

Training OOD Detectors in their Natural Habitats

Julian Katz-Samuels^{*1} Julia Nakhleh^{*2} Robert Nowak³ Yixuan Li²

Abstract

Out-of-distribution (OOD) detection is important for machine learning models deployed in the wild. Recent methods use auxiliary outlier data to regularize the model for improved OOD detection. However, these approaches make a strong distributional assumption that the auxiliary outlier data is completely separable from the in-distribution (ID) data. In this paper, we propose a novel framework that leverages wild mixture data—that naturally consists of both ID and OOD samples. Such wild data is abundant and arises freely upon deploying a machine learning classifier in their *natural habitats*. Our key idea is to formulate a constrained optimization problem and to show how to tractably solve it. Our learning objective maximizes the OOD detection rate, subject to constraints on the classification error of ID data and on the OOD error rate of ID examples. We extensively evaluate our approach on common OOD detection tasks and demonstrate superior performance. Code is available at https://github.com/jkatzsam/woods_ood.

1. Introduction

Out-of-distribution (OOD) detection has become a central challenge in safely deploying machine learning models in the wild, where test-time data can naturally arise from a mixture distribution of both knowns and unknowns (Bendale & Boult, 2015). Concerningly, modern neural networks are shown to produce overconfident and therefore untrustworthy predictions for unknown OOD inputs (Nguyen et al., 2015). To mitigate the issue, recent works have explored training with an auxiliary outlier dataset, where the model is regularized to produce lower confidence (Hendrycks et al., 2019) or higher energies (Liu et al., 2020) on the outlier data. These

methods have demonstrated encouraging OOD detection performance over the counterpart without auxiliary data.

Despite this promise, there are two primary limitations to the existing methods. First, the auxiliary data distribution collected offline may not be a good match for the true distribution of unknown data in the wild, thus the learned model may fail to detect deployment-time OOD data. Second, collecting such data can be very labor-intensive and inflexible, and necessitates careful data cleaning to ensure the auxiliary outlier data does not overlap with the ID data. We address these challenges by leveraging unlabeled “in-the-wild” data—which can be collected almost *for free* upon deploying a machine learning classifier in the open world, and has been largely overlooked for OOD learning purposes. Such data is available in abundance, does not require any human annotation, and is often a much better match to the true test time distribution than data collected offline. While this setting naturally suits many real-world applications, it also poses unique challenges since the wild data distribution is not pure and consists of both ID data and OOD data.

In this paper, we propose a novel framework that enables effectively exploiting unlabeled in-the-wild data for OOD detection. Unlabeled wild data is frequently available since it is produced essentially whenever deploying an existing classifier in a real-world system. This setting can be viewed as training OOD detectors in their *natural habitats*.

Our learning framework revolves around building the OOD classifier using only labeled ID data from \mathbb{P}_{in} and unlabeled wild data from \mathbb{P}_{wild} , which can be considered to be a mixture of \mathbb{P}_{in} and an unknown (OOD) distribution. To deal with the lack of a “clean” set of OOD examples, our key idea is to formulate a constrained optimization problem. In a nutshell, our learning objective aims to minimize the error of classifying data from \mathbb{P}_{wild} as ID, subject to two constraints: (i) the error of declaring an ID data from \mathbb{P}_{in} as OOD must be low, and (ii) the multi-class classification model must maintain the best-achievable accuracy (or close to it) of a baseline classifier designed without an OOD detection requirement. Even though our framework does not have access to a “clean” OOD dataset, we show both empirically and theoretically that it can learn to accurately detect OOD examples. Our work is inspired by the semi-supervised nov-

^{*}Equal contribution ¹Institute for Foundations of Data Science, University of Wisconsin, Madison ²Department of Computer Sciences, University of Wisconsin, Madison ³Department of Electrical and Computer Engineering, University of Wisconsin, Madison. Correspondence to: Julian Katz-Samuels <katzsamuels@wisc.edu>.

elty detection (SSND) framework (Blanchard et al., 2010), yet differs by simultaneously considering both classification and OOD detection tasks.

Beyond the mathematical framework, a key contribution of our paper is a constrained optimization implementation of the framework for deep neural networks. We propose a novel training procedure based on the augmented Lagrangian method, or ALM (Hestenes, 1969). While ALM is an established approach to optimization with functional constraints, its adaptation to modern deep learning is not straightforward or common. In particular, we adapt ALM to our problem setting with inequality constraints, obtaining an end-to-end training algorithm using stochastic gradient descent. Unlike methods that add a regularization term to the training objective, our method via constrained optimization offers strong guarantees (c.f. Section 3.1).

We extensively evaluate our approach on common OOD detection tasks and establish state-of-the-art performance. For completeness, we compare with two groups of approaches: (1) trained with only \mathbb{P}_{in} data, and (2) trained with both \mathbb{P}_{in} and an auxiliary dataset. On CIFAR-100, compared to a strong baseline using only \mathbb{P}_{in} , our method outperforms by **48.10%** (FPR95) on average. The performance gain precisely demonstrates the advantage of incorporating unlabeled wild data. Our method also outperforms Outlier Exposure (OE) (Hendrycks et al., 2019) by **7.36%** in FPR95. Our key contributions are summarized as follows:

- We propose a novel OOD detection framework via constrained optimization (dubbed WOODS, Wild OOD detection sans-Supervision), capable of exploiting unlabeled wild data. We show how to integrate constrained optimization into modern deep neural networks and solve it tractably.
- We provide novel theoretical insights that support WOODS, in particular for the choice of loss functions.
- We perform extensive ablations and comparisons under: (1) a diverse range of datasets, (2) different mixture ratios π of ID and OOD in \mathbb{P}_{wild} , and (3) different assumptions on the relationship between the wild distribution \mathbb{P}_{wild} and the test-time distribution. Our method establishes *state-of-the-art* results, significantly outperforming existing methods.

2. Problem Setup

Labeled In-distribution Data Let $\mathcal{X} = \mathbb{R}^d$ denote the input space and $\mathcal{Y} = \{1, \dots, K\}$ denote the label space. We assume access to the labeled training set $\mathcal{D}_{\text{in}}^{\text{train}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, drawn *i.i.d.* from the joint data distribution $\mathbb{P}_{\mathcal{X}\mathcal{Y}}$. Let \mathbb{P}_{in} denote the marginal distribution on \mathcal{X} , which is also referred to as the *in-distribution* (ID). Let

$f_\theta : \mathcal{X} \mapsto \mathbb{R}^{|\mathcal{Y}|}$ denote a function for the classification task, which predicts the label of an input sample.

Out-of-distribution Detection When deploying machine learning models in the real world, a reliable classifier should not only accurately classify known ID samples, but also identify as “unknown” any *out-of-distribution* (OOD) input—samples from a different distribution $\mathbb{P}_{\text{out}}^{\text{test}}$ that the model has not been exposed to during training. This can be achieved through having an OOD classifier, in addition to the multi-class classifier f_θ . Samples detected as OOD will be rejected; samples detected as ID will be classified by f_θ .

OOD detection can be formulated as a binary classification problem. At test time, the goal is to decide whether a test-time input $\mathbf{x} \in \mathcal{X}$ is from the in-distribution \mathbb{P}_{in} (ID) or not (OOD). We denote $g_\theta : \mathcal{X} \mapsto \{\text{in}, \text{out}\}$ as the function mapping for OOD detection.

Unlabeled in-the-wild Data A major challenge in OOD detection is the lack of labeled examples of OOD. In particular, the sample space for potential OOD data can be prohibitively large, making it expensive to collect labeled OOD data. In this paper, we incorporate unlabeled in-the-wild samples $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m$ into OOD detection. These samples consist of potentially both ID and OOD data, and can be collected almost for free upon deploying an existing classifier (say f_θ) in its natural habitats. We use the Huber contamination model (Huber, 1964) to model the marginal distribution of the wild data:

$$\mathbb{P}_{\text{wild}} := (1 - \pi)\mathbb{P}_{\text{in}} + \pi\mathbb{P}_{\text{out}},$$

where $\pi \in (0, 1]$. We note that the case $\pi = 0$ is technically possible in controlled deployment environments, but it is straightforward since no novelties occur and we therefore do not consider it.

Goal: Our learning framework revolves around building the OOD classifier g_θ and multi-class classifier f_θ by leveraging data from both \mathbb{P}_{in} and \mathbb{P}_{wild} . We use the shared parameters θ to indicate the fact that they may share the neural network parameters. In testing, we measure the following errors:

$$\begin{aligned} \downarrow \text{FPR}(g_\theta) &:= \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{out}}^{\text{test}}}(\mathbb{1}\{g_\theta(\mathbf{x}) = \text{in}\}), \\ \uparrow \text{TPR}(g_\theta) &:= \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{in}}}(\mathbb{1}\{g_\theta(\mathbf{x}) = \text{in}\}), \\ \uparrow \text{Acc}(f_\theta) &:= \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_{\mathcal{X}\mathcal{Y}}}(\mathbb{1}\{f_\theta(\mathbf{x}) = y\}), \end{aligned}$$

where $\mathbb{1}\{\cdot\}$ is the indicator function and the arrows indicate higher/lower is better. We distinguish between the OOD training distribution \mathbb{P}_{out} and the OOD test distribution $\mathbb{P}_{\text{out}}^{\text{test}}$ because as is common in out-of-distribution detection, the data available for training may be different from the data at test time. We call the case $\mathbb{P}_{\text{out}}^{\text{test}} = \mathbb{P}_{\text{out}}$ the *stationary wild*

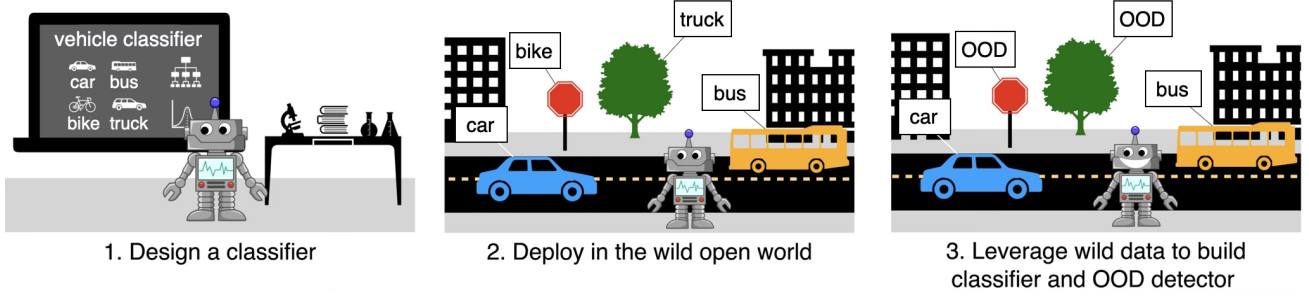


Figure 1: Overview of our learning framework. A model pre-trained to perform in-distribution (ID) classification can be deployed in the “wild” open world, where it will encounter large quantities of unlabeled ID *and* OOD data. Using our WOODS approach, the model can be fine-tuned using the “wild” data to perform accurate OOD detection as well as ID classification.

setting and the case with $\mathbb{P}_{\text{out}}^{\text{test}} \neq \mathbb{P}_{\text{out}}$ the *non-stationary wild setting* and will consider both later in Section 5.

3. Method: Out-of-distribution Learning via Constrained Optimization

In this section, we present a novel framework that performs out-of-distribution learning with unlabeled data in the wild. Our framework offers substantial advantages over the counterpart approaches that rely only on the ID data, and naturally suits many applications where machine learning models are deployed in the open world.

To exploit the in-the-wild data, our key idea is to formulate a constrained optimization problem (Section 3.1). Moreover, we show how to integrate this constrained optimization problem into modern neural networks and solve it tractably using the Augmented Lagrangian Method (Section 3.2).

