
Deep Discriminative to Kernel Generative Networks for Calibrated Inference

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The fight between discriminative versus generative goes deep, in both the study of
2 artificial and natural intelligence. In our view, both camps have complementary
3 values. So, we sought to synergistically combine them. Here, we propose a
4 methodology to convert deep discriminative networks to kernel generative networks.
5 We leveraged the fact that deep models, including both random forests and deep
6 networks, learn internal representations which are unions of polytopes with affine
7 activation functions to conceptualize them both as generalized partitioning rules.
8 We replace the affine function in each polytope populated by the training data with
9 Gaussian kernel that results in a generative model. Theoretically, we derive the
10 conditions under which our generative models are a consistent estimator of the
11 corresponding class conditional density. Moreover, our proposed models obtain
12 well calibrated posteriors for in-distribution, and extrapolate beyond the training
13 data to handle out-of-distribution inputs reasonably. We believe this approach
14 may be an important step in unifying the thinking and the approaches across the
15 discriminative and the generative divide.

16 1 Introduction

17 Machine learning methods, specially deep neural networks and random forests have shown excellent
18 performance in many real-world tasks, including drug discovery, autonomous driving and clinical
19 surgery. However, calibrating confidence over the whole feature space for these models remains a key
20 challenge in the field. Although these learning algorithms can achieve near optimal performance at
21 inferring on the samples lying in the high density regions of the training data [1-3], they yield highly
22 confident predictions for the samples lying far away from the training data [4]. Calibrated confidence
23 within the training or in-distribution (ID) region as well as in the out-of-distribution (OOD) region is
24 crucial for safety critical applications like autonomous driving and computer-assisted surgery, where
25 any aberrant reading should be detected and taken care of immediately [4, 5]. A well-calibrated model
26 capable of quantifying the uncertainty associated with inference for any points from the training
27 distribution as well as detecting OOD data can be a life-saver in these cases.

28 The approaches to calibrate OOD confidence for learning algorithms described in the literature can
29 be roughly divided into two groups: discriminative and generative. Discriminative approaches try to
30 scale the posteriors based on OOD detection or modify the learning loss function. Intuitively, the
31 easiest solution for OOD confidence calibration is to learn a function that gives higher scores for
32 in-distribution samples and lower scores for OOD samples, and thereby re-scale the posterior or
33 confidence score from the original model accordingly [6]. There are a number of approaches in the
34 literature which try to either modify the loss function [7-9] or adversarially train the network to be
35 less confident on OOD samples [10, 4]. However, one can adversarially manipulate an OOD sample
36 where the model is less confident to find another OOD sample where the model is overconfident
37 [11, 4, 12]. Recently, as shown by Hein et al. [4], the RELU networks produce arbitrarily high

38 confidence as the inference point moves far away from the training data. Therefore, calibrating
39 `ReLU` networks for the whole OOD region is not possible without fundamentally changing the
40 network architecture. As a result, all of the aforementioned algorithms are unable to provide any
41 guarantee about the performance of the network through out the whole feature space. On the other end
42 of the spectrum, the generative group tries to learn generative models for both the in-distribution as
43 well as the out-of-distribution samples. The general idea for the generative group is to get likelihoods
44 for a particular sample out of the generative models for both ID and OOD to do likelihood ratio test
45 [13] or control the likelihood for training distribution far away from the training data to detect OOD
46 samples by thresholding. However, it is not obvious how to control likelihoods far away from the
47 training data for powerful generative models like variational autoencoders (VAEs) [14] and generative
48 adversarial networks (GAN) [15]. Moreover, Nalisnick et al. [16] and Hendrycks et al. [10] showed
49 VAEs and GANs can also yield overconfident likelihoods far away from the training data.

50 The algorithms described so far are concerned with OOD confidence calibration for deep-nets only.
51 However, in this paper, we show that other approaches which partition the feature space, for example
52 random forest, can also suffer from poor confidence calibration both in the ID and the OOD regions.
53 Moreover, the algorithms described above are concerned about the confidence of the algorithms in the
54 OOD region only and they do not address the confidence calibration within the training distribution
55 at all. This issue is addressed separately in a different group of literature [17-19]. In this paper, we
56 consider both calibration problem jointly and propose an approach that achieves good calibration
57 throughout the whole feature space.

58 In this paper, we conceptualize both random forest and `ReLU` networks as generalized partitioning
59 rules with an affine activation over each polytope. We consider replacing the affine functions learned
60 over the polytopes with Gaussian kernels. We propose two novel kernel density estimation techniques
61 named *Kernel Generative Forest* (KGF) and *Kernel Generative Network* (KGN). We theoretically show
62 that they asymptotically converge to the true training distribution under certain conditions. At the
63 same time, the estimated likelihood from the kernel generative models decreases for samples far away
64 from the training samples. By adding a suitable bias to the kernel density estimate, we can achieve
65 calibrated posterior over the classes in the OOD region. It completely excludes the need for providing
66 OOD training examples to the model. We conduct several simulation and real data studies that show
67 both KGF and KGN are robust against OOD samples while they maintain good performance in the
68 in-distribution region.

