

# GENERALIZING CROSS ENTROPY LOSS WITH A BETA PROPER COMPOSITE LOSS: AN IMPROVED LOSS FUNCTION FOR OPEN SET RECOGNITION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Open set recognition involves identifying data instances encountered during test time that do not belong to known classes in the training set. The majority of recent deep learning approaches to open set recognition use a cross entropy loss to train their networks. Surprisingly, other loss functions are seldom used. In our work, we explore generalizing cross entropy with a Beta loss. This Beta loss is a proper composite loss with a Beta weight function. This weight function adds the flexibility of putting more emphasis on different parts of the observation-conditioned class probability (i.e.  $P(Y|\mathbf{X})$ ) range during training. We show that the flexibility gained through this is Beta loss function produces consistent improvements over cross entropy loss for open set recognition and produces state of the art results relative to recent methods.

## 1 INTRODUCTION

Machine learning systems deployed in the real world often encounter data instances that do not belong to any of the classes seen during training. For example, autonomous cars with visual perception systems must handle, to some degree, the unmodeled object classes that could be encountered. Failure to identify unknown classes not only presents a safety risk, but it also hinders performance if the input domain changes significantly. Thus, identifying these novel instances is a critical task for machine learning systems to achieve robustness in the real world, especially in high-stakes settings (Dietterich, 2017).

Open Set Recognition (OSR) (Scheirer et al., 2013; Boulton et al., 2019) attempts to address this problem of identifying instances belonging to previously unseen classes. Deep learning methods for OSR typically learn a representation such that data instances from novel classes are separable from known classes under this representation. Almost all of the recent deep learning OSR methods use a cross entropy loss combined with a function (e.g. softmax or sigmoid) to map unnormalized logits to probabilities. Surprisingly, other loss functions have not been extensively explored for OSR despite known issues with the use of cross entropy on representation learning for OSR (Goodfellow et al., 2015; Bendale & Boulton, 2016).

In order to improve on cross entropy, we explore its properties as a special case of a proper composite loss (Buja et al., 2005; Gneiting & Raftery, 2007; Reid & Williamson, 2010). Proper composite losses have an integral representation that reveals an implicit weight function over the observation-conditioned class probability range (i.e. the range of  $P(Y|\mathbf{X})$  where  $Y$  is the output and  $\mathbf{X}$  are the features). As we will show, cross-entropy has a U-shaped weight function that places higher weight at the extreme ends of the  $P(Y|\mathbf{X})$  range. It is unclear if this weight function is the best option for OSR. To our knowledge, there is no prior work that investigates changing this weight function to improve OSR.

In our work, we introduce a loss function specifically designed for OSR based on a proper composite loss with a Beta function used as the weight function. The Beta function is parameterized by  $\alpha$  and  $\beta$  parameters which results in a generalization of the weight function used by cross entropy. By changing these two parameters, we change the shape of the weight function, thereby creating a class of loss functions that can be tailored to the specific OSR task.

We call this loss a *Beta loss* and introduce a constrained optimization problem to set its parameters during training. The flexibility gained in tuning the weight function for a given OSR task results in consistent gains in performance over cross entropy. In addition, we demonstrate that these gains result in state of the art performance compared to recent OSR methods.

## 2 RELATED WORK

Scheirer et al. (2013) provided the first formalization of the OSR problem, including defining open set risk as the risk of misclassifying data that is far away from the training data. Since then, the literature on OSR has been extensive (see (Boult et al., 2019) for a survey) and we provide a brief survey with a focus on recent deep learning methods. One of the first deep learning OSR techniques is OpenMax (Bendale & Boult, 2016), which rescales activation vectors, particularly the extreme ones, through a Weibull distribution to reduce the open set risk. Another category of OSR methods leverage autoencoders for open set detection, including using Extreme Value Theory to model reconstruction error from class conditioned autoencoders (Oza & Patel, 2019), mapping classes into multivariate Gaussians in the latent space using a class-conditional VAE (CVAE) framework Sun et al. (2020) or combining CVAEs with a CapsuleNetwork (Guo et al., 2021). The Classification-Reconstruction learning for Open-Set Recognition (CROSR) method (Yoshihashi et al., 2019) trains a deep hierarchical reconstruction net to perform both reconstruction and supervised classification. Reciprocal Points Learning (RPL) (Chen et al., 2020) introduces the idea of reciprocal points, which are parts of the embedding space that are not from class  $k$  and thus more similar to the unknown classes. For each class, the RPL algorithm attempts to separate the space of known instances from the reciprocal points while regularizing this objective with a term that indirectly reduces the open space risk.