3.1. Learning Objective

In a nutshell, we formulate the learning objective by aiming to minimize the error of classifying data from \mathbb{P}_{out} as ID, subject to (i) the error of declaring an ID as OOD is at most α , and (ii) the multi-class classification model meets some error threshold τ . Mathematically, this can be formalized as a constrained optimization problem:

Objective Overview Given $\alpha, \tau \in [0, 1]$, we aim to optimize:

$$\begin{aligned} \inf_{\theta} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{out}}} (\mathbb{1}\{g_{\theta}(\mathbf{x}) = \text{in}\}) \\ \text{s.t. } \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{in}}} (\mathbb{1}\{g_{\theta}(\mathbf{x}) = \text{out}\}) \leq \alpha \\ \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_{\mathcal{X}\mathcal{Y}}} (\mathbb{1}\{f_{\theta}(\mathbf{x}) \neq y\}) \leq \tau. \end{aligned} \quad (1)$$

However, we never observe a clean dataset from \mathbb{P}_{out} , making it difficult to directly solve (1). To sidestep this issue,

we reformulate the learning objective as follows:

$$\begin{aligned} \inf_{\theta} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{wild}}} (\mathbb{1}\{g_{\theta}(\mathbf{x}) = \text{in}\}) \\ \text{s.t. } \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{in}}} (\mathbb{1}\{g_{\theta}(\mathbf{x}) = \text{out}\}) \leq \alpha \\ \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_{\mathcal{X}\mathcal{Y}}} (\mathbb{1}\{f_{\theta}(\mathbf{x}) \neq y\}) \leq \tau. \end{aligned} \quad (2)$$

where we replaced the OOD distribution \mathbb{P}_{out} in the objective with \mathbb{P}_{wild} , a distribution that we observe an *i.i.d.* dataset from. As we explain in more detail shortly, optimizing (2) is sufficient under mild conditions to solve the original optimization problem (1).

Empirically, we can solve the optimization problem (2) by minimizing the number of samples $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m$ from the wild distribution \mathbb{P}_{wild} that are labeled as ID, subject to (i) labeling at least $1 - \alpha$ of the ID samples $\mathbf{x}_1, \dots, \mathbf{x}_n$ correctly, and (ii) achieving the classification performance threshold. Equivalently, we consider solving:

$$\begin{aligned} \inf_{\theta} \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{g_{\theta}(\tilde{\mathbf{x}}_i) = \text{in}\} \\ \text{s.t. } \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{g_{\theta}(\mathbf{x}_i) = \text{out}\} \leq \alpha \\ \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{f_{\theta}(\mathbf{x}_i) \neq y_i\} \leq \tau. \end{aligned} \quad (3)$$

Surrogate Problem Note that the above objective in (3) is intractable due to the 0/1 loss and here we propose a tractable relaxation, replacing the 0/1 loss with a surrogate

loss as follows:

$$\begin{aligned} & \inf_{\theta} \frac{1}{m} \sum_{i=1}^m \mathcal{L}_{\text{ood}}(g_{\theta}(\tilde{\mathbf{x}}_i), \text{in}) \\ & \text{s.t. } \frac{1}{n} \sum_{j=1}^n \mathcal{L}_{\text{ood}}(g_{\theta}(\mathbf{x}_j), \text{out}) \leq \alpha \\ & \quad \frac{1}{n} \sum_{j=1}^n \mathcal{L}_{\text{cls}}(f_{\theta}(\mathbf{x}_j), y_j) \leq \tau \end{aligned} \quad (4)$$

where \mathcal{L}_{ood} denotes the loss of the binary OOD classifier and \mathcal{L}_{cls} denotes a loss for the classification task. Our framework is general and can be instantiated by different forms of loss functions, for which we describe details later in Section 4.

Here, we state a theoretical result justifying the optimization problem (4) where we use tractable losses, specifically using the sigmoid loss $\sigma(t) = \frac{1}{1+e^{-t}}$ for \mathcal{L}_{ood} and the hinge loss for \mathcal{L}_{cls} . We suppose that the weights are p -dimensional and belong to a subset $\theta \in \Theta \subset \mathbb{R}^p$. Let opt denote the value to the population-level optimization problem of interest:

$$\begin{aligned} & \inf_{\theta \in \Theta} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{out}}} \mathcal{L}_{\text{ood}}(g_{\theta}(\tilde{\mathbf{x}}_i), \text{in}) \\ & \text{s.t. } \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{in}}} \mathcal{L}_{\text{ood}}(g_{\theta}(\mathbf{x}_j), \text{out}) \leq \alpha \\ & \quad \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_{\mathcal{X}\mathcal{Y}}} \mathcal{L}_{\text{cls}}(f_{\theta}(\mathbf{x}_j), y_j) \leq \tau. \end{aligned} \quad (5)$$

Proposition 3.1. *Suppose $K = 2$. Suppose $f_{\theta}(x) : \mathbb{R}^d \mapsto [0, 1]$. Suppose $\mathcal{L}_{\text{ood}}(t, \text{in}) = \sigma(-t)$, $\mathcal{L}_{\text{ood}}(t, \text{out}) = \sigma(t)$, and $\mathcal{L}_{\text{cls}}(s, y) = y \log(\frac{1}{s}) + (1 - y) \log(\frac{1}{1-s})$ for $s \in [0, 1]$. Define $\epsilon_k := \sqrt{\frac{2 \ln(6/\delta)}{k}} + 2 \max_{h \in \{f, g\}} \max_{\mathbb{P} \in \{\mathbb{P}_1, \mathbb{P}_{\text{wild}}\}} \mathbb{E}_{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m \sim \mathbb{P}} \mathbb{E}_{\eta_1, \dots, \eta_m \theta \in \Theta} \sup_{\theta \in \Theta} \frac{1}{k} \sum_{i=1}^k \eta_i h_{\theta}(\tilde{\mathbf{x}}_i)$ where η_1, \dots, η_k are i.i.d. and $\mathbb{P}(\eta_i = 1) = \mathbb{P}(\eta_i = -1) = 1/2$. Let $\hat{\theta}_{\epsilon}$ solve (4) with some tolerance ϵ (see the Appendix for a formal statement). Then, there exist universal positive constants c_1, c_2, c_3 such that under a mild condition (see the appendix), with probability at least $1 - \delta$*

1. $\mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{out}}} \sigma(-\hat{g}_{\hat{\theta}_{\epsilon}}(\mathbf{x})) \leq \text{opt} + c_1 \pi^{-1}(\epsilon_n + \epsilon_m),$
2. $\mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{in}}} \sigma(g_{\hat{\theta}_{\epsilon}}(\mathbf{x})) \leq \alpha + c_2 \epsilon_n, \text{ and}$
3. $\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_{\mathcal{X}\mathcal{Y}}} \mathcal{L}_{\text{cls}}(f_{\hat{\theta}_{\epsilon}}(\mathbf{x}), y) \leq \tau + c_3 \epsilon_n.$

The above Proposition shows that as long as the Rademacher complexities of the function classes $\{g_{\theta} : \theta \in \Theta\}$ and $\{f_{\theta} : \theta \in \Theta\}$ decay at a suitable rate, then solving (4) gives a solution that approaches feasibility and optimality for (5), the optimization problem of interest. Due to space constraints, we defer a full Proposition statement with additional details and proof to the Appendix. The above

Proposition extends Theorem 1 of (Blanchard et al., 2010) from the computationally intractable 0/1 loss to the sigmoid loss, a tractable, differentiable loss that we use in our experiments. In addition, we replace the VC dimension in their Theorem with the Rademacher complexity of the function class, a much more fine-grained measure of the function class complexity. We note that it is also possible to prove an analogue of Proposition 3.1 for the 0/1 loss.

Although deep neural networks tend to have a trivial Rademacher complexity scaling with the number of data points due to their ability to interpolate random labels (Zhang et al., 2021), our training approach has several regularizing features that limit the expressivity of the deep neural network. First, in our training procedure described in Section 5, we add a weight decay term that penalizes deep neural networks with very large weights and limits the network’s expressiveness. Second, we perform the classification task and the OOD detection task at once and we conjecture that requiring our algorithm to have good classification performance regularizes the model for the OOD task. Finally, we also fine-tune from a network pre-trained on the classification task with a relatively small learning rate, limiting the optimization procedure’s ability to arrive at solutions interpolating random labels and therefore adding further regularization.

Having theoretically justified the optimization problem (4), we next show how to optimize it.

3.2. Solving the Constrained Optimization

In this subsection, we first provide background on the Augmented Lagrangian method, and then discuss how it can solve our constrained optimization problem.

Augmented Lagrangian Method (ALM) Augmented Lagrangian method (Hestenes, 1969) is an established approach to optimization with functional constraints. ALM improves over two other related methods for constrained optimization: the penalty method and the method of Lagrangian multipliers. While the penalty method suffers from ill-conditioning (Nocedal & Wright, 2006), the method of Lagrangian multipliers is specific to the convex case (Rockafellar, 1973). In this paper, we adapt ALM to our setting with inequality constraints, and later show that it can be optimized end-to-end with modern neural networks.

To provide some background, we consider the following constrained optimization problem as an example:

$$\begin{aligned} & \min_{\theta \in \mathbb{R}^p} f(\theta) \\ & \text{s.t. } c_i(\theta) \leq 0 \forall i \in [q], \end{aligned} \quad (6)$$

where f and c_1, \dots, c_q are convex. ALM solves the constrained optimization problem in (6) by converting it into a

Algorithm 1 WOODS (Wild OOD detection sans-Supervision)

```

1: Input:  $\theta_{(1)}^{(1)}, \lambda_{(1)}^{(1)}, \beta_1, \beta_2$ , epoch length  $T$ , batch size  $B$ ,
   learning rate  $\mu_1$ , learning rate  $\mu_2$ , penalty multiplier  $\gamma$ ,
   tol
2: for epoch = 1, 2, ... do
3:   for  $t = 1, 2, \dots, T - 1$  do
4:     Sample a batch of data, calculate  $\mathcal{L}_{\beta}^{\text{batch}}(\theta, \lambda)$ 
5:      $\theta_{(\text{epoch})}^{(t+1)} \leftarrow \theta_{(\text{epoch})}^{(t)} - \mu_1 \nabla_{\theta} \mathcal{L}_{\beta}^{\text{batch}}(\theta, \lambda)$ 
6:   end for
7:    $\lambda^{(\text{epoch}+1)} \leftarrow \lambda^{(\text{epoch})} + \mu_2 \nabla_{\lambda} \mathcal{L}_{\beta}(\theta_{(\text{epoch})}^{(T)}, \lambda^{(\text{epoch})})$ 
8:   if  $\frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\text{ood}}(g_{\theta^{(T)}}(\mathbf{x}_i), \text{out}) > \alpha + \text{tol}$  then
9:      $\beta_1 \leftarrow \gamma \beta_1$ 
10:  end if
11:  if  $\frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\text{cls}}(f_{\theta^{(T)}}(\mathbf{x}_i), y_i) > \tau + \text{tol}$  then
12:     $\beta_2 \leftarrow \gamma \beta_2$ 
13:  end if
14:   $\theta_{(\text{epoch}+1)}^{(1)} \leftarrow \theta_{(\text{epoch})}^{(T)}$ 
15: end for
    
```

sequence of unconstrained optimization problems.

Define the the classical augmented Lagrangian (AL) function

$$\mathcal{L}_{\beta}(\theta, \lambda) = f(\theta) + \sum_{i=1}^q \psi_{\beta}(c_i(\theta), \lambda_i)$$

where

$$\psi_{\beta}(u, v) = \begin{cases} uv + \frac{\beta}{2}u^2 & \beta u + v \geq 0 \\ -\frac{v^2}{2\beta} & \text{o/w} \end{cases},$$

$\lambda = (\lambda_1, \dots, \lambda_q)^{\top}$, and $\beta > 0$. At iteration k , ALM minimizes the augmented Lagrangian function with respect to θ and then performs a gradient ascent update step on λ (Xu, 2021a;b):

1. $\theta^{(k+1)} \leftarrow \text{argmin}_{\theta} \mathcal{L}_{\beta_k}(\theta, \lambda^{(k)})$
2. $\lambda^{(k+1)} \leftarrow \lambda^{(k)} + \rho \nabla_{\lambda} \mathcal{L}_{\beta_k}(\theta^{(k+1)}, \lambda)$,

where ρ is a learning rate for the dual variable λ and $\{\beta_k\}_k$ is a sequence of penalty parameters. The sequence of penalty parameters $\{\beta_k\}_k$ may be chosen beforehand or adapted based on the optimization process.

Our Algorithm Algorithm 1 presents our approach to using ALM to solve (4). We define the augmented Lagrangian

function as:

$$\begin{aligned} \mathcal{L}_{\beta}(\theta, \lambda) = & \frac{1}{m} \sum_{i=1}^m \mathcal{L}_{\text{ood}}(g_{\theta}(\tilde{\mathbf{x}}_i), \text{in}) \\ & + \psi_{\beta_1} \left(\frac{1}{n} \sum_{j=1}^n \mathcal{L}_{\text{ood}}(g_{\theta}(\mathbf{x}_j), \text{out}) - \alpha, \lambda_1 \right) \\ & + \psi_{\beta_2} \left(\frac{1}{n} \sum_{j=1}^n \mathcal{L}_{\text{cls}}(f_{\theta}(\mathbf{x}_j), y_j) - \tau, \lambda_2 \right), \end{aligned}$$

where $\beta = (\beta_1, \beta_2)^{\top}$.