69 2 Related Works and Our Contributions

70 There are a number of approaches in the literature which attempt to learn a generative model and
71 control the likelihoods far away from the training data. For example, Ren et al. [13] employed
72 likelihood ratio test for detecting OOD samples. Wan et al. [8] modify the training loss so that the
73 downstream projected features follow a Gaussian distribution. However, there is no guarantee of
74 performance for OOD detection for the above methods. To the best of our knowledge, only Meinke
75 et al. [5] has proposed an approach to guarantee asymptotic performance for OOD detection. They
76 model the training and the OOD distribution using Gaussian mixture models which enable them
77 to control the class conditional posteriors far away. Compared to the aforementioned methods, our
78 approach differs in several ways:

- 79 • We address the confidence calibration problems for both `ReLU`-nets and random forests from
80 a common ground.
- 81 • We address in-distribution (ID) and out-of-distribution (OOD) calibration problem as a
82 continuum rather than two separate problems.
- 83 • We provide guarantees for asymptotic convergence of our proposed approach under certain
84 conditions for both ID and OOD regions.
- 85 • We propose an unsupervised OOD calibration approach, i.e., we do not need to train
86 exhaustively on different OOD samples.

87 **3 Methods**

88 **3.1 Setting**

89 Consider a supervised learning problem with independent and identically distributed training samples
 90 $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ such that $(X, Y) \sim P_{X,Y}$, where $X \sim P_X$ is a $\mathcal{X} \subseteq \mathbb{R}^d$ valued input and $Y \sim P_Y$ is
 91 a $\mathcal{Y} = \{1, \dots, K\}$ valued class label. We define in-distribution region as the high density region
 92 of $P_{X,Y}$ and denote it by \mathcal{S} . Here the goal is to learn a confidence score, $\mathbf{g} : \mathbb{R}^d \rightarrow [0, 1]^K$,
 93 $\mathbf{g}(\mathbf{x}) = [g_1(\mathbf{x}), g_2(\mathbf{x}), \dots, g_K(\mathbf{x})]$ such that,

$$g_y(\mathbf{x}) = \begin{cases} P_{Y|X}(y|\mathbf{x}), & \text{if } \mathbf{x} \in \mathcal{S} \\ P_Y(y), & \text{if } \mathbf{x} \notin \mathcal{S} \end{cases}, \quad \forall y \in \mathcal{Y} \quad (1)$$

94 where $P_{Y|X}(y|\mathbf{x})$ is the posterior probability for class y given by the Bayes formula:

$$P_{Y|X}(y|\mathbf{x}) = \frac{P_{X|Y}(\mathbf{x}|y)P_Y(y)}{\sum_{k=1}^K P_{X|Y}(\mathbf{x}|k)P_Y(k)}, \quad \forall y \in \mathcal{Y}. \quad (2)$$

95 Here $P_{X|Y}(\mathbf{x}|y)$ is the class conditional density for the training data which we will refer as $f_y(\mathbf{x})$
 96 hereafter for brevity.

97 **3.2 Background and Main Idea**

98 Deep discriminative networks partition the feature space \mathbb{R}^d into a union of p affine polytopes Q_r
 99 such that $\bigcup_{r=1}^p Q_r = \mathbb{R}^d$, and learn an affine function over each polytope [4, 20]. Mathematically, the
 100 class-conditional density for the label y estimated by these deep discriminative models at a particular
 101 point \mathbf{x} can be expressed as:

$$\hat{f}_y(\mathbf{x}) = \sum_{r=1}^p (\mathbf{a}_r^\top \mathbf{x} + b_r) \mathbb{1}(\mathbf{x} \in Q_r). \quad (3)$$

102 For example, in the case of a decision tree, $\mathbf{a}_r = \mathbf{0}$, i.e., decision tree assumes uniform distribution
 103 for the class-conditional densities over the leaf nodes. Among these polytopes, the ones that lie on
 104 the boundary of the training data extend to the whole feature space and hence encompass all the OOD
 105 samples. Since the posterior probability for a class is determined by the affine activation over each of
 106 these polytopes, the algorithms tend to be overconfident when making predictions on the OOD inputs.
 107 Moreover, there exist some polytopes that are not populated with training data. These unpopulated
 108 polytopes serve to interpolate between the training sample points. If we replace the affine activation
 109 function of the populated polytopes with Gaussian kernel \mathcal{G} learned using maximum likelihood
 110 approach on the training points within the corresponding polytope and prune the unpopulated ones,
 111 the tail of the kernel will help interpolate between the training sample points while assigning lower
 112 likelihood to the low density or unpopulated polytope regions of the feature space. This may result in
 113 better confidence calibration for the proposed modified approach.

114 **3.3 Proposed Model**

115 Consider the collection of polytope indices \mathcal{P} which contains the indices of total \tilde{p} polytopes populated
 116 by the training data. We consider replacing the affine function over the populated polytopes with a
 117 Gaussian kernel $\mathcal{G}(\cdot; \hat{\mu}_r, \hat{\Sigma}_r)$. For a particular inference point \mathbf{x} , we consider the Gaussian kernel
 118 with the minimum distance from the center of the kernel to the corresponding point:

$$r_{\mathbf{x}}^* = \operatorname{argmin}_r \|\mu_r - \mathbf{x}\|, \quad (4)$$

119 where $\|\cdot\|$ denotes a suitable distance measure. We use Euclidean distance metric while conducting
 120 the simulation and the benchmark datasets experiments in this paper for simplicity. In short, we
 121 modify Equation 3 from the parent ReLU-net or random forest to estimate the class-conditional density
 122 as:

$$\tilde{f}_y(\mathbf{x}) = \frac{1}{n_y} \sum_{r \in \mathcal{P}} n_{ry} \mathcal{G}(\mathbf{x}; \mu_r, \Sigma_r) \mathbb{1}(r = r_{\mathbf{x}}^*), \quad (5)$$