A closely related task to OSR is Out Of Distribution (OOD) detection. The task in OOD detection is to detect data instances that are not in the training distribution. In past work, this task often involves detecting data instances that come from a totally different dataset, rather than detecting unseen classes from the same dataset. A simple OOD technique is to detect open set instances by thresholding the maximum softmax probability (Hendrycks & Gimpel, 2017). ODIN builds on Hendrycks & Gimpel (2017) by applying temperature scaling to softmax and adding small perturbations to data instances to improve the max softmax score. Generalized ODIN (G-Odin) (Hsu et al., 2020) further improves on ODIN by removing the need for OOD data to tune hyperparameters. We believe that OOD can be cast as a subproblem under the more general OSR framework, and we refer to the general problem of identifying previously unseen test data as OSR.

The most closely related approach to our work is the focal loss (Lin et al., 2020), which was originally intended to handle imbalanced classes. The focal loss modifies cross entropy to be  $(1 - p)^\gamma \log(p)$ , where  $\gamma$  is a focusing parameter that focuses training on hard instances that are misclassified. We will show that our Beta loss is different and substantially outperforms focal loss.

The majority of these prior approaches use cross entropy for the loss function. To our knowledge, none of the previous work has investigated changing the weight function of cross entropy to improve OSR. Thus, our work investigates a previously unexplored direction for OSR. Finally, we mention other strategies from OSR and OOD that are orthogonal to our approach and could be combined with our work to improve performance even further. These strategies used include incorporating an auxiliary dataset of outliers (Hendrycks et al., 2019; Liu et al., 2020), using generative models to create informative synthetic examples that augment the training data (Ge et al., 2017; Neal et al., 2018; Yue et al., 2021) and using self-supervision to guide the latent representation learning (Perera et al., 2020; Tack et al., 2020).

## 3 BACKGROUND

We first introduce the notation that we will use. Vectors and matrices are denoted with boldface font while scalars are not. Let the training set  $\mathcal{D} = \{(\mathbf{X}_1, \mathbf{Y}_1), \dots, (\mathbf{X}_N, \mathbf{Y}_N)\}$  be made up of  $N$  training examples. The  $i$ th training instance  $(\mathbf{X}_i, \mathbf{Y}_i)$  consists of a  $D$ -dimensional feature vector  $\mathbf{X}_i = (X_i^1, \dots, X_i^D)$  and a class label vector  $\mathbf{Y}_i = (Y_i^1, \dots, Y_i^K)$  that is a one-hot encoding vector of size  $K$ , corresponding to the  $K$  known classes during training. In this vector, if the  $i$ th instance belongs to class  $k$ , then  $Y_i^k = 1$  and all other components have value 0. If we refer to features in

general, we use the notation  $\mathbf{X}$  (without the subscript index). Similarly, we use  $\mathbf{Y}$  to refer to class labels when we do not refer to a specific instance. Finally, the notation  $-k$  indicates "not from class  $k$ ".

During testing, there are an additional  $U$  classes not seen during training that data instances can be drawn from. The OSR task is to classify whether a data instance belongs to one of the  $K$  known classes or whether the instances come from an unseen class.