Adaptation to Stochastic Gradient Descent We show that our framework can be adapted to the stochastic case, which is more amenable for training with modern neural networks. We outline the full process in Algorithm 1. In each iteration, we calculate the per-batch loss as follows:

$$\begin{aligned} \mathcal{L}_{\beta}^{\text{batch}}(\theta, \lambda) = & \frac{1}{B} \sum_{i \in I} \mathcal{L}_{\text{ood}}(g_{\theta}(\tilde{\mathbf{x}}_i), \text{in}) \\ & + \psi_{\beta_1} \left(\frac{1}{B} \sum_{j \in J} \mathcal{L}_{\text{ood}}(g_{\theta}(\mathbf{x}_j), \text{out}) - \alpha, \lambda_1^{(\text{epoch})} \right) \\ & + \psi_{\beta_2} \left(\frac{1}{B} \sum_{j \in J} \mathcal{L}_{\text{cls}}(f_{\theta}(\mathbf{x}_j), y_j) - \tau, \lambda_2^{(\text{epoch})} \right), \end{aligned} \quad (7)$$

where I and J denote the set of mini-batch of size B , randomly sampled from the wild data and ID data respectively. Since $\psi(u, v)$ is convex in u , by Jensen's inequality, the objective function in (7) is an upper bound on $\mathcal{L}_{\beta}(\theta, \lambda^{(\text{epoch})})$. This step therefore approximates ALM; indeed, it is not straightforward to adapt ALM to the stochastic case (Yan & Xu, 2020). At the end of the epoch, it performs a gradient ascent update on λ (see line 7). Finally, in lines 9 and 12, it increases the constraints weight penalties β_1 and β_2 by a penalty multiplier $\gamma > 1$.

4. Loss Functions with Neural Networks

In this section, we discuss how to realize our learning framework in the context of modern neural networks. Concretely, we address how to define the loss functions \mathcal{L}_{cls} and \mathcal{L}_{ood} .

Classification Loss \mathcal{L}_{cls} We consider a neural network parameterized by θ , which encodes an input $\mathbf{x} \in \mathbb{R}^d$ to a feature space with dimension r . We denote by $h_{\theta}(\mathbf{x}) \in \mathbb{R}^r$ the feature vector from the penultimate layer of the network. A weight matrix $\mathbf{W} \in \mathbb{R}^{r \times K}$ connects the feature $h_{\theta}(\mathbf{x})$ to the output $f_{\theta}(\mathbf{x}) = \mathbf{W}^{\top} h_{\theta}(\mathbf{x}) \in \mathbb{R}^K$. The per-sample classification loss \mathcal{L}_{cls} can be defined using the cross-entropy (CE) loss:

$$\mathcal{L}_{\text{cls}}(f_{\theta}(\mathbf{x}), y) = -\log \frac{e^{f_{\theta}^{(y)}(\mathbf{x})}}{\sum_{j=1}^K e^{f_{\theta}^{(j)}(\mathbf{x})}}, \quad (8)$$

where $f_{\theta}^{(y)}(\mathbf{x})$ denotes the y -th element of $f_{\theta}(\mathbf{x})$ corresponding to the label y .

Binary Loss \mathcal{L}_{ood} The loss function \mathcal{L}_{ood} should ideally optimize for the separability between the ID vs. OOD data under some function that captures the data density. However, directly estimating $\log p(\mathbf{x})$ can be computationally intractable as it requires sampling from the entire space \mathcal{X} . We note that the log partition function $E_{\theta}(\mathbf{x}) := -\log \sum_{j=1}^K e^{f_{\theta}^{(j)}(\mathbf{x})}$ is proportional to $-\log p(\mathbf{x})$ with some unknown factor, which can be seen from the following:

$$p(y|\mathbf{x}) = \frac{p(\mathbf{x}, y)}{p(\mathbf{x})} = \frac{e^{f_{\theta}^{(y)}(\mathbf{x})}}{\sum_{j=1}^K e^{f_{\theta}^{(j)}(\mathbf{x})}}.$$

The negative log partition function is also known as the free energy, which was shown to be an effective uncertainty measurement for OOD detection (Liu et al., 2020).

Our idea is to explicitly optimize for a level-set estimation based on the energy function (threshold at 0), where the ID data has negative energy values and vice versa.

$$\begin{aligned} \arg\min_{\theta} & \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{E_{\theta}(\tilde{\mathbf{x}}_i) \leq 0\} \\ \text{s.t.} & \frac{1}{n} \sum_{j=1}^n \mathbb{1}\{E_{\theta}(\mathbf{x}_j) \geq 0\} \leq \alpha \\ & \frac{1}{n} \sum_{j=1}^n \mathcal{L}_{\text{cls}}(f_{\theta}(\mathbf{x}_j), y_j) \leq \tau \end{aligned}$$

Since the 0/1 loss is intractable, we replace it with the binary sigmoid loss, a smooth approximation of the 0/1 loss, yielding the following optimization problem:

$$\begin{aligned} \arg\min_{\theta, w \in \mathbb{R}} & \frac{1}{m} \sum_{i=1}^m \frac{1}{1 + \exp(-w \cdot E_{\theta}(\tilde{\mathbf{x}}_i))} \\ \text{s.t.} & \frac{1}{n} \sum_{j=1}^n \frac{1}{1 + \exp(w \cdot E_{\theta}(\mathbf{x}_j))} \leq \alpha \\ & \frac{1}{n} \sum_{j=1}^n \mathcal{L}_{\text{cls}}(f_{\theta}(\mathbf{x}_j), y_j) \leq \tau. \end{aligned} \quad (9)$$

Here w is a learnable parameter modulating the slope of the sigmoid function. Now we may apply the same approach as in section 3 to solve the constrained optimization problem (9). In effect, we have

$$\mathcal{L}_{\text{ood}}(g_{\theta}(\tilde{\mathbf{x}}_i), \text{in}) = \frac{1}{1 + \exp(-w \cdot E_{\theta}(\tilde{\mathbf{x}}_i))}.$$

This loss function is originally developed in (Du et al., 2022b) for model regularization. Our loss has two notable

advancements over energy-regularized learning (ERL) (Liu et al., 2020): (1) We consider a more general unsupervised setting where the wild data distribution \mathbb{P}_{wild} contains both ID and OOD data, whereas ERL assumes having access to an auxiliary outlier dataset that is completely separable from the ID data. Methods like OE (Hendrycks et al., 2019) in general require performing manual data collection and cleaning, which is more restrictive. (2) We formalize the learning objective as a constrained optimization, which offers strong guarantees (see Section 3.1).

5. Experiments

5.1. Experimental Setup

Datasets Following the common benchmarks in literature, we use CIFAR-10 and CIFAR-100 (Krizhevsky et al., 2009) as ID datasets (\mathbb{P}_{in}). For OOD test datasets $\mathbb{P}_{\text{out}}^{\text{test}}$, we use a suite of natural image datasets including SVHN (Netzer et al., 2011), Textures (Cimpoi et al., 2014), Places365 (Zhou et al., 2018), LSUN-Crop (Yu et al., 2016), and LSUN-Resize (Yu et al., 2016).

To simulate the wild data \mathbb{P}_{wild} , we mix a subset of ID data (as \mathbb{P}_{in}) with the outlier dataset (as \mathbb{P}_{out}) under various $\pi \in \{0.05, 0.1, 0.2, 0.5, 1.0\}$. We consider the stationary setting where $\mathbb{P}_{\text{out}}^{\text{test}} = \mathbb{P}_{\text{out}}$, and the non-stationary setting where $\mathbb{P}_{\text{out}}^{\text{test}} \neq \mathbb{P}_{\text{out}}$. When $\mathbb{P}_{\text{out}}^{\text{test}} = \mathbb{P}_{\text{out}}$ (say from SVHN), we train the model using CIFAR+SVHN as the wild data and test on SVHN as OOD. When $\mathbb{P}_{\text{out}}^{\text{test}} \neq \mathbb{P}_{\text{out}}$, for \mathbb{P}_{out} , we use the publicly available 300K Random Images (Hendrycks et al., 2019), a subset of the original 80 Million Tiny Images dataset, and test on the 5 OOD datasets mentioned above. Note that we split CIFAR datasets into two halves: 25,000 images as ID training data, and remainder 25,000 used to create the wild mixture data.

Evaluation Metrics To evaluate the methods, we use the standard measures in the literature: the false positive rate of declaring OOD examples as ID when 95% of ID data points are declared as ID (FPR95) and the area under the receiver operating characteristic curve (AUROC).

Training Details For all experiments and methods, we use the Wide ResNet (Zagoruyko & Komodakis, 2016) architecture with 40 layers and widen factor of 2. The model is optimized using stochastic gradient descent with Nesterov momentum (Duchi et al., 2011). We set the weight decay coefficient to be 0.0005, and momentum to be 0.09. Models are initialized with a model pre-trained on the CIFAR-10/CIFAR-100 data and trained for 100 epochs in the $\mathbb{P}_{\text{out}}^{\text{test}} = \mathbb{P}_{\text{out}}$ setting and for 50 epochs in the $\mathbb{P}_{\text{out}}^{\text{test}} \neq \mathbb{P}_{\text{out}}$ setting. Our initialization scheme from a pre-trained model naturally suits our setting (c.f. Section 2), where an existing classifier in deployment is available.

Training OOD Detectors in their Natural Habitats

Method	SVHN		LSUN-R		LSUN-C		Textures		Places365		Average		Acc.
	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	
With \mathbb{P}_{in} only													
MSP	84.59	71.44	82.42	75.38	66.54	83.79	83.29	73.34	82.84	73.78	79.94	75.55	75.96
ODIN	84.66	67.26	71.96	81.82	55.55	87.73	79.27	73.45	87.88	71.63	75.86	76.38	75.96
Energy	85.82	73.99	79.47	79.23	35.32	93.53	79.41	76.28	80.56	75.44	72.12	79.69	75.96
Mahalanobis	57.52	86.01	21.23	96.00	91.18	69.69	39.39	90.57	88.83	67.87	59.63	82.03	75.96
GODIN	83.38	84.05	62.24	88.22	72.86	83.84	83.83	78.91	80.56	76.14	76.57	82.23	75.33
CSI	64.70	84.97	91.55	63.42	38.10	92.52	74.70	92.66	82.25	73.63	70.26	81.44	69.90
With \mathbb{P}_{in} and \mathbb{P}_{wild}													
OE	1.57 \pm 0.1	99.63 \pm 0.0	0.93 \pm 0.2	99.79 \pm 0.0	3.83 \pm 0.4	99.26 \pm 0.1	27.89 \pm 0.5	93.35 \pm 0.2	60.24 \pm 0.6	83.43 \pm 0.6	18.89 \pm 0.4	95.09 \pm 0.2	71.65 \pm 0.4
Energy (w/ OE)	1.47 \pm 0.3	99.68 \pm 0.0	2.68 \pm 1.9	99.50 \pm 0.3	2.52 \pm 0.4	99.44 \pm 0.1	37.26 \pm 9.1	91.26 \pm 2.5	54.67 \pm 1.0	86.09 \pm 0.4	19.72 \pm 2.5	95.19 \pm 0.7	73.46 \pm 0.8
WOODS (ours)	0.52 \pm 0.1	99.88 \pm 0.0	0.38 \pm 0.1	99.92 \pm 0.0	0.93 \pm 0.2	99.77 \pm 0.0	17.92 \pm 0.5	96.44 \pm 0.2	37.90 \pm 0.6	91.22 \pm 0.3	11.53 \pm 0.3	97.45 \pm 0.1	74.79 \pm 0.2
WOODS-alt (ours)	0.12 \pm 0.0	99.96 \pm 0.0	0.07 \pm 0.1	99.96 \pm 0.0	0.11 \pm 0.0	99.96 \pm 0.0	9.12 \pm 0.3	96.65 \pm 0.1	29.58 \pm 0.4	90.60 \pm 0.3	7.80 \pm 0.5	97.43 \pm 0.5	75.22 \pm 0.2

Table 1. Main results when $\mathbb{P}_{out}^{test} = \mathbb{P}_{out}$. Comparison with competitive OOD detection methods on CIFAR-100. For methods using \mathbb{P}_{wild} , we train under the same dataset and same $\pi = 0.1$. For each dataset, we create corresponding wild mixture distribution $\mathbb{P}_{wild} := (1 - \pi)\mathbb{P}_{in} + \pi\mathbb{P}_{out}$ for training and test on the corresponding OOD dataset. \uparrow indicates larger values are better and vice versa. $\pm x$ denotes the standard error, rounded to the first decimal point.