123 where n_y is the total number of samples with label y and n_{r_y} is the number of samples from class y
 124 that end up in polytope Q_r . We add a bias to the class conditional density \tilde{f}_y :

$$\hat{f}_y(\mathbf{x}) = \tilde{f}_y(\mathbf{x}) + \frac{b}{\log(n)}. \quad (6)$$

125 Note that in Equation 6, $\frac{b}{\log(n)} \rightarrow 0$ as the total training points, $n \rightarrow \infty$. The class posterior
 126 probability (confidence) $\hat{g}_y(\mathbf{x})$ of class y for a test point \mathbf{x} is estimated using the Bayes rule as
 127 follows:

$$\hat{g}_y(\mathbf{x}) = \frac{\hat{f}_y(\mathbf{x})\hat{P}_Y(y)}{\sum_{k=1}^K \hat{f}_k(\mathbf{x})\hat{P}_Y(k)}, \quad (7)$$

128 where $\hat{P}_Y(y)$ is the empirical prior probability of class y estimated from the training data. We
 129 estimate the class for a particular inference point \mathbf{x} as:

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} \hat{g}_y(\mathbf{x}). \quad (8)$$

130 3.4 Desiderata

131 We desire our proposed model to estimate confidence score \hat{g}_y to satisfy the following two desiderata:

1. **Asymptotic Performance:** We want point-wise convergence for our estimated confidence as $n \rightarrow \infty$, i.e.,

$$\max_{y \in \mathcal{Y}} \sup_{\mathbf{x} \in \mathbb{R}^d} |g_y(\mathbf{x}) - \hat{g}_y(\mathbf{x})| \rightarrow 0.$$

- 132 2. **Finite Sample Performance:** We want better posterior calibration for $\hat{g}_y(\mathbf{x})$ both in ID and
 133 OOD region compared to that of its parent model.

134 We theoretically derive the conditions under which we achieve Desiderata 1 in Section 4. However, we
 135 run extensive experiments on various simulation and benchmark datasets in Section 6 to empirically
 136 verify that our proposed approach achieves Desiderata 2.

137 4 Theoretical Results

138 **Theorem 1** (Asymptotic Convergence to the True Distribution). *Consider a partition rule that*
 139 *partitions \mathbb{R}^d into hypercubes of the same size $h_n > 0$. Formally, let $\mathcal{P}_n = \{Q_1, Q_2, \dots\}$ be a*
 140 *partition of \mathbb{R}^d , that is, it partitions \mathbb{R}^d into sets of the type $\prod_{i=1}^d [\psi_i h_n, (\psi_i + 1)h_n)$, where ψ_i 's are*
 141 *integers. Let n be the total number of samples and n_r be the number of data points within polytope*
 142 *Q_r . Consider the probability density f estimated for the samples populating the polytopes using*
 143 *Equation 5 denoted as \hat{f} . The conditions for choosing the Gaussian kernel parameters are:*

- 144 1. *The center of the kernel can be any point z_r within the polytope Q_r as $n \rightarrow \infty$,*
- 145 2. *The kernel bandwidth along any dimension σ_r is any positive number always bounded by*
 146 *the polytope bandwidth h_n as $n \rightarrow \infty$, i.e., $\sigma_r = C_r h_n$, where $0 < C_r \leq 1$.*

147 *Consider the following assumptions as well:*

- 148 1. *The polytope bandwidth $h_n \rightarrow 0$ as $n \rightarrow \infty$.*
- 149 2. *n grows faster than the shrinkage of h_n , i.e., $nh_n \rightarrow \infty$ as $h_n \rightarrow 0$ in probability.*

150 *Given these assumptions, we have that as $n \rightarrow \infty$:*

$$\sup_{\mathbf{x} \in \mathbb{R}^d} |f(\mathbf{x}) - \hat{f}(\mathbf{x})| \rightarrow 0,$$

151 *where $|\cdot|$ denotes absolute value of the scalar it operates on.*

152 *Proof.* Please see Appendix A for the proof. □

153 **Theorem 2** (Asymptotic OOD Convergence). *Given n independent and identically distributed*
 154 *training samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, we define the distance of an inference point \mathbf{x} from the training points*
 155 *as:*

$$d_{\mathbf{x}} = \min_{i=1, \dots, n} \|\mathbf{x} - \mathbf{x}_i\|. \quad (9)$$

156 *Here $\|\cdot\|$ denotes a suitable distance measure as mentioned in Equation 4. Given non-zero and*
 157 *bounded bandwidth of the Gaussians, then we have almost sure convergence for \hat{g}_y as: $\hat{g}_y(\mathbf{x}) \xrightarrow{a.s.}$*
 158 *$\hat{P}_Y(y)$ as $d_{\mathbf{x}} \rightarrow \infty$.*

159 *Proof.* Please see Appendix A for the proof. □

160 **Corollary 2.1.** *Given the conditions in Theorem 1 and 2 we have:*

$$\max_{y \in \mathcal{Y}} \sup_{\mathbf{x} \in \mathbb{R}^d} |g_y(\mathbf{x}) - \hat{g}_y(\mathbf{x})| \rightarrow 0.$$

161 *Proof.* Using the law of large numbers, we have $\hat{P}_Y(y) = \frac{n_y}{n} \xrightarrow{a.s.} P_Y(y)$ as $n_y \rightarrow \infty$. The rest of the
 162 proof follows from Theorem 1 and 2. □