### 3.1 PROPER COMPOSITE LOSSES

We briefly provide an explanation of proper composite losses, following the presentation by Reid & Williamson (2010) with some minor notational modifications. For more details, we refer the reader to (Reid & Williamson, 2010; Buja et al., 2005). To simplify the explanation of proper composite losses, for the entirety of Section 3.1,  $Y$  is a binary output variable taking value 0 or 1 (rather than dealing with the  $K$ -component class label vector  $\mathbf{Y}$ ).

**Proper losses** Let  $\eta(\mathbf{X}, \theta) = P(Y = 1 | \mathbf{X}, \theta)$  be the class probability function. Here,  $\theta$  are the parameters of the model for the class probability function. We will drop  $\theta$  from the notation when it is not important for the discussion and simply denote the class probability function as  $\eta(\mathbf{X})$ . We make  $\eta$  an explicit function of the features  $\mathbf{X}$  to emphasize that it takes features as input.

We want an estimate of the class probability  $\eta(\mathbf{X})$  from training data by fitting a model  $q(\mathbf{X})$  using a loss function  $\ell(Y, q(\mathbf{X}))$ , where  $Y = \{0, 1\}$  is a binary response variable. The loss function  $\ell(Y, V)$  can be defined in terms of two partial loss functions  $\ell_1(V) = \ell(1, V)$  and  $\ell_0(V) = \ell(0, V)$  for when the ground truth label is 1 and 0 respectively. Specifically  $\ell(Y, V) = \mathbb{I}[y = 1] \ell_1(V) + \mathbb{I}[y = 0] \ell_0(V)$ . For example, the log loss used in binary cross entropy is defined as  $\ell(Y, \eta(\mathbf{X})) = -Y \log(\eta(\mathbf{X})) - (1 - Y) \log(1 - \eta(\mathbf{X}))$  with partial losses  $\ell_1(\eta(\mathbf{X})) = -\log(\eta(\mathbf{X}))$  and  $\ell_0(\eta(\mathbf{X})) = -\log(1 - \eta(\mathbf{X}))$ .

Given a loss function, we can define the *point-wise risk* as:

$$L(\eta(\mathbf{X}), V) = E_{Y \sim \eta(\mathbf{X})}[\ell(Y, V)] = \eta(\mathbf{X}) \ell_1(V) + (1 - \eta(\mathbf{X})) \ell_0(V)$$

The point-wise risk is the expectation of the loss with the labels  $Y$  being drawn from a Bernoulli distribution with parameter  $\eta(\mathbf{X})$ . Let  $\eta(\mathbf{X})$  be the true value of the class probability conditioned on  $\mathbf{X}$ . A *proper loss* (also known as a proper scoring rule in Statistics) is a loss  $\ell(Y, V)$  with the property that the point-wise risk  $L(\eta(\mathbf{X}), \hat{\eta}(\mathbf{X}))$  is minimized by  $\hat{\eta}(\mathbf{X}) = \eta(\mathbf{X})$  and is said to be Fisher consistent (Gneiting & Raftery, 2007).

Proper losses have an integral representation that reveals the use of a weight function (Shuford et al., 1966; Savage, 1971; Schervish, 1989). In order to describe this integral, we first define the cost-weighted misclassification loss, with cost parameter  $c \in (0, 1)$  as:

$$\ell_c(0, \hat{\eta}(\mathbf{X})) = c \mathbb{I}[\hat{\eta}(\mathbf{X}) \geq c] \quad (1)$$

$$\ell_c(1, \hat{\eta}(\mathbf{X})) = (1 - c) \mathbb{I}[\hat{\eta}(\mathbf{X}) < c] \quad (2)$$

Equation 1 is the cost of a false positive (i.e.  $c$ ) while Equation 2 is the cost of a false negative (i.e.  $(1 - c)$ ).