Method	OOD Dataset												Acc.
	SVHN		LSUN-R		LSUN-C		Textures		Places365		300K Rand. Img.		
	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	
$\pi = 0.05$													
OE	80.21±1.7	77.47±1.8	77.97±2.3	78.68±1.7	61.27±1.4	86.27±0.4	77.15±1.2	77.94±0.5	80.24±0.3	74.86±0.2	75.33±0.3	77.16±0.3	73.63±0.3
Energy (w/ OE)	77.47±2.0	80.48±1.2	70.83±3.1	82.86±2.0	29.42±4.3	94.61±0.8	72.05±0.8	80.73±0.5	74.69±0.6	78.60±0.4	66.91±0.7	80.44±0.5	75.77±0.1
WOODS (ours)	74.54 ±1.7	82.01 ±1.3	66.29 ±3.9	84.46 ±2.3	19.07 ±1.6	96.48 ±0.3	65.75 ±0.6	83.71 ±0.2	69.97 ±1.1	80.82 ±0.5	62.48 ±1.1	82.92 ±0.5	75.92 ±0.1
$\pi = 0.1$													
OE	79.56±1.6	77.00±1.2	76.86±2.1	78.75±1.2	58.53±2.8	86.92±0.8	74.63±1.2	79.13±0.5	78.52±0.3	75.68±0.1	72.18±0.2	78.48±0.3	73.53±0.4
Energy (w/ OE)	77.45±2.1	80.94±1.4	67.13±3.6	83.68±2.4	27.08±2.1	94.97±0.4	70.15±1.0	81.59±0.6	71.71±1.1	79.89±0.6	64.24±2.3	82.28±1.1	75.27±0.2
WOODS (ours)	71.67 ±1.9	84.11 ±1.4	59.27 ±3.9	86.80 ±1.9	15.03 ±1.4	97.24 ±0.3	61.38 ±0.7	85.57 ±0.2	64.19 ±1.0	83.12 ±0.5	55.51 ±1.3	85.72 ±0.4	75.64 ±0.3
$\pi = 0.2$													
OE	72.59±3.9	81.38±1.9	65.04±3.8	82.65±1.8	48.62±3.1	89.52±0.8	65.95±1.2	82.43±0.3	71.29±0.7	78.71±0.4	65.40±0.8	81.99±0.1	72.89±0.3
Energy (w/ OE)	72.76±2.5	83.48±1.2	62.53±5.7	84.46±2.8	22.49±1.2	95.84±0.2	64.93±0.5	83.87±0.4	64.62±0.2	82.72±0.2	56.07±1.2	85.50±0.2	75.00±0.3
WOODS (ours)	71.61 ±2.3	84.99 ±1.2	51.66 ±2.8	89.68 ±1.2	12.63 ±0.6	97.67 ±0.1	59.77 ±0.5	86.74 ±0.1	58.29 ±0.4	85.22 ±0.1	49.87 ±1.8	88.25 ±0.2	75.26 ±0.2
$\pi = 0.5$													
OE	68.80 ±2.8	82.89±1.1	47.64±4.7	88.84±1.8	30.86±1.9	93.91±0.4	56.18 ±1.6	86.11±0.4	62.24±0.5	82.53±0.3	53.70±1.6	86.58±0.2	73.00±0.3
Energy (w/ OE)	69.81±2.4	85.59±1.0	56.11±3.1	87.41±1.5	16.23±0.6	97.02±0.1	58.41±0.9	86.70±0.1	58.31±0.5	85.36±0.4	48.12±1.3	88.76±0.3	74.87±0.4
WOODS (ours)	69.41±2.7	86.76 ±0.8	44.60 ±2.6	91.72 ±0.7	12.70 ±0.4	97.71 ±0.1	57.60±0.6	87.74 ±0.1	55.03 ±0.3	86.82 ±0.1	45.00 ±0.7	89.85 ±0.2	75.72 ±0.0
$\pi = 1.0$													
OE	46.45 ±2.7	91.82 ±0.5	51.26±3.6	88.47±1.2	20.08±0.7	96.42±0.1	51.31 ±0.8	88.81±0.2	55.66±0.4	87.28±0.1	44.29±0.6	90.44±0.1	74.99±0.1
Energy (w/ OE)	56.40±4.0	89.48±1.2	54.41±2.5	88.77±0.8	17.14±0.9	96.91±0.1	52.36±1.3	89.38 ±0.3	54.11 ±0.9	88.35 ±0.2	43.42 ±1.0	90.88 ±0.1	74.85±0.2
WOODS (ours)	62.13±4.4	88.89±1.4	45.87 ±1.1	91.64 ±0.2	13.48 ±1.1	97.58 ±0.2	56.83±0.7	88.19±0.3	54.57±0.3	87.43±0.3	45.61±3.0	89.78±1.0	75.60 ±0.2

Table 2. Effect of π . ID dataset is CIFAR-100, and the auxiliary outlier training data is 300K Random Images. \uparrow indicates larger values are better and vice versa. $\pm x$ denotes the standard error, rounded to the first decimal point.

The initial learning rate is set to be 0.001 and decayed by a factor of 2 after 50%, 75%, and 90% of the epochs. We use a batch size of 128 and a dropout rate of 0.3. All training is performed in PyTorch using NVIDIA GeForce RTX 2080 Ti GPUs. For optimization in WOODS, we vary the penalty multiplier $\gamma \in \{1.1, 1.5\}$ and the dual update learning rate $\mu_2 \in \{0.1, 1, 2\}$. We set $\text{tol} = 0.05$, $\alpha = 0.05$, and set τ to be twice the loss of the pre-trained model. For validation, we use subsets of the ID data and of the \mathbb{P}_{out} data. Further details are included in Appendix A.3 and code is made publicly available online.

5.2. Results

WOODS Achieves Superior Performance We present results in Table 1, where WOODS outperforms the state-of-the-art results in the stationary wild setting $\mathbb{P}_{out}^{test} = \mathbb{P}_{out}$. This setting reflects the practical scenario when the OOD

data remains stable. Our comparison covers an extensive collection of competitive OOD detection methods. For clarity, we divide the baseline methods into two categories: trained with and without in-the-wild data. For methods using ID data \mathbb{P}_{in} only, we compare with methods such as MSP (Hendrycks & Gimpel, 2017), ODIN (Liang et al., 2018), Mahalanobis (Lee et al., 2018b), and Energy (Liu et al., 2020), the model is trained with softmax CE loss, same as in Equation 8. GODIN (Hsu et al., 2020) is trained using a DeConf-C loss, which does not involve auxiliary data loss either. We also include the latest development based on self-supervised losses, namely CSI (Tack et al., 2020).¹

¹We do not include standard errors for the methods that only use the ID data. We note the large difference in performance is statistically significant. We also include the pre-trained model in the Github repository.

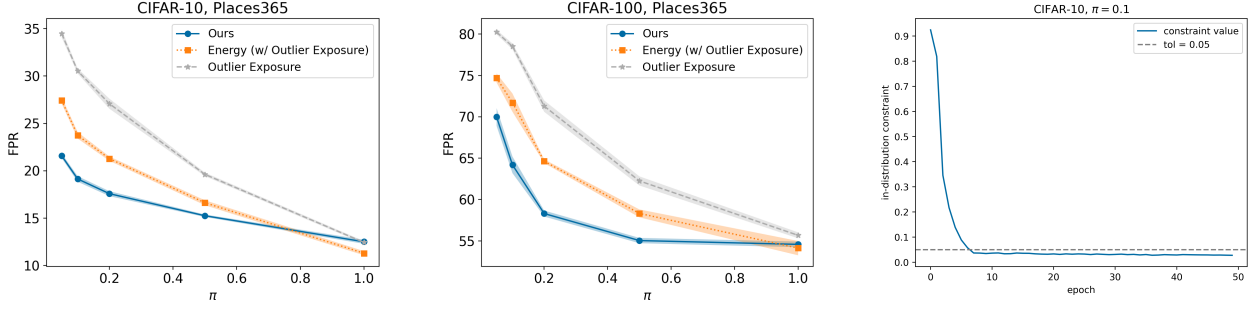


Figure 2: **Left and middle:** Ablation on π for the OOD setting, using CIFAR-10 (left) and CIFAR-100 (middle) as the ID dataset and Places365 as the OOD dataset. Our method WOODS is more reliable as π decreases. **Right:** Value of the ID constraint term $\frac{1}{m} \sum_{j=1}^m \frac{1}{1+\exp(w \cdot E_{\theta}(\mathbf{x}_i))} - \alpha$ from (9) over different training epochs. Our method is effective in satisfying this constraint, reducing it to zero (within a tolerance of 0.05).

Closest to ours are Outlier Exposure (OE) (Hendrycks et al., 2019) and energy-based OOD learning method (Liu et al., 2020). These are among the strongest OOD detection baselines, which regularize the classification model by producing lower confidence or higher energy on the auxiliary outlier data. For a fair comparison, all the methods in this group are trained using the same ID and in-the-wild data, under the same mixture ratio $\pi = 0.1$.

We highlight several observations: (1) Methods using wild data \mathbb{P}_{wild} , in general, show strong OOD detection performance over the counterpart (without \mathbb{P}_{wild}). Compared to the strongest baseline Mahalanobis in the first group, our method outperforms by **48.10%** in FPR95, averaged across all test datasets. The performance gain precisely demonstrates the advantage of our setting, which incorporates in-the-wild data for effective OOD learning. (2) Compared to methods using \mathbb{P}_{wild} , WOODS outperforms OE by **7.36%** in FPR95, with particularly strong performance on the Textures and Places365 datasets. In particular, OE makes a strong distributional assumption that the auxiliary outlier data does not overlap with the ID data. This assumption is violated when using wild mixture data. Indeed, our framework does not make this assumption, and demonstrates superior performance. (3) Lastly, the ID accuracy of the model trained with our method is comparable to that using the CE loss alone. Due to space constraints, we provide results on CIFAR-10 in the Appendix A.1, where our method’s strong performance holds.

Effect of π In Table 4 and Figure 2, we ablate the effect under different π , which modulates the fraction of OOD data in the mixture distribution \mathbb{P}_{wild} . Here we consistently use the 300k Random Images as the auxiliary outlier training dataset. Recall our definition in Section 2, a smaller π indicates more ID data and less OOD data—this reflects the practical scenario that the majority of test data may remain

ID. Note that the table covers both settings of $\mathbb{P}_{\text{out}}^{\text{test}} = \mathbb{P}_{\text{out}}$ (i.e., the column of 300K Rand Img) and $\mathbb{P}_{\text{out}}^{\text{test}} \neq \mathbb{P}_{\text{out}}$ (other columns). We highlight a few interesting observations: (1) The OOD detection performance for all methods (including OE and energy regularized learning) generally degrades as with decreasing π . In particular, a smaller π translates into a harder learning problem, because \mathbb{P}_{in} and \mathbb{P}_{wild} become largely overlapping. For example, on the 300K Random Images dataset, the FPR95 of OE increases from 44.29% ($\pi = 1.0$) to 75.33% ($\pi = 0.05$). (3) Our method WOODS is overall more robust under small π settings than the baselines. In a challenging case with $\pi = 0.05$, our method outperforms OE by **12.85%** in FPR95 on 300K Random Images. This demonstrates the benefits of WOODS performing constrained optimization.

WOODS Satisfies the Constraints in Optimization We also perform a sanity check on whether WOODS satisfies the constraints of the optimization objective in (9). As shown in Figure 2 (right), the ID constraint value $\frac{1}{m} \sum_{j=1}^m \frac{1}{1+\exp(w \cdot E_{\theta}(\mathbf{x}_i))} - \alpha$ is reduced to zero, within a specified tolerance of 0.05. This indeed verifies the efficacy of our constrained optimization framework.