163 5 Model Parameter Estimation

164 5.1 Gaussian Kernel Parameter Estimation

165 Theorem 1 implies that the Gaussian kernel parameters need to maintain two key properties. We
 166 use the training data within the polytopes to estimate the Gaussian parameters in a way that we
 167 asymptotically satisfy the above two conditions for consistency. To satisfy the first condition, we set
 168 the kernel center as:

$$\hat{\mu}_r = \frac{1}{n_r} \sum_{i=1}^n \mathbf{x}_i \mathbb{1}(\mathbf{x}_i \in Q_r). \quad (10)$$

169 Note that $\hat{\mu}_r$ in Equation 10 resides always within the corresponding polytope Q_r . For improving the
 170 estimates for the kernel bandwidth, we incorporate the samples from other polytopes Q_s based on the
 171 similarity w_{rs} between Q_r and Q_s . Moreover, We constrain our estimated Gaussian kernels to have
 172 diagonal covariance matrix. We use weighted likelihood estimation to estimate the variance Σ_r for a
 173 particular polytope Q_r . For simplicity, we will describe the estimation procedure for w_{rs} later. The
 174 weighted likelihood estimation for Σ_r can be written as:

$$\hat{\Sigma}_r = \operatorname{argmin}_{\Sigma} - \sum_{i=1}^n \sum_{s \in \mathcal{P}} w_{rs} \mathbb{1}(\mathbf{x}_i \in Q_s) \log \mathcal{G}(\mathbf{x}_i; \hat{\mu}_r, \Sigma) + \lambda \|\Sigma^{-1}\|_F^2, \quad (11)$$

175 where we regularize the Frobenius norm of precision matrix Σ^{-1} so that Σ does not become singular
 176 and λ is the regularization parameter. By solving Equation 11, we find:

$$\hat{\Sigma}_r = \frac{\sum_{s \in \mathcal{P}} \sum_{i=1}^n w_{rs} \mathbb{1}(\mathbf{x}_i \in Q_s) (\mathbf{x}_i - \hat{\mu}_r)(\mathbf{x}_i - \hat{\mu}_r)^\top + \lambda I_d}{\sum_{s \in \mathcal{P}} \sum_{i=1}^n w_{rs} \mathbb{1}(\mathbf{x}_i \in Q_s)}, \quad (12)$$

177 where, I_d is a d dimensional identity matrix. However, we want Σ_r to be estimated based on the
 178 samples within Q_r so that the second condition for the Gaussian parameters is satisfied. Therefore,
 179 as $n \rightarrow \infty$ and $h_n \rightarrow 0$, the estimated weights w_{rs} should satisfy the condition:

$$w_{rs} \rightarrow \begin{cases} 0, & \text{if } Q_r \neq Q_s \\ 1, & \text{if } Q_r = Q_s. \end{cases} \quad (13)$$

180 We need Condition 13 as we will be only using the data within the polytope Q_r as $n \rightarrow \infty$ to estimate
 181 the Gaussian bandwidth and the estimated Gaussian bandwidth will be bounded by the polytope
 182 bandwidth. Additionally, we use weighted samples to replace the ratio $\frac{n_{ry}}{n_y}$ in Equation 5 as:

$$\frac{\tilde{w}_{ry}}{\tilde{w}_y} = \frac{\tilde{w}_{ry}}{\sum_{r \in \mathcal{P}} \tilde{w}_{ry}} = \frac{\sum_{s \in \mathcal{P}} \sum_{i=1}^n w_{rs} \mathbb{1}(\mathbf{x}_i \in Q_s) \mathbb{1}(y_i = y)}{\sum_{r \in \mathcal{P}} \sum_{s \in \mathcal{P}} \sum_{i=1}^n w_{rs} \mathbb{1}(\mathbf{x}_i \in Q_s) \mathbb{1}(y_i = y)}. \quad (14)$$

183 Note that if we satisfy Condition [13](#), then we have $\frac{\tilde{w}_{ry}}{\tilde{w}_y} \rightarrow \frac{n_{ry}}{n_y}$ as $n \rightarrow \infty$. Therefore, we modify
 184 Equation [5](#) as:

$$\tilde{f}_y(\mathbf{x}) = \frac{1}{\tilde{w}_y} \sum_{r \in \mathcal{P}} \tilde{w}_{ry} \mathcal{G}(\mathbf{x}; \hat{\mu}_r, \hat{\Sigma}_r) \mathbb{1}(r = \hat{r}_\mathbf{x}^*), \quad (15)$$

185 where $\hat{r}_\mathbf{x}^* = \operatorname{argmin}_r \|\hat{\mu}_r - \mathbf{x}\|$. Below, we describe how we estimate w_{rs} for KGF and KGN.

186 5.2 Kernel Generative Forest

187 Consider T number of decision trees in a random forest trained on n i.i.d training samples
 188 $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$. Each tree t partitions the feature space into p_t polytopes resulting in a set of polytopes:
 189 $\{\{Q_{t,r}\}_{r=1}^{p_t}\}_{t=1}^T$. The intersection of these polytopes gives a new set of polytopes $\{Q_r\}_{r=1}^p$ for the
 190 forest. For any point $\mathbf{x}_r \in Q_r$, we push every other sample point $\mathbf{x}_s \in Q_s$ down the trees. Now, we
 191 define the weight w'_{rs} as:

$$w'_{rs} = \frac{t_{rs}}{T}, \quad (16)$$

192 where t_{rs} is the total number trees \mathbf{x}_r and \mathbf{x}_s end up in the same leaf node. Note that $0 \leq w'_{rs} \leq 1$.
 193 *If the two samples end up in the same leaves in all the trees, they belong to the same polytope, i.e.*
 194 $Q_r = Q_s$.

