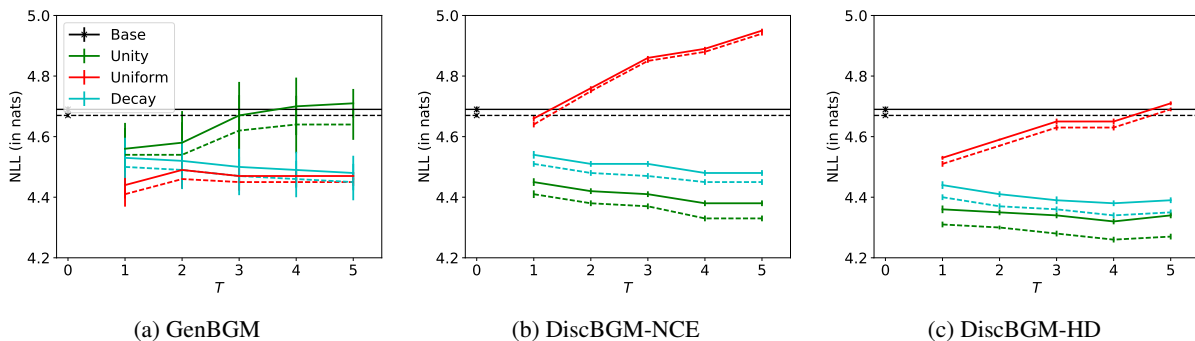


Figure 5: Train (dashed curves) and test (bold curves) NLL (in nats) of various heuristic strategies on density estimation for a mixture of Gaussians. T is the number of rounds of boosting. The base model is shown as a black cross at $T = 0$.

Our algorithmic framework requires as hyperparameters the number of rounds of boosting (*i.e.*, T), as well as a set of weights associated with each model in the ensemble (*i.e.*, α 's). For any practical setting, these hyperparameters are specific to the dataset and task under consideration and should be set based on cross-validation. Here, we propose some heuristic weighting strategies for assigning the model weights in Table 3 and study their effect on density estimation for the aforementioned setting of fitting a boosted generative model to a mixture of Gaussians. The partition function is estimated using importance sampling with a Gaussian fitted to the training data as the proposal.

In Figure 5, we observe that the performance of the algorithms can however be sensitive to the weighting strategies. In particular, DiscBGM can produce worse estimates as T increases for the “uniform” (red) strategy. GenBGM also slightly degrades in performance with increasing T for the “unity” (green) strategy. For all other strategies, the algorithms are fairly robust to the number of rounds of boosting. Notably, the “decay” (cyan) strategy achieves stable performance for both the algorithms. Intuitively, this heuristic follows the rationale of reducing the step size in gradient based stochastic optimization algorithms, and we expect this strategy to work better even in other settings. However, this strategy could potentially result in slower convergence as opposed to the unity strategy.

B.2. MNIST

VAE architecture and learning procedure details. Only the last layer in every VAE is stochastic, rest are deterministic. The inference network specifying the posterior contains the same architecture for the hidden layer as the generative network. The prior over the latent variables is standard Gaussian, the hidden layer activations are ReLU, and learning is done using Adam (Kingma & Ba, 2015) with a learning rate of 10^{-3} and mini-batches of size 100. The convergence curves for the baseline VAEs and the smaller VAE used in the boosted generative model is shown in Figure 6.

CNN architecture and learning procedure details. The CNN contains two convolutional layers and a single full connected layer with 1024 units. Convolution layers have kernel size 5×5 , and 32 and 64 output channels, respectively. We apply ReLUs and 2×2 max pooling after each convolution. The net is randomly initialized prior to training, and learning is done for 2 epochs using Adam (Kingma & Ba, 2015) with a learning rate of 10^{-3} and mini-batches of size 100.

Sampling procedure for BGM sequences. Samples from the BGM sequences are drawn from a Markov chain run using

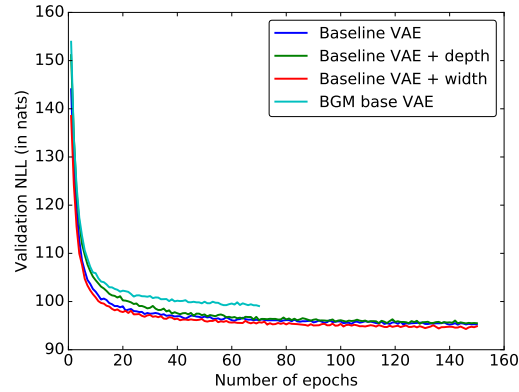


Figure 6: The baselines are run up till convergence as measured by validation NLL. The training of the VAEs used in the BGMs is terminated early to ensure fairness in total computation time.

the Metropolis-Hastings algorithm with a discrete, uniformly random proposal and the BGM distribution as the stationary distribution for the chain. Every sample in Figure 4 (d-f) is drawn from an independent Markov chain with a burn-in period of 100,000 samples and a different start seed state.

B.3. CIFAR-10

Model architecture. The overall architecture contains three scales. Every scale is composed of three different kinds of layers: coupling layers (C), squeezing layers (S), and factoring layers (F). Please refer to (Dinh et al., 2016) for a complete descriptions of these layers. In increasing order starting from the scale containing the layer expressing the prior density, the scale compositions are as follows,

1. Scale 1: F-C-C-C-C
2. Scale 2: F-C-C-C-S-C-C-C
3. Scale 3: F-C-C-S-C-C

Our implementation details follow (Dinh et al., 2016) closely. The batch size was set to 12 and the images was scaled to lie between $[-1, 1]$ before feeding them into the network. Following common practice when modeling densities, we add real-valued noise to the pixel values to dequantize the data (Uria et al., 2013). We use the ADAM optimizer (Kingma & Ba, 2015) with a learning rate of 0.001. For the boosted models, the learning rate was reset for training the additional scale (Scale 3). We used weight normalization (Salimans & Kingma, 2016) and batch normalization (Ioffe & Szegedy, 2015) to improve the training signal across layers.

Samples. Samples generated from the best performing boosted model and the jointly trained baseline model are shown in Figure 7 and Figure 8 respectively.

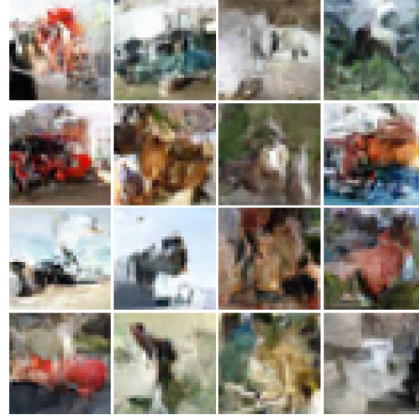


Figure 7: Samples generated by the boosted normalizing flow model ($\beta_2 = 0.5$) for CIFAR-10.

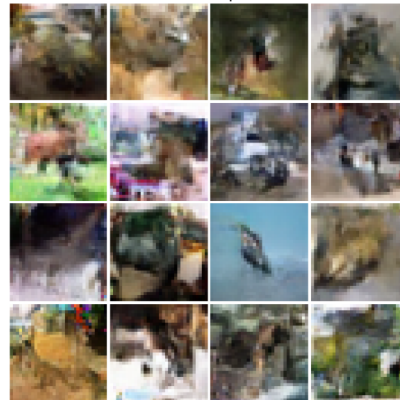


Figure 8: Samples generated by the jointly trained baseline model for the CIFAR-10 dataset.