Alternative Loss We note that our framework could also be compatible with alternative forms of loss function \mathcal{L}_{ood} . For example, we also explored another variant of our method, WOODS-alt, which uses a one hidden layer neural network applied to the penultimate layer for the OOD detection task. Section A.2 for more details of the loss function. The performance is summarized in Table 1. While WOODS-alt has strong performance in the setting where $\mathbb{P}_{\text{out}} = \mathbb{P}_{\text{out}}^{\text{test}}$, we do not expect it to perform as well as WOODS in the setting studied previously where $\mathbb{P}_{\text{out}} \neq \mathbb{P}_{\text{out}}^{\text{test}}$. WOODS uses the energy score to build its classifier, which already has reasonable performance even without any additional auxiliary dataset. On the other hand, WOODS-alt

would do no better than random guessing without an auxiliary dataset. In this way, WOODS has a prior given by the energy score that we believe helps in the $\mathbb{P}_{\text{out}} \neq \mathbb{P}_{\text{out}}^{\text{test}}$ setting. Indeed, we observed this in some experiments.

6. Related Work

OOD Detection OOD detection is an essential topic for safely deploying machine learning models in the open world, attracting much recent interest.

1) A rich line of methods design scoring functions for detecting samples from outside training categories, such as OpenMax (Bendale & Boulton, 2015), Maximum Softmax Probability (Hendrycks & Gimpel, 2017), ODIN score (Liang et al., 2018), Mahalanobis distance-based score (Lee et al., 2018b), Energy score (Liu et al., 2020; Wang et al., 2021; Morteza & Li, 2022; Sun et al., 2021), gradient-based score (Huang et al., 2021), and non-parametric KNN score (Sun et al., 2022). The aforementioned methods utilize ID data only. We show that by leveraging wild data that exists naturally in the model’s habitats, one can in fact build a stronger OOD detector.

2) Another line of works address the OOD detection problem by training-time regularization (Lee et al., 2018a; Bevandić et al., 2018; Hendrycks et al., 2019; Malinin & Gales, 2018; Liu et al., 2020; Chen et al., 2021). For example, models are encouraged to give predictions with lower confidence (Hendrycks et al., 2019) or higher energies (Liu et al., 2020; Ming et al., 2022). These methods typically require access to auxiliary outlier dataset. To circumvent this, recent work also explored synthesizing virtual outliers (Du et al., 2022b;a). In this work, we instead explore a more realistic setting by training OOD detectors using wild mixtures data containing both ID and OOD data—which can be easily obtainable upon deploying a machine learning classifier. We formulate a novel constrained optimization problem and show how to solve it tractably with modern neural networks.

Anomaly Detection Anomaly detection has received much attention in recent years (e.g., (Ruff et al., 2018; Chalapathy et al., 2018; Ergen & Kozat, 2019; Perera & Patel, 2019; Song et al., 2017)). In anomaly detection, a dataset is drawn *i.i.d.* from \mathbb{P}_{in} and the goal is to identify whether new data points are anomalous in the sense that they are not realizations from \mathbb{P}_{in} . In semi-supervised anomaly detection, an additional clean OOD dataset drawn from \mathbb{P}_{out} is observed (e.g., (Ruff et al., 2020; Daniel et al., 2019; Hendrycks et al., 2019)). An important difference between anomaly detection and the OOD detection literature is that OOD detection additionally requires learning a classifier for the distribution $\mathbb{P}_{\mathcal{X}\mathcal{Y}}$. We refer the reader to (Ruff et al., 2021; Chalapathy & Chawla, 2019) for detailed surveys on anomaly detection.

A closely related paper to our work is (Blanchard et al.,

2010), which studies the setting where samples from \mathbb{P}_{in} and \mathbb{P}_{wild} are observed and the goal is to find a θ minimizing $\mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{out}}}(\mathbb{I}\{g_{\theta}(\mathbf{x}) = \text{in}\})$ subject to the constraint that $\mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{in}}}(\mathbb{I}\{g_{\theta}(\mathbf{x}) = \text{out}\}) \leq \alpha$. This work has several important differences with ours. First, they do not consider the out-of-distribution problem, that is, where the distribution at test time $\mathbb{P}_{\text{out}}^{\text{test}}$ differs from \mathbb{P}_{out} , whereas our energy-based approach (9) does. Second, their formulation only considers the task of distinguishing \mathbb{P}_{out} and \mathbb{P}_{in} , not the task of doing classification simultaneously. Third, our work uses and theoretically analyzes the sigmoid loss for OOD detection, a differentiable loss that can be used in deep learning, whereas their work focuses on the computationally intractable 0-1 loss. Finally, their work is mainly statistical, only implementing their algorithm using a plug-in kernel-density estimator whereas we leverage neural networks using a computational approach based on ALM.

Constrained Optimization The augmented Lagrangian method is a popular approach to constrained optimization. It improves over two other related methods: the penalty method and the method of Lagrangian multipliers. While the penalty method suffers from ill-conditioning (Nocedal & Wright, 2006), the method of Lagrangian multipliers is specific to the convex case (Rockafellar, 1973). In this paper, we adapt ALM to our setting from a recent version proposed and analyzed for the case of nonlinear inequality constraints (Xu, 2021a). There are only a limited number of examples of adapting ALM to modern neural networks (e.g., (Sangalli et al., 2021) for class imbalance). To the best of our knowledge, our work is the first to explore constrained optimization for OOD detection problem

7. Conclusion

In this paper, we propose a novel framework for OOD detection using wild data. Wild data has significant promise since (i) it is abundant, (ii) can be collected essentially for free upon deploying an ML system, and (iii) is often a much better match to the test-time distribution than data collected offline. At the same time, it is challenging to leverage because it naturally consists of both ID and OOD examples. To overcome this challenge, we propose a framework based on constrained optimization and solve it tractably by adapting the augmented Lagrangian method to deep neural networks. We believe that wild data has the potential to dramatically advance OOD detection in practice, thereby helping to accelerate the deployment of safe and reliable machine learning.

Acknowledgments

We would like to acknowledge Xuefeng Du’s help in verifying the baseline performance. We would also like to thank Ahmet Alacaoglu for very useful conversations.

References

- Bendale, A. and Boulton, T. Towards open world recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1893–1902, 2015.
- Bevandić, P., Krešo, I., Oršić, M., and Šegvić, S. Discriminative out-of-distribution detection for semantic segmentation. *arXiv preprint arXiv:1808.07703*, 2018.
- Blanchard, G., Lee, G., and Scott, C. Semi-supervised novelty detection. *The Journal of Machine Learning Research*, 11:2973–3009, 2010.
- Chalapathy, R. and Chawla, S. Deep learning for anomaly detection: A survey. *arXiv preprint arXiv:1901.03407*, 2019.
- Chalapathy, R., Menon, A. K., and Chawla, S. Anomaly detection using one-class neural networks. *arXiv preprint arXiv:1802.06360*, 2018.
- Chen, J., Li, Y., Wu, X., Liang, Y., and Jha, S. Atom: Robustifying out-of-distribution detection using outlier mining. In *Proceedings of European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PKDD)*, 2021.
- Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., and Vedaldi, A. Describing Textures in the Wild. In *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- Daniel, T., Kurutach, T., and Tamar, A. Deep variational semi-supervised novelty detection. *arXiv preprint arXiv:1911.04971*, 2019.
- Du, X., Wang, X., Gozum, G., and Li, Y. Unknown-aware object detection: Learning what you don’t know from videos in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022a.
- Du, X., Wang, Z., Cai, M., and Li, Y. Vos: Learning what you don’t know by virtual outlier synthesis. *Proceedings of the International Conference on Learning Representations*, 2022b.
- Duchi, J., Hazan, E., and Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- Ergen, T. and Kozat, S. S. Unsupervised anomaly detection with lstm neural networks. *IEEE transactions on neural networks and learning systems*, 31(8):3127–3141, 2019.
- Hendrycks, D. and Gimpel, K. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *Proceedings of International Conference on Learning Representations*, 2017.
- Hendrycks, D., Mazeika, M., and Dietterich, T. Deep anomaly detection with outlier exposure. In *International Conference on Learning Representations*, 2019.
- Hestenes, M. R. Multiplier and gradient methods. *Journal of optimization theory and applications*, 4(5):303–320, 1969.
- Hsu, Y.-C., Shen, Y., Jin, H., and Kira, Z. Generalized odin: Detecting out-of-distribution image without learning from out-of-distribution data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10951–10960, 2020.
- Huang, R., Geng, A., and Li, Y. On the importance of gradients for detecting distributional shifts in the wild. In *Advances in Neural Information Processing Systems*, 2021.
- Huber, P. J. Robust estimation of a location parameter. *Annals of Mathematical Statistics*, 35:73–101, March 1964.
- Krizhevsky, A., Hinton, G., and others. Learning multiple layers of features from tiny images. 2009. Publisher: Citeseer.
- Lee, K., Lee, H., Lee, K., and Shin, J. Training confidence-calibrated classifiers for detecting out-of-distribution samples. *International Conference on Learning Representations (ICLR)*, 2018a.
- Lee, K., Lee, K., Lee, H., and Shin, J. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In *Advances in Neural Information Processing Systems*, pp. 7167–7177, 2018b.
- Liang, S., Li, Y., and Srikant, R. Enhancing the reliability of out-of-distribution image detection in neural networks. In *6th International Conference on Learning Representations, ICLR 2018*, 2018.
- Liu, W., Wang, X., Owens, J., and Li, Y. Energy-based out-of-distribution detection. *Advances in Neural Information Processing Systems*, 2020.
- Malinin, A. and Gales, M. Predictive uncertainty estimation via prior networks. *arXiv preprint arXiv:1802.10501*, 2018.
- Ming, Y., Fan, Y., and Li, Y. Poem: Out-of-distribution detection with posterior sampling. In *International Conference on Machine Learning (ICML)*. PMLR, 2022.
- Morteza, P. and Li, Y. Provable guarantees for understanding out-of-distribution detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022.

- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. Reading Digits in Natural Images with Unsupervised Feature Learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011.
- Nguyen, A., Yosinski, J., and Clune, J. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 427–436, 2015.
- Nocedal, J. and Wright, S. *Numerical optimization*. Springer Science & Business Media, 2006.
- Perera, P. and Patel, V. M. Learning deep features for one-class classification. *IEEE Transactions on Image Processing*, 28(11):5450–5463, 2019.
- Rockafellar, R. T. A dual approach to solving nonlinear programming problems by unconstrained optimization. *Mathematical programming*, 5(1):354–373, 1973.
- Ruff, L., Vandermeulen, R., Goernitz, N., Deecke, L., Siddiqui, S. A., Binder, A., Müller, E., and Kloft, M. Deep one-class classification. In *International conference on machine learning*, pp. 4393–4402. PMLR, 2018.
- Ruff, L., Vandermeulen, R. A., Görnitz, N., Binder, A., Müller, E., Müller, K.-R., and Kloft, M. Deep semi-supervised anomaly detection. In *International Conference on Learning Representations*, 2020.
- Ruff, L., Kauffmann, J. R., Vandermeulen, R. A., Montavon, G., Samek, W., Kloft, M., Dietterich, T. G., and Müller, K.-R. A unifying review of deep and shallow anomaly detection. *Proceedings of the IEEE*, 2021.
- Sangalli, S., Erdil, E., Hötker, A., Donati, O., and Konukoglu, E. Constrained optimization to train neural networks on critical and under-represented classes. *Advances in Neural Information Processing Systems*, 34, 2021.
- Song, H., Jiang, Z., Men, A., and Yang, B. A hybrid semi-supervised anomaly detection model for high-dimensional data. *Computational intelligence and neuroscience*, 2017, 2017.
- Sun, Y., Guo, C., and Li, Y. React: Out-of-distribution detection with rectified activations. In *Advances in Neural Information Processing Systems*, 2021.
- Sun, Y., Ming, Y., Zhu, X., and Li, Y. Out-of-distribution detection with deep nearest neighbors. In *International Conference on Machine Learning (ICML)*. PMLR, 2022.
- Tack, J., Mo, S., Jeong, J., and Shin, J. Csi: Novelty detection via contrastive learning on distributionally shifted instances. In *Advances in Neural Information Processing Systems*, 2020.
- Wang, H., Liu, W., Bocchieri, A., and Li, Y. Can multi-label classification networks know what they don’t know? *Advances in Neural Information Processing Systems*, 34, 2021.
- Xu, Y. First-order methods for constrained convex programming based on linearized augmented lagrangian function. *Inform Journal on Optimization*, 3(1):89–117, 2021a.
- Xu, Y. Iteration complexity of inexact augmented lagrangian methods for constrained convex programming. *Mathematical Programming*, 185(1):199–244, 2021b.
- Yan, Y. and Xu, Y. Adaptive primal-dual stochastic gradient method for expectation-constrained convex stochastic programs. *arXiv preprint arXiv:2012.14943*, 2020.
- Yu, F., Seff, A., Zhang, Y., Song, S., Funkhouser, T., and Xiao, J. LSUN: Construction of a Large-scale Image Dataset using Deep Learning with Humans in the Loop. *arXiv:1506.03365 [cs]*, June 2016. arXiv: 1506.03365.
- Zagoruyko, S. and Komodakis, N. Wide Residual Networks. In *Proceedings of the British Machine Vision Conference*, 2016.
- Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.
- Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., and Torralba, A. Places: A 10 Million Image Database for Scene Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(6):1452–1464, June 2018.