195 In short, w'_{rs} is the fraction of total trees where the two samples follow the same path from the root to
 196 a leaf node. We exponentiate w'_{rs} with a suitable function of n which grows with n so that Condition
 197 [13](#) is satisfied:

$$w_{rs} = (w'_{rs})^{O(n)}. \quad (17)$$

198 5.3 Kernel Generative Network

199 Consider a fully connected `ReLU`-net trained on n i.i.d training samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$. We have the
 200 set of all nodes denoted by \mathcal{N}_l at a particular layer l . We can randomly pick a node $n_l \in \mathcal{N}_l$ from \mathcal{A}_l
 201 at each layer l , and construct a sequence of nodes starting at the input layer and ending at the output
 202 layer which we call an **activation path**: $m = \{n_l \in \mathcal{N}_l\}_{l=1}^L$. Note that there are $N = \prod_{l=1}^L |\mathcal{N}_l|$
 203 possible activation paths for a sample in the `ReLU`-net, where $|\cdot|$ denotes the cardinality or the number
 204 of elements in the set. We index each path by a unique identifier number $z \in \mathbb{N}$ and construct a
 205 sequence of activation paths as: $\mathcal{M} = \{m_z\}_{z=1, \dots, N}$. Therefore, \mathcal{M} contains all possible activation
 206 pathways from the input to the output of the network.

207 While pushing a training sample \mathbf{x}_i through the network, we define the activation from a `ReLU` unit at
 208 any node as ‘1’ when it has non-negative input and ‘0’ otherwise. Therefore, the activation indicates
 209 on which side of the affine function at each node the sample falls. The activation for all nodes in an
 210 activation path m_z for a particular sample creates an **activation mode** $a_z \in \{0, 1\}^L$. If we evaluate
 211 the activation mode for all activation paths in \mathcal{M} while pushing a sample through the network, we
 212 get a sequence of activation modes: $\mathcal{A}_r = \{a_z^r\}_{z=1}^N$. Here r is the index of the polytope where the
 213 sample falls in.

214 *If the two sequences of activation modes for two different training samples are identical, they belong*
 215 *to the same polytope.* In other words, if $\mathcal{A}_r = \mathcal{A}_s$, then $Q_r = Q_s$. This statement holds because the
 216 above samples will lie on the same side of the affine function at each node in different layers of the
 217 network. Now, we define the weight w'_{rs} as:

$$w'_{rs} = \frac{\sum_{z=1}^N \mathbb{1}(a_z^r = a_z^s)}{N}. \quad (18)$$

218 Note that $0 \leq w'_{rs} \leq 1$. In short, w'_{rs} is the fraction of total activation paths which are identically
 219 activated for two samples in two different polytopes r and s . We exponentiate the weights using
 220 Equation [17](#)

221 Pseudocodes outlining the two algorithms are provided in Appendix [C](#).

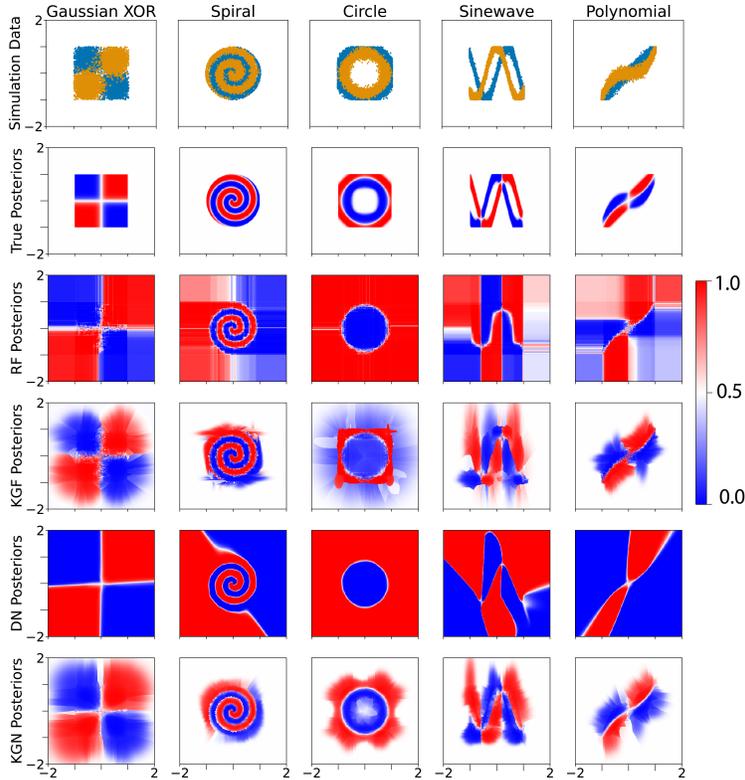


Figure 1: **Visualization of true and estimated posteriors for class 0 from five binary class simulation experiments.** *Row 1:* 10,000 training points with 5,000 samples per class sampled from 5 different simulation setups for binary class classification. The class labels are indicated by yellow and blue colors. *Row 2:* True class conditional posteriors. *Row 3:* Estimated posteriors from random forest. *Row 4:* Estimated posteriors from KGF. *Row 5:* Estimated posteriors from Deep-net. *Row 6:* Estimated posteriors from KGN. The posteriors estimated from KGN and KGF are better calibrated for both in- and out-of-distribution regions compared to those of their parent algorithms.