With this definition, a proper loss  $\ell(Y, \hat{\eta}(\mathbf{X}))$  can be represented as a weighted integral of cost-weighted misclassification losses (Theorem 6 of Reid & Williamson (2010)):

$$\ell(Y, \hat{\eta}(\mathbf{X})) = \int_0^1 \ell_c(Y, \hat{\eta}(\mathbf{X})) w(c) dc \quad (3)$$

In Equation 3,  $w(c)$  is a weight function  $w : (0, 1) \mapsto [0, \infty)$  that satisfies  $\int_c^{1-c} w(c) dc < \infty, \forall \epsilon > 0$ . A more intuitive representation of this property can be found in the corollary to Theorem 6 in Reid & Williamson (2010), which illustrates the effect of the weight function on the partial losses:

$$\ell_1(\hat{\eta}(\mathbf{X})) = \int_{\hat{\eta}(\mathbf{X})}^1 (1 - c) w(c) dc \quad (4)$$

$$\ell_0(\hat{\eta}(\mathbf{X})) = \int_0^{\hat{\eta}(\mathbf{X})} c w(c) dc \quad (5)$$

Thus, if certain values of  $c$  produce high weight according to  $w(c)$ , they will have more influence on the partial losses, thereby emphasizing the importance of these values. Since the range of  $c$  is the same as that of  $P(Y|\mathbf{X})$ , we can view  $w(c)$  as focusing in different parts of this probability range. A large variety of commonly used loss functions can be produced by changing the weight function (see Table 1 of Reid & Williamson (2010)). The key to our approach will be specializing this weight function to the specific task of OSR.

**Proper Composite Losses** An inverse link function  $\Psi^{-1}(q(\mathbf{X}))$  is often needed to map the real-valued output of a model (i.e.  $q(\mathbf{X})$ ) to the range  $[0, 1]$ . For example, in logistic regression  $q(\mathbf{X}) = \beta^T \mathbf{X}$  where the inverse link function for log loss is the sigmoid function  $\Psi^{-1}(p) = \frac{1}{1+e^{-p}}$ , with forward link being the logistic function  $\Psi(p) = \log(\frac{p}{1-p})$ . A loss function that uses an inverse link function is called a *composite loss* because it composes the inverse link function with  $q(\mathbf{X})$  i.e.  $\eta(\mathbf{X}) = \Psi^{-1}(q(\mathbf{X}))$ . We denote a composite loss as  $\ell^\Psi(Y, V) = \ell(Y, \Psi^{-1}(V))$ , along with partial composite losses  $\ell_1^\Psi(V)$  and  $\ell_0^\Psi(V)$ .

Making a composite loss proper requires accounting for the effects of the link function on the partial losses. If  $\Psi$  is a strictly monotone link function and  $\ell_1^\Psi(V)$  and  $\ell_0^\Psi(V)$  are differentiable partial losses, then Theorem 10 of Reid & Williamson (2010) states that  $\ell^\Psi$  is a *proper* composite loss if and only if:

$$\frac{-\ell_1^{\Psi'}(\hat{\eta}(\mathbf{X}))}{1 - \hat{\eta}(\mathbf{X})} = \frac{\ell_0^{\Psi'}(\hat{\eta}(\mathbf{X}))}{\hat{\eta}(\mathbf{X})} = \frac{w(\hat{\eta}(\mathbf{X}))}{\Psi'(\hat{\eta}(\mathbf{X}))} \quad (6)$$

for all  $\hat{\eta}(\mathbf{X}) \in (0, 1)$ . Equation 6 shows that once the weight function and the link function are specified, then the partial loss functions are also specified (except for additive constants). More importantly, the weight and link function pairing has implications on the convexity of the loss function (see Sections 15 and 16 of Buja et al. (2005) and Section 6 of Reid & Williamson (2010)).