A. Additional Experimental Results

A.1. Main Experiments: CIFAR-10

Table 3 shows a comparison of our method’s performance with OOD baseline methods on CIFAR-10. On average, WOODS outperforms CSI, the best of the baseline methods, by **13.59%**, and OE by **3.61%**. Since all of the methods that use \mathbb{P}_{wild} achieve very good performance on most of the auxiliary test datasets, we highlight in particular the results for Textures (where our method outperforms CSI by **14.50%** and OE by **4.47%**) and Places365 (where our method outperforms CSI by **22.46%** and OE by **10.98%**). Table 4 shows the results of ablation on π using CIFAR-10 as the ID dataset, where WOODS demonstrates similarly strong performance across the board.

Table 4 shows the results of ablation on π using CIFAR-10 as the ID dataset. In general, WOODS outperforms OE and Energy by an even larger margin than on CIFAR-10 and across many values of π .

Method	SVHN		LSUN-R		OOD Dataset LSUN-C		Textures		Places365		Average		Acc.
	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	
With \mathbb{P}_{in} only													
MSP	48.49	91.89	52.15	91.37	30.80	95.65	59.28	88.50	59.48	88.20	50.04	91.12	94.84
ODIN	33.35	91.96	26.62	94.57	15.52	97.04	49.12	84.97	57.40	84.49	36.40	90.61	94.84
Energy	35.59	90.96	27.58	94.24	8.26	98.35	52.79	85.22	40.14	89.89	32.87	91.73	94.84
Mahalanobis	12.89	97.62	42.62	93.23	39.22	94.15	15.00	97.33	68.57	84.61	35.66	93.34	94.84
GODIN	13.55	97.61	17.93	96.86	17.68	96.93	29.43	94.87	41.27	91.49	23.97	95.55	94.48
CSI	17.30	97.40	12.15	98.01	1.95	99.55	20.45	95.93	34.95	93.64	17.36	96.91	94.17
With \mathbb{P}_{in} and \mathbb{P}_{wild}													
OE	0.85 \pm 0.1	99.82 \pm 0.0	0.33 \pm 0.0	99.93 \pm 0.0	1.84 \pm 0.2	99.65 \pm 0.0	10.42 \pm 0.6	98.01 \pm 0.0	23.47 \pm 0.5	94.62 \pm 0.0	7.38 \pm 0.3	98.41 \pm 0.0	94.07 \pm 0.2
Energy (w/ OE)	4.95 \pm 4.5	98.92 \pm 1.0	5.04 \pm 4.9	98.83 \pm 1.1	1.93 \pm 1.4	99.49 \pm 0.3	13.43 \pm 6.5	96.69 \pm 1.8	17.26 \pm 0.4	95.84 \pm 0.1	8.52 \pm 3.5	97.95 \pm 0.9	94.81 \pm 0.1
WOODS (ours)	0.15 \pm 0.0	99.97 \pm 0.0	0.03 \pm 0.0	99.99 \pm 0.0	0.22 \pm 0.0	99.94 \pm 0.0	5.95 \pm 0.6	98.79 \pm 0.1	12.49 \pm 0.3	97.00 \pm 0.0	3.77 \pm 0.2	99.14 \pm 0.0	94.84 \pm 0.1
WOODS-alt (ours)	0.10 \pm 0.0	99.96 \pm 0.0	0.02 \pm 0.0	99.99 \pm 0.0	0.08 \pm 0.0	99.96 \pm 0.0	2.11 \pm 0.2	99.32 \pm 0.0	11.51 \pm 0.2	96.25 \pm 0.2	2.76 \pm 0.0	99.10 \pm 0.0	94.71 \pm 0.1

Table 3. **Main results when $\mathbb{P}_{\text{out}}^{\text{test}} = \mathbb{P}_{\text{out}}$.** Comparison with competitive OOD detection methods on CIFAR-10. For methods using \mathbb{P}_{wild} , we train under the same dataset and same $\pi = 0.1$. \uparrow indicates larger values are better and vice versa. $\pm x$ denotes the standard error, rounded to the first decimal point.

Method	OOD Dataset												Acc.
	SVHN		LSUN-R		LSUN-C		Textures		Places365		300K Rand. Img.		
	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	
$\pi = 0.05$													
OE	17.23 \pm 1.7	96.66 \pm 0.3	18.31 \pm 1.0	96.80 \pm 0.2	7.38 \pm 0.4	98.59 \pm 0.1	29.47 \pm 1.1	93.63 \pm 0.3	34.47 \pm 0.4	92.44 \pm 0.2	33.69 \pm 0.3	91.89 \pm 0.2	94.38 \pm 0.1
Energy (w/ OE)	14.19 \pm 2.0	97.04 \pm 0.5	12.54 \pm 0.1	97.47 \pm 0.0	3.76 \pm 0.2	99.16 \pm 0.0	29.44 \pm 1.6	92.73 \pm 0.5	27.42 \pm 0.3	93.21 \pm 0.1	29.34 \pm 1.2	92.23 \pm 0.5	94.71 \pm 0.1
WOODS (ours)	9.23 \pm 1.0	97.98 \pm 0.2	6.89 \pm 0.6	98.53 \pm 0.1	1.82 \pm 0.1	99.56 \pm 0.0	23.12 \pm 1.5	94.98 \pm 0.3	21.58 \pm 0.2	95.07 \pm 0.1	24.39 \pm 0.9	94.16 \pm 0.4	94.83 \pm 0.1
$\pi = 0.1$													
OE	13.18 \pm 1.3	97.34 \pm 0.2	14.32 \pm 0.4	97.44 \pm 0.1	5.87 \pm 0.2	98.86 \pm 0.0	25.69 \pm 0.6	94.35 \pm 0.1	30.54 \pm 0.3	93.31 \pm 0.1	30.69 \pm 0.7	92.80 \pm 0.2	94.21 \pm 0.1
Energy (w/ OE)	8.52 \pm 1.8	98.13 \pm 0.3	9.05 \pm 0.7	98.13 \pm 0.1	2.78 \pm 0.2	99.38 \pm 0.0	22.32 \pm 1.3	94.72 \pm 0.4	23.74 \pm 0.3	94.26 \pm 0.1	24.59 \pm 0.8	93.99 \pm 0.2	94.54 \pm 0.1
WOODS (ours)	5.70 \pm 1.0	98.54 \pm 0.2	4.13 \pm 0.4	99.01 \pm 0.1	1.31 \pm 0.1	99.66 \pm 0.0	17.92 \pm 1.3	96.43 \pm 0.2	19.14 \pm 0.3	95.74 \pm 0.0	19.82 \pm 0.5	95.52 \pm 0.1	94.74 \pm 0.0
$\pi = 0.2$													
OE	9.94 \pm 1.0	98.04 \pm 0.1	10.37 \pm 0.5	98.10 \pm 0.1	4.82 \pm 0.3	99.08 \pm 0.0	20.34 \pm 0.7	95.67 \pm 0.2	27.08 \pm 0.5	93.98 \pm 0.2	24.94 \pm 0.8	94.25 \pm 0.1	94.10 \pm 0.0
Energy (w/ OE)	5.88 \pm 1.1	98.50 \pm 0.2	6.18 \pm 0.7	98.65 \pm 0.1	2.38 \pm 0.1	99.45 \pm 0.0	17.28 \pm 0.8	96.14 \pm 0.3	21.28 \pm 0.2	94.93 \pm 0.1	19.54 \pm 0.6	94.99 \pm 0.1	94.30 \pm 0.1
WOODS (ours)	4.20 \pm 0.8	98.68 \pm 0.2	3.35 \pm 0.5	99.19 \pm 0.1	1.31 \pm 0.1	99.64 \pm 0.0	13.17 \pm 0.7	97.45 \pm 0.1	17.58 \pm 0.3	96.12 \pm 0.0	15.98 \pm 0.4	96.29 \pm 0.1	94.75 \pm 0.1
$\pi = 0.5$													
OE	5.07 \pm 1.1	98.75 \pm 0.1	5.23 \pm 0.8	98.88 \pm 0.1	2.94 \pm 0.2	99.36 \pm 0.0	12.06 \pm 0.7	97.45 \pm 0.2	19.62 \pm 0.2	95.47 \pm 0.1	16.32 \pm 0.6	95.90 \pm 0.0	94.16 \pm 0.1
Energy (w/ OE)	2.57 \pm 0.4	99.16 \pm 0.1	4.64 \pm 0.7	98.94 \pm 0.1	1.57 \pm 0.1	99.60 \pm 0.0	11.77 \pm 0.9	97.50 \pm 0.2	16.63 \pm 0.3	96.04 \pm 0.0	13.86 \pm 0.1	96.34 \pm 0.1	94.57 \pm 0.1
WOODS (ours)	3.58 \pm 1.4	98.69 \pm 0.2	2.78 \pm 0.2	99.21 \pm 0.0	1.44 \pm 0.1	99.56 \pm 0.0	9.20 \pm 0.3	98.09 \pm 0.1	15.25 \pm 0.1	96.47 \pm 0.0	11.82 \pm 0.1	96.84 \pm 0.1	94.72 \pm 0.0
$\pi = 1.0$													
OE	1.72 \pm 0.1	99.21 \pm 0.1	2.81 \pm 0.2	99.20 \pm 0.0	1.72 \pm 0.1	99.50 \pm 0.0	6.70 \pm 0.3	98.51 \pm 0.0	12.43 \pm 0.1	96.98 \pm 0.0	10.43 \pm 0.4	97.32 \pm 0.0	94.62 \pm 0.1
Energy (w/ OE)	8.20 \pm 2.4	97.85 \pm 0.4	2.72 \pm 0.2	99.05 \pm 0.1	1.62 \pm 0.1	99.36 \pm 0.0	5.77 \pm 0.2	98.50 \pm 0.1	11.29 \pm 0.2	97.09 \pm 0.0	9.82 \pm 0.1	97.33 \pm 0.0	94.23 \pm 0.1
WOODS (ours)	2.35 \pm 0.5	98.42 \pm 0.2	2.27 \pm 0.2	99.14 \pm 0.1	1.34 \pm 0.0	99.50 \pm 0.0	6.31 \pm 0.4	98.42 \pm 0.1	12.50 \pm 0.2	96.78 \pm 0.0	10.01 \pm 0.3	97.11 \pm 0.1	94.85 \pm 0.0

Table 4. **Effect of π .** A larger π indicates more OOD data in the mixture distribution \mathbb{P}_{wild} . ID dataset is CIFAR-10, and the auxiliary outlier training data is 300K Random Images. \uparrow indicates larger values are better and vice versa. $\pm x$ denotes the standard error, rounded to the first decimal point.

A.2. WOODS-alt Description

Here, we describe WOODS-alt, for which an OOD confidence score is not extracted directly from the output of the ID classifier, but rather learned by a separate neural network attached to the ID classifier’s penultimate layer. The additional neural network has one fully-connected hidden layer with 300 neurons, followed by a ReLU activation and a single output

logit, which provides an OOD confidence score, denoted $g_\theta(\cdot)$. With this architecture, we apply the same ALM algorithm to solve the following constrained optimization problem:

$$\begin{aligned} & \inf_{\theta} \frac{1}{m} \sum_{i=1}^m \max(1 - g_\theta(\tilde{\mathbf{x}}_i), 0) \\ & \text{s.t. } \frac{1}{n} \sum_{j=1}^n \max(1 + g_\theta(\mathbf{x}_j), 0) \leq \alpha \\ & \quad \frac{1}{n} \sum_{j=1}^n \mathcal{L}_{\text{cls}}(f_\theta(\mathbf{x}_j), y_j) \leq \tau. \end{aligned}$$

A.3. Additional Experimental Details

Here we give additional experimental details, presented in Sections 5 and A.1. Energy and OE both optimize an objective of the form

$$\min_{\theta} \mathcal{L}_{\text{cls}} + \lambda \mathcal{L}_{\text{OOD}}$$

For energy, we varied $\lambda \in \{0.1, 1, 5\}$ and for OE we varied $\lambda \in \{0.1, 0.5, 1\}$.