222 6 Experimental Results

223 We conduct several experiments on two dimensional simulated datasets and OpenML-CC18 data
 224 suite [21] to gain insights on the finite sample performance of KGF and KGN. The details of the
 225 simulation datasets and hyperparameters used for all the experiments are provided in Appendix B.
 226 For the simulation setups, we use classification error, hellinger distance [22, 23] from the true class
 227 conditional posteriors and mean max confidence or posterior [4] as performance statistics. While
 228 measuring in-distribution calibration for the datasets in OpenML-CC18 data suite, as we do not know
 229 the true distribution, we used adaptive calibration error as defined by Nixon et al. [24] with a fixed
 230 bin number of $R = 15$ across all the datasets. Given n OOD samples, we define OOD calibration
 231 error to measure OOD performance for the benchmark datasets as:

$$\left| \frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}} (\hat{P}_{Y|X}(y|\mathbf{x}_i)) - \max_{y \in \mathcal{Y}} (\hat{P}_Y(y)) \right|.$$

232 6.1 Simulation Study

233 Figure 1 top row shows 10000 training samples with 5000 samples per class sampled within the
 234 region $[-1, 1] \times [-1, 1]$ from the five simulation setups described in Appendix B. Therefore, the
 235 empty annular region between $[-1, 1] \times [-1, 1]$ and $[-2, 2] \times [-2, 2]$ is the low density or OOD

<https://www.openml.org/s/99>

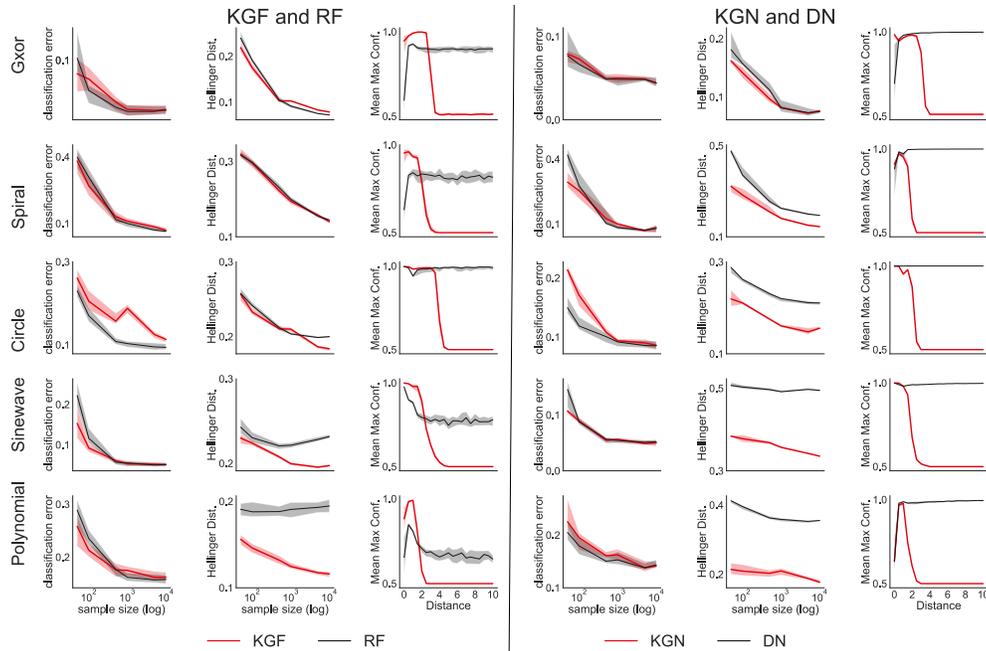


Figure 2: **Classification error, Hellinger distance from true posteriors, mean max confidence or posterior for the simulation experiments.** The median performance is shown as a dark curve with shaded region as error bars showing the 25-th and the 75-th percentile. KGF (**Left block**) and KGN (**Right block**) improve both in- and out-of-distribution calibration of their respective parent algorithms while maintaining nearly similar classification accuracy on the simulation datasets.

236 region in Figure 1. The corresponding true posteriors $\mathbb{P}[Y = 0|X = \mathbf{x}]$ are shown in the second
 237 row of Figure 1. As shown in Row 3 and 5, RF and DN are really good at estimating the high density
 238 regions of training distribution. However, they overestimate the posteriors in the low density regions
 239 of training distribution. Row 4 and 6 of Figure 1 show KGF and KGN improves the posterior estimation
 240 specially in the low density of the training distribution or OOD regions of the feature space. Because
 241 of axis aligned split in random forest RF and thereby, KGF are less efficient in learning non-linear
 242 decision boundaries like spiral, circle and sinewave simulations than ReLU -net and KGN. Figure 2
 243 quantifies the performance of the algorithms which are visually represented in Figure 1. KGF and KGN
 244 maintain similar classification accuracy to those of their parent algorithms. We measure hellinger
 245 distance from the true distribution for increasing training sample size within $[-1, 1] \times [-1, 1]$ region
 246 as an index for in-distribution calibration. Column 2 of left and right block in Figure 2 show KGF
 247 and KGN are better at estimating the high density region of training distribution compared to their
 248 parent methods. For measuring OOD performance, we normalize the training data by the maximum
 249 of their l_2 norm so that the training data is confined within a unit circle. For inference, we sample
 250 1000 inference points uniformly from a circle where the circles have increasing radius and plot
 251 mean max posterior for increasing distance from the origin. Therefore, for distance up to 1 we have
 252 in-distribution samples and distances farther than 1 can be considered as OOD region. As shown in
 253 Column 3 of Figure 2, mean max posteriors or confidence for KGF and KGN converge to the maximum
 254 of the class priors, i.e., 0.5 as we go farther away from the training data origin.