### 3.2 THE BETA FAMILY OF PROPER COMPOSITE LOSSES

Buja et al. (2005) show how to create a family of proper losses by using a Beta function as the weight function. The beta function is parameterized by  $\alpha$  and  $\beta$  and has the following form:

$$w(c) = c^{\alpha-1}(1-c)^{\beta-1} \quad (7)$$

Using Equations 4 and 5, the corresponding partial loss functions are:

$$\ell_1(q) = \int_q^1 c^{\alpha-1}(1-c)^{\beta} dc \quad (8)$$

$$\ell_0(q) = \int_0^q c^{\alpha}(1-c)^{\beta-1} dc \quad (9)$$

Setting  $\alpha$  and  $\beta$  to specific values produces many commonly used losses such as boosting loss, cross entropy, squared error loss and misclassification loss. For instance, setting  $\alpha = 0$  and  $\beta = 0$  produces  $\ell_0(q) = -\log(1-q)$  and  $\ell_1(1-q) = -\log(q)$ , which results in binary cross entropy. When  $\alpha$  and  $\beta$  are integer multiples of  $\frac{1}{2}$ , the integrals in Equations 9 and 8 have closed form expressions.

Section 15 of Buja et al. (2005) describes conditions for convex loss functions involving the Beta family for weight functions and scaled logistic links  $\Psi^{-1}(V) = \frac{1}{1+e^{-V/\sigma}}$  where  $\sigma$  is a scaling parameter. We restate an important corollary from that section as a theorem below.

**Theorem 1** *Proper losses in the Beta family and scaled logistic links result in convex composite losses iff  $\alpha, \beta \in [-1, 0]$  and  $\sigma > 0$ .*

## 4 METHODOLOGY

Our aim is to validate the Beta family of proper composite losses for the task of Open Set Recognition. We apply the Beta family to high dimensional datasets, of images, under the framework of deep neural networks for classification.

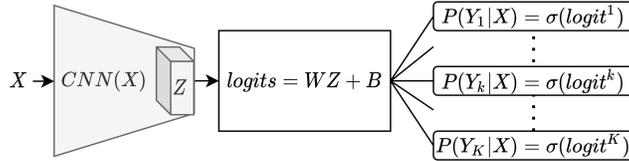


Figure 1: A diagram of the architecture used for Beta loss with a  $\sigma$  =sigmoid inverse link function.

#### 4.1 THE NEURAL NETWORK ARCHITECTURE

Figure 1 illustrates the neural network architecture that we use in our experiments. The architecture uses a CNN as its backbone and has  $K$  heads, corresponding to the  $K$  known classes in our training set, connected to a common latent embedding layer. These  $K$  heads correspond to  $K$  one-vs-all predictions, with the  $k$ th output predicting class- $k$  membership of the input and a given output’s prediction is unaffected by the value of any other output  $logit^{-k}$ . The output of the  $k$ th head is computed by applying the logistic inverse link function (i.e.  $\sigma(X) = (1 + \exp(-X))^{-1}$ ) to the  $k$ th logit, resulting in the output  $P(Y^k = 1|X) = \sigma(logit^k)$ .

This architecture is much simpler than some of the other deep architectures used for OSR. As we will show in Section 5, a simpler architecture with our novel loss function outperforms OSR methods with more sophisticated architectures.

#### 4.2 BETA LOSS

During training, each output contributes a one-vs-all loss term to the overall loss. This one-vs-all loss term is calculated using the Beta weight function (Equation 7) combined with a scaled logistic link function with  $\sigma = 1$ , which results in a sigmoid function. The partial losses  $\ell_0(q)$  and  $\ell_1(q)$  are shown in Equations 9 and 8. If we keep  $\alpha, \beta \in [0, 1]$ , then according to Theorem 1, this weight-link function combination results in a convex loss for a given output  $k$ .

The overall Beta loss (Equation 10) sums up the one-vs-all losses over all the  $K$  outputs. Since the sum of convex functions is also convex, the overall Beta loss is also convex.

$$L_{beta}(\mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K Y_i^k \ell_1(1 - \hat{\eta}^k(\mathbf{X}, \boldsymbol{\theta})) + (1 - Y_i^k) \ell_0(\hat{\eta}^k(\mathbf{X}, \boldsymbol{\theta})) \quad (10)$$