We simulate the mixture distribution as follows. For each iteration at training, for the ID dataset we draw one batch of size 128 and for the wild dataset \mathbb{P}_{wild} we draw another batch of size 128 where each example is drawn from \mathbb{P}_{out} with probability π and from \mathbb{P}_{in} with probability $1 - \pi$.

We split the data as follows. For the stationary setting $\mathbb{P}_{\text{out}}^{\text{test}} = \mathbb{P}_{\text{out}}$ (say from SVHN), we use 70% of the SVHN for the mixture training dataset and for the validation dataset. We use the remaining examples for the test set. Of the data for training/validation, we use 30% for validation and the remaining for training. In the nonstationary setting $\mathbb{P}_{\text{out}}^{\text{test}} \neq \mathbb{P}_{\text{out}}$, we use the same splitting approach except in the initial split into train/validation data and test, we use 99% of the 300K Random Images dataset (due to its large size). For the ID data, we use 50% for in-distribution training samples and 50% for generating the mixture.

We repeat each experiment 5 times with 5 separate seeds.

B. Proof of Proposition 3.1

In this Section, we prove Proposition 3.1. We begin by giving notation and proving an important Lemma. To ease notation, we write $\mathbb{P}_{\text{out}} =: \mathbb{P}_1$ and $\mathbb{P}_{\text{in}} =: \mathbb{P}_0$. Define the sigmoid loss for the OOD task

$$R_y(g_\theta) := \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{out}}}(\sigma(-g_\theta(\mathbf{x})) \cdot y).$$

Define

$$\begin{aligned} R_1^* &:= \inf_{\theta} R_1(g_\theta) \\ &\text{s.t. } R_0(g_\theta) \leq \alpha. \end{aligned}$$

Define

$$\begin{aligned} R_{\text{wild}} &:= \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{wild}}}(\sigma(-g_\theta(\mathbf{x}))) \\ &= \pi R_1(g_\theta) + (1 - \pi) \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_0}(\sigma(-g_\theta(\mathbf{x}))) \\ &= \pi R_1(g_\theta) + (1 - \pi)(1 - R_0(g_\theta)) \end{aligned}$$

where we used the symmetry of the sigmoid function, that is, $\sigma(z) + \sigma(-z) = 1$ for $z \in \mathbb{R}$. Now, define

$$\begin{aligned} R_{\text{wild}}^* &:= \inf_{\theta} R_{\text{wild}}(g_\theta) \\ &\text{s.t. } R_0(g_\theta) \leq \alpha. \end{aligned}$$

Next, we prove a key Lemma for our proof. This Lemma has a similar proof to Theorem 1 in (Blanchard et al., 2010), which applies to the 0/1 loss. We establish an analogous result for the sigmoid loss. The key observation is that the symmetry property of the sigmoid loss enables a similar proof.

Lemma B.1. *Suppose that there exists θ^* such that $R_1(g_\theta) = R_1^*(g_{\theta^*})$ and $R_0(g_{\theta^*}) = \alpha$. Then,*

$$R_1(g_\theta) - R_1^* \leq \frac{1}{\pi}(R_{\text{wild}}(g_\theta) - R_1^* + (1 - \pi)(R_0(g_\theta) - \alpha)).$$

Proof. We begin by showing that for any θ , $R_{\text{wild}}(g_\theta) = R_{\text{wild}}^*$ and $R_0(g_\theta) \leq \alpha$ if and only if $R_1(g_\theta) = R_1^*$ and $R_0(g_\theta) = \alpha$.

\implies : First, suppose θ satisfies $R_{\text{wild}}(g_\theta) = R_{\text{wild}}^*$ and $R_0(g_\theta) \leq \alpha$. Suppose that either $R_0(g_\theta) < \alpha$ or $R_1(g_\theta) > R_1^*$. By the assumption in the Proposition, there exists θ^* such that $R_1(g_\theta) = R_1^*(g_{\theta^*})$ and $R_0(g_{\theta^*}) = \alpha$. Then, we have that

$$\begin{aligned} R_{\text{wild}}^*(g_{\theta^*}) &= \pi R_1(g_{\theta^*}) + (1 - \pi)(1 - R_0(g_{\theta^*})) \\ &= \pi R_1(g_\theta) + (1 - \pi)(1 - \alpha) \\ &< \pi R_1(g_\theta) + (1 - \pi)(1 - R_0(g_\theta)) \\ &= R_{\text{wild}}(g_\theta). \end{aligned}$$

But, this contradicts the assumption that $R_{\text{wild}}(g_\theta) = R_{\text{wild}}^*$, completing this direction of the claim.

\impliedby : Suppose θ satisfies $R_1(g_\theta) = R_1^*$ and $R_0(g_\theta) = \alpha$. By the assumption in the Lemma, there exists θ^* such that $R_1(g_\theta) = R_1^*(g_{\theta^*})$ and $R_0(g_{\theta^*}) = \alpha$. Towards a contradiction, suppose that $R_{\text{wild}}(g_{\theta^*}) < R_{\text{wild}}(g_\theta)$. Then, using $\pi > 0$, we have that

$$\begin{aligned} R_1(g_{\theta^*}) &= \frac{1}{\pi}(R_{\text{wild}}(g_{\theta^*}) - (1 - \pi)(1 - R_0(g_{\theta^*}))) \\ &< \frac{1}{\pi}(R_{\text{wild}}(g_\theta) - (1 - \pi)(1 - \alpha)) \\ &= R_1(g_\theta) \end{aligned}$$

but this contradicts our assumption on g_θ . This completes the proof of the claim.

The established claim implies that $R_{\text{wild}}^* = \pi R_1^* + (1 - \pi)(1 - \alpha)$. The result now follows by subtracting this equality from $R_{\text{wild}}^*(g_\theta) = \pi R_1(g_\theta) + (1 - \pi)(1 - R_0(g_\theta))$.

□

Here we restate Proposition 3.1 with all technical details included. Let $\Theta \subset \mathbb{R}^p$ where we have that $\theta \in \Theta$. Recall that for the purposes of this proposition we assume that $f_\theta(x) : \mathbb{R}^d \mapsto [0, 1]$ and define $\mathcal{L}_{\text{cls}}(s, y) = y \log(\frac{1}{s}) + (1 - y) \log(\frac{1}{1-s})$ for $s \in [0, 1]$. Define

$$\begin{aligned} \hat{\theta}_\epsilon &\leftarrow \operatorname{argmin}_{\theta \in \Theta} \frac{1}{m} \sum_{i=1}^m \sigma(-g_\theta(\tilde{\mathbf{x}}_i)) \\ \text{s.t. } &\frac{1}{n} \sum_{j=1}^n \sigma(g_\theta(\mathbf{x}_j)) \leq \alpha + \epsilon \\ &\frac{1}{n} \sum_{j=1}^n \mathcal{L}_{\text{cls}}(f_\theta(\mathbf{x}_j), y_j) \leq \tau + \epsilon \end{aligned} \tag{10}$$

Recall the optimization problem of interest:

$$\begin{aligned} \inf_{\theta \in \Theta} & \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{out}}} \sigma(-g_\theta(\tilde{\mathbf{x}}_i)) \\ \text{s.t. } & \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{in}}} \sigma(g_\theta(\mathbf{x}_j)) \leq \alpha \\ & \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_{\mathcal{X} \times \mathcal{Y}}} [\mathcal{L}_{\text{cls}}(f_\theta(\mathbf{x}), y)] \leq \tau. \end{aligned} \tag{11}$$

and let opt denote its value.

We will make the following mild assumption and describe settings where it holds later.

Assumption B.2. There exists θ^* such that $R_1(g_\theta) = R_1^*(g_{\theta^*})$, $R_0(g_{\theta^*}) = \alpha$, and $\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_{\mathcal{X}\mathcal{Y}}}[\mathcal{L}_{\text{cls}}(f_{\theta^*}(\mathbf{x}), y)] \leq \tau$.

Proposition B.3. Suppose $K = 2$. Suppose Assumption (B.2) holds. Define $\epsilon_k := \sqrt{\frac{2 \ln(6/\delta)}{k}} + 2 \max_{h \in \{f, g\}} \max_{\mathbb{P} \in \{\mathbb{P}_1, \mathbb{P}_{\text{wild}}\}} \mathbb{E}_{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m \sim \mathbb{P}} \mathbb{E}_{\eta_1, \dots, \eta_m} \sup_{\theta \in \Theta} \frac{1}{k} \sum_{i=1}^k \eta_i h_\theta(\tilde{\mathbf{x}}_i)$ where η_1, \dots, η_k are i.i.d. and $\mathbb{P}(\eta_i = 1) = \mathbb{P}(\eta_i = -1) = 1/2$. Let $\hat{\theta}_\epsilon$ solve (10) with tolerance $\epsilon = c\epsilon_n$ where c is a universal positive constant. Then, with probability at least $1 - \delta$

1. $\mathbb{E}_{\text{out}} \sigma(-g_{\hat{\theta}_\epsilon}(\mathbf{x})) \leq \text{opt} + c_1 \pi^{-1}(\epsilon_n + \epsilon_m)$,
2. $\mathbb{E}_{\text{in}} \sigma(g_{\hat{\theta}_\epsilon}(\mathbf{x})) \leq \alpha + c_2 \epsilon_n$, and
3. $\mathbb{E}_{\mathcal{X}\mathcal{Y}}[\mathcal{L}_{\text{cls}}(f_{\hat{\theta}_\epsilon}(\mathbf{x}), y)] \leq \tau + c_3 \epsilon_n$.

Proof. Define the following events

$$\begin{aligned} \Sigma_1 &= \left\{ \left| \frac{1}{m} \sum_{i=1}^m \sigma(-g_\theta(\tilde{\mathbf{x}}_i)) - R_{\text{wild}}(g_\theta) \right| \leq 2 \mathbb{E}_{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m \sim \mathbb{P}_{\text{wild}}} \mathbb{E}_{\eta_1, \dots, \eta_m} \sup_{\theta \in \Theta} \frac{1}{m} \sum_{i=1}^m \eta_i g_\theta(\tilde{\mathbf{x}}_i) + c \sqrt{\frac{2 \ln(6/\delta)}{m}} : \forall \theta \in \Theta \right\} \\ \Sigma_2 &= \left\{ \left| \frac{1}{n} \sum_{i=1}^n \sigma(g_\theta(\mathbf{x}_i)) - R_0(g_\theta) \right| \leq 2 \mathbb{E}_{\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathbb{P}_0} \mathbb{E}_{\eta_1, \dots, \eta_n} \sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \eta_i g_\theta(\mathbf{x}_i) + c \sqrt{\frac{2 \ln(6/\delta)}{n}} : \forall \theta \in \Theta \right\} \\ \Sigma_3 &= \left\{ \left| \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\text{cls}}(f_\theta(\mathbf{x}_i), y_i) - \mathbb{E}_{\mathcal{X}\mathcal{Y}}[\mathcal{L}_{\text{cls}}(f_\theta(\mathbf{x}), y)] \right| \leq 2 \mathbb{E}_{\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathbb{P}_0} \mathbb{E}_{\eta_1, \dots, \eta_n} \sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \eta_i f_\theta(\mathbf{x}_i) \right. \\ &\quad \left. + c \sqrt{\frac{2 \ln(6/\delta)}{n}} : \forall \theta \in \Theta \right\} \end{aligned}$$

By Lemma B.4, we have that $\mathbb{P}(\Sigma_i) \geq 1 - \delta/3$ for all $i = 1, 2, 3$. Then, by the union bound, we have that $\Sigma := \Sigma_1 \cap \Sigma_2 \cap \Sigma_3$ holds with probability at least $1 - \delta$. Assume Σ holds for the remainder of the proof.