255 6.2 Benchmark Data Study

256 We used OpenML-CC18 data suite for benchmark dataset study. We exclude any dataset which
 257 contain categorical features or NaN values and conduct our experiments on 46 datasets with varying
 258 dimensions and sample sizes. For the OOD experiments, we follow a similar setup as that of the
 259 simulation data. We normalize the training data by their maximum l_2 norm and sample 1000 testing
 260 samples uniformly from each hypersphere where the hyperspheres have increasing radius starting
 261 from 1 to 5. Figure 3 shows the summary of performance of the algorithms. The extended results for

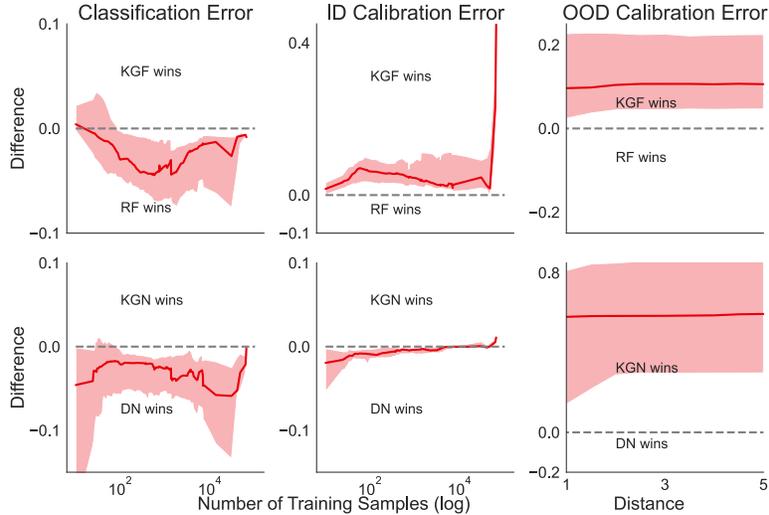


Figure 3: **Performance summary of KGF and KGN on OpenML-CC18 data suite.** The dark red curve in the middle shows the median of performance on 46 datasets. The shaded region shows the error bar consisting of the 25-th and the 75-th percentile of the performance statistics. **Left:** KGF and KGN maintains performance close to their parent algorithms for classification. **Middle:** KGF significantly improves the in-distribution calibration for random forest and KGN improves ReLU-net’s in-distribution calibration for high training sample sizes. **Right:** Both of the proposed approaches yield highly calibrated confidence in the OOD region.

262 each dataset is shown separately in appendix Figure 4, 5, 6, 7, 8 and 9. Figure 3 left column shows on
 263 average KGF and KGN has nearly similar classification accuracy to their respective parent algorithm.
 264 However, according to Figure 3 middle column, KGF improves the in-distribution calibration for
 265 random forest by a huge margin. On the contrary, KGN maintains similar in-distribution calibration
 266 performance to that of its parent ReLU-net. Most interestingly, Figure 3 right column KGN and KGF
 267 improves OOD calibration of their respective parent algorithms by a huge margin.

268 7 Discussion

269 In this paper, we convert deep discriminative models to deep generative models by replacing the affine
 270 function over the polytopes in the discriminative models with a Gaussian kernel. This replacement of
 271 affine function results in better in- and out-of-distribution calibration for our proposed approaches
 272 while maintaining classification accuracy close to the parent algorithm. Theoretically, we show under
 273 certain conditions our approaches asymptotically converge to the true training distribution and this
 274 establishes confidence calibration for learning algorithms in in- and out-of-distribution regions as a
 275 continuum rather than two different problems.

276 For a feature space densely partitioned with small polytopes, we can use Euclidean distance metric in
 277 Equation 4. This is because the Euclidean manifold approximation holds locally for the corresponding
 278 polytope with index r_x^* . On the contrary, for a feature space partitioned with large polytopes,
 279 Euclidean distance measure may be a wrong notion of distance, specially when the underlying
 280 manifold is non-Euclidean. Note that the indicator function in Equation 5 and pruning of the
 281 unpopulated polytopes result in an enlargement of the polytopes in our proposed method compared
 282 to that of the parent model in Equation 3. Therefore, our proposed approach while using Euclidean
 283 metric in Equation 4 may have less classification accuracy compared to that of its parent algorithm.
 284 A correct measure of distance in all the cases including the aforementioned non-Euclidean one would
 285 be the geodesic distance as explored by Madhyastha et al. [25]. We will explore convolutional
 286 neural nets (CNN) trained on image and language benchmark datasets using geodesic distance in
 287 Equation 4 in our future work. Additionally, the proposed approach needs benchmarking against
 288 other calibration approaches in the literature which are mainly based on image datasets and we will
 289 pursue the benchmarking task in our future work.