In Equation 10,  $\hat{\eta}^k(\mathbf{X}, \boldsymbol{\theta})$  is the prediction of  $P(Y_i^k = 1|\mathbf{X}, \boldsymbol{\theta})$  from the output of the neural network. The Beta loss could involve values of  $\alpha$  and  $\beta$  that are not multiples of  $\frac{1}{2}$  and thus do not have close form expressions for  $\ell_1$  and  $\ell_0$ . Fortunately, the integrals are one-dimensional and approximating these integrals to an arbitrary precision can be done by applying a simple numerical integration technique like the trapezoidal rule. However, in some cases, a large number of trapezoids may be needed to achieve the desired precision due to the rapid changes at the endpoints of the partial losses for a large number of classes  $K$ . To mitigate the problem of high memory usage, we have experimented with training a multi-layer perceptron to approximate the integral; the approximation is promising and we will explore it more fully in future work.

**Selecting  $\alpha$  and  $\beta$**  Selecting the values of  $\alpha$  and  $\beta$  before seeing the data can be challenging. Naively optimizing  $\alpha$ ,  $\beta$  and the neural network parameters  $\boldsymbol{\theta}$  simultaneously will produce a meaningless loss function as the optimization will minimize the loss by flattening the weight function towards zero in as many places as possible; in effect, this will cause  $\alpha$  and  $\beta$  to approach  $\infty$ .

In order to understand the effects of  $\alpha$  and  $\beta$ , we used CIFAR10 as a development dataset and apply a posthoc grid search over the two parameters to see their effect on AUC on an OOD task with CIFAR10 data as inliers (see Section 4.6). Table 1 provides a brief overview of our grid search for

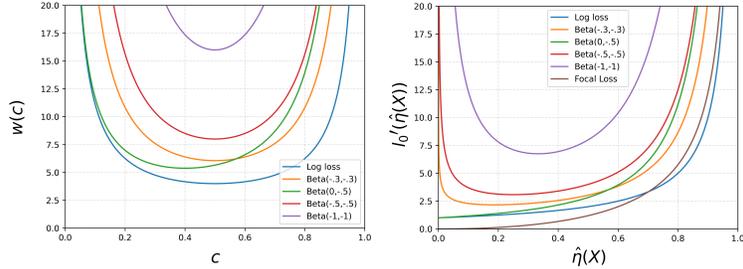


Figure 2: **(Left)**: The Beta weight function for a few values of  $\alpha$  and  $\beta$ . **(Right)**: The derivative of partial loss  $\ell_0$  (the plot for the derivative of  $\ell_1$  is similar but reflected). Beta losses with lower values of  $\alpha$  and  $\beta$  place much more emphasis on instances that are close to 0 (i.e. close to correct predictions) than Focal and Log loss.

the cases where  $\alpha = \beta$ . In general, our experiments showed that values near  $\alpha, \beta = -0.4$  produce the highest AUC values on CIFAR10.

Figure 2 (left) shows a plot of the weight function for a few select parameters of  $\alpha$  and  $\beta$ . A plot of the derivative of the partial loss function  $\ell_0$  (Figure 2 (right)) is also instructive as it shows how Beta functions with parameters close to  $-0.5$  skew the curve shape to be higher towards 0, which effectively increases the area under the curve towards the left for the integral. This shape effectively places more weight on examples that are close to correct i.e. having  $\hat{\eta}(\mathbf{X}) = P(Y_i^k = 1 | \mathbf{X})$  be close to 0 for an instance with true label 0. In effect, this weight function emphasizes "core" instances for that class.

We use this information to create a constrained optimization problem that learns the neural network parameters  $\theta$  as well as the values for  $\alpha$  and  $\beta$ . First, we constrain the values of  $\alpha, \beta \in [-1, 0]$  to preserve the convexity requirements from Theorem 1. Second, we add an L1 norm regularization term that penalizes  $\alpha$  and  $\beta$  from deviating too much from  $-1$ . Without this regularization term, the optimization will always return  $\alpha = 0$  and  $\beta = 0$  since these values minimize the objective. The  $\lambda$  parameter handles the tension between tethering the  $\alpha$  and  $\beta$  parameters to  $-1$  and minimizing the objective, thereby resulting in  $\alpha$  and  $\beta$  parameters in-between the endpoints.