Note that we have

$$\begin{aligned} R_0(g_{\hat{\theta}_\epsilon}) - R_0(g_{\theta^*}) &= R_0(g_{\hat{\theta}_\epsilon}) - \frac{1}{n} \sum_{i=1}^n \sigma(-g_{\hat{\theta}_\epsilon}(\mathbf{x}_i)) \\ &\quad + \frac{1}{n} \sum_{i=1}^n \sigma(-g_{\hat{\theta}_\epsilon}(\mathbf{x}_i)) - \frac{1}{n} \sum_{i=1}^n \sigma(-g_{\theta^*}(\mathbf{x}_i)) \\ &\quad + \frac{1}{n} \sum_{i=1}^n \sigma(-g_{\theta^*}(\mathbf{x}_i)) - R_0(g_{\theta^*}) \\ &\leq R_0(g_{\hat{\theta}_\epsilon}) - \frac{1}{n} \sum_{i=1}^n \sigma(-g_{\hat{\theta}_\epsilon}(\mathbf{x}_i)) \\ &\quad + \frac{1}{n} \sum_{i=1}^n \sigma(-g_{\theta^*}(\mathbf{x}_i)) - R_0(g_{\theta^*}) \end{aligned} \tag{12}$$

$$\leq 4 \mathbb{E}_{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m \sim \mathbb{P}_{\text{wild}}} \mathbb{E}_{\eta_1, \dots, \eta_m} \sup_{\theta \in \Theta} \frac{1}{m} \sum_{i=1}^m \sigma_i g_\theta(\tilde{\mathbf{x}}_i) + 2c \sqrt{\frac{2 \ln(2/\delta)}{m}}. \tag{13}$$

(12) follows since $\hat{\theta}_\epsilon$ is feasible for the optimization problem (10) because, by the choice of ϵ and Σ ,

$$\begin{aligned} R_0(g_{\theta^*}) &\leq \alpha \\ \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_0} \max(1 - f_{\theta^*}(\mathbf{x})y, 0) &\leq \tau \end{aligned}$$

Therefore, by definition of $\hat{\theta}_\epsilon$ as the minimizer of (10), we have that

$$\frac{1}{n} \sum_{i=1}^n \sigma(-g_{\hat{\theta}_\epsilon}(\mathbf{x}_i)) - \frac{1}{n} \sum_{i=1}^n \sigma(-g_{\theta^*}(\mathbf{x}_i)) \leq 0.$$

(13) follows by the event Σ_1 .

Similarly,

$$\begin{aligned} \mathbb{E}_{\text{in}} \sigma(g_{\hat{\theta}_\epsilon}(\mathbf{x})) &= \frac{1}{n} \sum_{i=1}^n \sigma(g_\theta(\mathbf{x}_i)) + \mathbb{E}_{\text{in}} \sigma(g_{\hat{\theta}_\epsilon}(\mathbf{x})) - \frac{1}{n} \sum_{i=1}^n \sigma(g_\theta(\mathbf{x}_i)) \\ &\leq \alpha + c\epsilon_n + \mathbb{E}_{\text{in}} \sigma(g_{\hat{\theta}_\epsilon}(\mathbf{x})) - \frac{1}{n} \sum_{i=1}^n \sigma(g_\theta(\mathbf{x}_i)) \end{aligned} \quad (14)$$

$$\leq \alpha + c_2 \epsilon_n \quad (15)$$

where (14) follows from the definition of $\hat{\theta}_\epsilon$ and (13) follows from Σ_2 . This establishes claim 2 in the Proposition.

Claim 3 follows by a similar argument to claim 2. Finally, claim 1 follows by (13), (15), and Lemma B.1. \square

Lemma B.4. *Let $\delta \in (0, 1)$. Let $\mathbf{x}_1, \dots, \mathbf{x}_k \sim \mathbb{P}$ and $y_i \sim \text{Bernoulli}(p(\mathbf{x}_i))$ for $i \in [k]$. Let η_i be a Rademacher random variable, i.e., $\mathbb{P}(\eta_i = 1) = \mathbb{P}(\eta_i = -1) = 1/2$. Then,*

- With probability at least $1 - \delta$, for all $\theta \in \Theta$

$$\left| \frac{1}{k} \sum_{i=1}^k \sigma(-g_\theta(\mathbf{x}_i)) - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}} \sigma(-g_\theta(\mathbf{x})) \right| \leq 2 \mathbb{E}_{\mathbf{x}_1, \dots, \mathbf{x}_k \sim \mathbb{P}} \mathbb{E}_{\eta_1, \dots, \eta_m} \sup_{\theta \in \Theta} \frac{1}{k} \sum_{i=1}^k \eta_i g_\theta(\mathbf{x}_i) + c \sqrt{\frac{2 \ln(2/\delta)}{k}}$$

- With probability at least $1 - \delta$, for all $\theta \in \Theta$

$$\left| \frac{1}{k} \sum_{i=1}^k \sigma(g_\theta(\mathbf{x}_i)) - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}} \sigma(g_\theta(\mathbf{x})) \right| \leq 2 \mathbb{E}_{\mathbf{x}_1, \dots, \mathbf{x}_k \sim \mathbb{P}} \mathbb{E}_{\eta_1, \dots, \eta_m} \sup_{\theta \in \Theta} \frac{1}{k} \sum_{i=1}^k \eta_i g_\theta(\mathbf{x}_i) + c \sqrt{\frac{2 \ln(2/\delta)}{k}}$$

- With probability at least $1 - \delta$, for all $\theta \in \Theta$

$$\begin{aligned} \left| \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\text{cls}}(f_\theta(\mathbf{x}_i), y_i) - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}, y \sim \text{Bernoulli}(p(\mathbf{x}))} [\mathcal{L}_{\text{cls}}(f_\theta(\mathbf{x}), y)] \right| &\leq 2 \mathbb{E}_{\mathbf{x}_1, \dots, \mathbf{x}_k \sim \mathbb{P}} \mathbb{E}_{\eta_1, \dots, \eta_m} \sup_{\theta \in \Theta} \frac{1}{k} \sum_{i=1}^k \eta_i f_\theta(\mathbf{x}_i) \\ &\quad + c \sqrt{\frac{2 \ln(2/\delta)}{k}}. \end{aligned}$$

Proof. We show the first bullet point. The second and third bullet points follow by a similar argument. Using Mcdiarmid's inequality, we have that with probability at least $1 - \delta$ for all $\theta \in \Theta$

$$\left| \frac{1}{k} \sum_{i=1}^k \sigma(-g_\theta(\mathbf{x}_i)) - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}} \sigma(-g_\theta(\mathbf{x})) \right| \leq 2 \mathbb{E}_{\mathbf{x}_1, \dots, \mathbf{x}_k \sim \mathbb{P}} \mathbb{E}_{\eta_1, \dots, \eta_m} \sup_{\theta \in \Theta} \frac{1}{k} \sum_{i=1}^k \eta_i \sigma(g_\theta(\mathbf{x}_i)) + c \sqrt{\frac{2 \ln(2/\delta)}{k}}$$

Then, using the contraction Lemma and the fact that the sigmoid function σ is 1-Lipschitz, we have that

$$\mathbb{E}_{\mathbf{x}_1, \dots, \mathbf{x}_k \sim \mathbb{P}} \mathbb{E}_{\eta_1, \dots, \eta_m} \sup_{\theta \in \Theta} \frac{1}{k} \sum_{i=1}^k \eta_i \sigma(g_\theta(\mathbf{x}_i)) \leq \mathbb{E}_{\mathbf{x}_1, \dots, \mathbf{x}_k \sim \mathbb{P}} \mathbb{E}_{\eta_1, \dots, \eta_m} \sup_{\theta \in \Theta} \frac{1}{k} \sum_{i=1}^k \eta_i g_\theta(\mathbf{x}_i).$$

The result follows by combining the above two inequalities. \square

As an example where Assumption B.2 holds, consider for instance when $\theta = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix}$ with $w_1, w_2 \in \mathbb{R}^d$ and $g_\theta(x) = w_1^\top x$ and $f_\theta(x) = w_2^\top x$. Then if \mathbb{P}_0 is absolutely continuous wrt the Lebesgue measure, Assumption B.2 holds. We could similar replace the linear maps g and f with neural networks that share a penultimate layer. See Blanchard et al. (2010) for a more detailed discussion and for more examples.

C. Validation using \mathbb{P}_{wild} data

In this Section, we discuss how to use data from \mathbb{P}_{wild} for a validation procedure and demonstrate its feasibility. For simplicity, we focus on the OOD task since it is standard to have a clean ID validation set for classification and therefore this captures the main difficulty. To ease notation, we write $\mathbb{P}_{\text{out}} = \mathbb{P}_1$ and $\mathbb{P}_{\text{in}} = \mathbb{P}_0$. Here, overloading notation, we assume access to a holdout set from \mathbb{P}_0 and the mixture \mathbb{P}_{wild} :

- $\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathbb{P}_0$
- $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m \sim \mathbb{P}_{\text{wild}} := (1 - \pi)\mathbb{P}_0 + \pi\mathbb{P}_1$ ($\pi \in (0, 1]$ unknown)

We suppose that we have access to a small, finite set of models $\mathcal{G} \subset \{g : \mathbb{R}^d \mapsto \mathbb{R}\}$. \mathcal{G} is typically obtained from training a model with a set of distinct hyperparameters, generating one $g \in \mathcal{G}$ for each hyperparameter configuration. Note that \mathcal{G} is totally generic. As is typical in the OOD literature, we obtain from \mathcal{G} a set of OOD predictors by thresholding each $g \in \mathcal{G}$ as follows:

$$\mathcal{H} := \{\text{sign}(g(x) - \tau) : g \in \mathcal{G}, \tau \in \mathbb{R}\}.$$

We now introduce some notation, overloading notation from Section B. Define

$$R_y(g_\theta) := \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{out}}}(\mathbb{1}\{h(\mathbf{x}) \neq y\}).$$

Define the optimization problem

$$\begin{aligned} R_1^* &:= \inf_{h \in \mathcal{H}} R_1(h) \\ \text{s.t. } R_0(h) &\leq \alpha. \end{aligned}$$

Define risk for \mathbb{P}_{wild}

$$\begin{aligned} R_{\text{wild}} &:= \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{wild}}}(\mathbb{1}\{h(\mathbf{x}) \neq 1\}) \\ &= \pi R_1(h) + (1 - \pi)(1 - R_0(h)) \end{aligned}$$

Now, define another similar optimization, only changing the objective:

$$\begin{aligned} R_{\text{wild}}^* &:= \inf_{\theta} R_{\text{wild}}(h) \\ \text{s.t. } R_0(h) &\leq \alpha. \end{aligned}$$

We choose the $\hat{h} \in \mathcal{H}$ that minimizes the FNR@95 on the holdout set:

$$\begin{aligned} \hat{h}_\epsilon &\in \text{argmin}_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{h(\tilde{\mathbf{x}}_i) \neq 1\} \\ \text{s.t. } \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{h(\mathbf{x}_i) \neq -1\} &\leq \alpha + \epsilon \end{aligned}$$

where we write $\hat{h} := \hat{h}_0$. We emphasize that this procedure is not only intuitive; it is also justified theoretically by applying Theorem 2 from (Blanchard et al., 2010). Theorem 2 requires that the following condition is satisfied:

Assumption C.1. For any $\alpha \in (0, 1)$, there exists $h^* \in \mathcal{G}$ such that $R_0(h^*) = \alpha$ and $R_1(h^*) = R_{1,\alpha}^*(\mathcal{G})$.

Here, we show that \mathcal{H} satisfies Assumption C.1 if P_0 is absolutely continuous with respect to the Lebesgue measure. Fix some $g \in \mathcal{G}$ and define $h_\tau(x) := \mathbb{1}\{h(x) > \tau\}$. Notice that if $\tau > \tau'$, then

$$h_\tau(x) \leq h_{\tau'}(x).$$

Thus, as discussed in the Remark of page 2978 in (Blanchard et al., 2010), we have that if for a given τ we have that $R_0(h_\tau) < \alpha$, using the absolute continuity of P_0 , we can find a τ' such that $R_0(h_{\tau'}) = \alpha$ and $R_1(h_{\tau'}) \leq R_1(h_\tau)$. Since this holds for any $g \in \mathcal{G}$, this implies that Assumption C.1 holds. Then, as a Corollary from Theorem 2 of (Blanchard et al., 2010), we obtain

Corollary C.2. Let $\epsilon_k := \sqrt{\frac{\mathcal{VC}(\mathcal{H}) - \log(\delta)}{k}}$ where $\mathcal{VC}(\mathcal{H})$ denotes the VC dimension of \mathcal{H} . If $\epsilon = c\epsilon_n$, with probability at least $1 - \delta$

1. $R_1(\hat{h}_\epsilon) \leq R_1^* + c\pi^{-1}(\epsilon_n + \epsilon_m)$, and
2. $R_0(\hat{h}_\epsilon) \leq \alpha + c\pi^{-1}\epsilon_n$.