References

- 290
- 291 [1] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural
292 networks. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International
293 Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*,
294 pages 1321–1330. PMLR, 06–11 Aug 2017.
- 295 [2] Agustinus Kristiadi, Matthias Hein, and Philipp Hennig. Being bayesian, even just a bit, fixes
296 overconfidence in ReLU networks. In Hal Daumé III and Aarti Singh, editors, *Proceedings
297 of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of
298 Machine Learning Research*, pages 5436–5446. PMLR, 13–18 Jul 2020.
- 299 [3] Haoyin Xu, Kaleab A. Kinfu, Will LeVine, Sambit Panda, Jayanta Dey, Michael Ainsworth,
300 Yu-Chung Peng, Madi Kusmanov, Florian Engert, Christopher M. White, Joshua T. Vogelstein,
301 and Carey E. Priebe. When are Deep Networks really better than Decision Forests at small
302 sample sizes, and how? *arXiv preprint arXiv:2108.13637*, 2021.
- 303 [4] Matthias Hein, Maksym Andriushchenko, and Julian Bitterwolf. Why relu networks yield
304 high-confidence predictions far away from the training data and how to mitigate the problem. In
305 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
306 41–50, 2019.
- 307 [5] Alexander Meinke, Julian Bitterwolf, and Matthias Hein. Provably robust detection of out-of-
308 distribution data (almost) for free. *arXiv preprint arXiv:2106.04260*, 2021.
- 309 [6] Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-of-
310 distribution image detection in neural networks. *arXiv preprint arXiv:1706.02690*, 2017.
- 311 [7] Jay Nandy, Wynne Hsu, and Mong Li Lee. Towards maximizing the representation gap between
312 in-domain & out-of-distribution examples. *Advances in Neural Information Processing Systems*,
313 33:9239–9250, 2020.
- 314 [8] Weitao Wan, Yuanyi Zhong, Tianpeng Li, and Jiansheng Chen. Rethinking feature distribution
315 for loss functions in image classification. In *Proceedings of the IEEE conference on computer
316 vision and pattern recognition*, pages 9117–9126, 2018.
- 317 [9] Terrance DeVries and Graham W Taylor. Learning confidence for out-of-distribution detection
318 in neural networks. *arXiv preprint arXiv:1802.04865*, 2018.
- 319 [10] Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier
320 exposure. *arXiv preprint arXiv:1812.04606*, 2018.
- 321 [11] Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High
322 confidence predictions for unrecognizable images. In *Proceedings of the IEEE conference on
323 computer vision and pattern recognition*, pages 427–436, 2015.
- 324 [12] Vikash Sehwal, Arjun Nitin Bhagoji, Liwei Song, Chawin Sitawarin, Daniel Cullina, Mung
325 Chiang, and Prateek Mittal. Better the devil you know: An analysis of evasion attacks using
326 out-of-distribution adversarial examples. *arXiv preprint arXiv:1905.01726*, 2019.
- 327 [13] Jie Ren, Peter J Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark Depristo, Joshua Dillon,
328 and Balaji Lakshminarayanan. Likelihood ratios for out-of-distribution detection. *Advances in
329 neural information processing systems*, 32, 2019.
- 330 [14] Diederik P Kingma, Max Welling, et al. An introduction to variational autoencoders. *Founda-
331 tions and Trends® in Machine Learning*, 12(4):307–392, 2019.
- 332 [15] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil
333 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications
334 of the ACM*, 63(11):139–144, 2020.
- 335 [16] Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, and Balaji Lakshminarayanan.
336 Do deep generative models know what they don’t know? *arXiv preprint arXiv:1810.09136*,
337 2018.

- 338 [17] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural
339 networks. In *International conference on machine learning*, pages 1321–1330. PMLR, 2017.
- 340 [18] Richard Guo, Ronak Mehta, Jesus Arroyo, Hayden Helm, Cencheng Shen, and Joshua T
341 Vogelstein. Estimating information-theoretic quantities with uncertainty forests. *arXiv*, pages
342 arXiv–1907, 2019.
- 343 [19] Meelis Kull, Miquel Perello Nieto, Markus Kängsepp, Telmo Silva Filho, Hao Song, and Peter
344 Flach. Beyond temperature scaling: Obtaining well-calibrated multi-class probabilities with
345 dirichlet calibration. *Advances in neural information processing systems*, 32, 2019.
- 346 [20] Haoyin Xu, Kaleab A Kinfu, Will LeVine, Sambit Panda, Jayanta Dey, Michael Ainsworth,
347 Yu-Chung Peng, Madi Kusmanov, Florian Engert, Christopher M White, et al. When are deep
348 networks really better than decision forests at small sample sizes, and how? *arXiv preprint*
349 *arXiv:2108.13637*, 2021.
- 350 [21] Bernd Bischl, Giuseppe Casalicchio, Matthias Feurer, Pieter Gijsbers, Frank Hutter, Michel
351 Lang, Rafael G Mantovani, Jan N van Rijn, and Joaquin Vanschoren. Openml benchmarking
352 suites. *arXiv preprint arXiv:1708.03731*, 2017.
- 353 [22] Thomas Kailath. The divergence and bhattacharyya distance measures in signal selection. *IEEE*
354 *transactions on communication technology*, 15(1):52–60, 1967.
- 355 [23] C Radhakrishna Rao. A review of canonical coordinates and an alternative to correspondence
356 analysis using hellinger distance. *Qüestió: quaderns d'estadística i investigació operativa*,
357 1995.
- 358 [24] Jeremy Nixon, Michael W Dusenberry, Linchuan Zhang, Ghassen Jerfel, and Dustin Tran.
359 Measuring calibration in deep learning. In *CVPR workshops*, volume 2, 2019.
- 360 [25] Meghana Madhyastha, Gongkai Li, Veronika Strnadová-Neeley, James Browne, Joshua T
361 Vogelstein, Randal Burns, and Carey E Priebe. Geodesic forests. In *Proceedings of the 26th*
362 *ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages
363 513–523, 2020.