$$\begin{aligned} \min_{\alpha, \beta, \theta} \quad & L_{beta}(\mathbf{Y}, \mathbf{X}, \theta) + \lambda * (|\alpha + 1| + |\beta + 1|) \\ \text{s.t.} \quad & -1 \leq \alpha \leq 0, -1 \leq \beta \leq 0 \end{aligned} \tag{11}$$

To set the value of  $\lambda$ , we again perform posthoc tuning on CIFAR10 for the same OOD task to get the value of  $\lambda$  that yields the best AUC for the Beta loss algorithm under investigation (see Table 2). We then apply this value of  $\lambda$  (tuned over CIFAR10) to the other datasets in our evaluation.

### 4.3 BASELINE METHODS AND DETECTION FUNCTIONS

We compare our OSR results with a large variety of methods, including: Sigmoid ( binary cross entropy with a sigmoid link ), Softmax (cross entropy followed applied over softmax), Focal Loss (Lin et al., 2020), Openmax (Bendale & Boult, 2016), G-Openmax (Ge et al., 2017), OSRCI (Neal et al., 2018), CROSR (Yoshihashi et al., 2019), C2AE (Oza & Patel, 2019), GFROR (Perera et al., 2020), CGDL (Sun et al., 2020), RPL (Chen et al., 2020), CVCap (Guo et al., 2021), G-Odin (Hsu et al., 2020), and CSI (Tack et al., 2020). For Sigmoid, Softmax, Focal loss, G-Odin and CSI, we ran these algorithms ourselves and report the AUC averaged over 5 random splits. For the other algorithms, we present results as reported in Table 1 from Guo et al. (2021).

$\alpha = \beta$	0.0	-0.1	-0.2	-0.3	-0.4	-0.5	-0.6	-0.7	-0.8	-0.9	-1.0
AUC	0.952	0.955	0.990	0.979	0.990	0.989	0.988	0.981	0.980	0.986	0.978

Table 1: The AUC for different values of  $\alpha$  and  $\beta$  on the AUC of CIFAR-10, averaged over the OOD datasets. Due to space constraints, we only include results for  $\alpha = \beta$ .

$\lambda$	0.01	0.05	0.1	0.5	1	5	10
Sigmoid $\beta$	0.953	0.967	0.982	0.988	<b>0.992</b>	0.975	0.974
Softmax $\beta$	0.987	0.989	0.990	0.990	<b>0.991</b>	0.979	0.970
G-Odin $\beta$	<b>0.993</b>	0.979	0.985	0.990	0.987	0.981	0.971
CSI $\beta$	0.966	0.966	<b>0.967</b>	0.966	0.962	0.964	0.960

Table 2: A posthoc sensitivity analysis of the effect of Beta loss  $\lambda$  on various algorithms’ AUC for CIFAR-10, averaged over the five OOD datasets.

Previous methods (Neal et al., 2018; Hendrycks & Gimpel, 2017) that use softmax simply use the probability of the max predicted class to detect open set instances, but we find using  $L_2$  norm of the latent representation  $z$ ,  $S(x) = \|CNN(x)\|_2$  improves performance. Unless otherwise noted, our experiments use the  $L_2$  metric as the detecting function for all methods.

For both Generalized ODIN and CSI, we use their proposed functions for OSR as we found these metrics were consistently better than  $L_2$ . Generalized ODIN uses the largest  $f(x)$  i.e. the largest logit of the  $K$  classes before the learned temperature is applied. CSI simply uses the max softmax probability,  $S(x) = \underset{y}{\operatorname{argmax}} P(y|x)$ , averaged over four rotations of an input image.

#### 4.4 REPLACING LOSS FUNCTIONS IN EXISTING METHODS WITH BETA LOSS

We also investigate replacing the cross entropy loss in existing OSR methods with Beta loss. Although this swap of loss functions may violate the theoretical assumptions of the Beta loss, it gives us an idea of the effect of the Beta loss while keep the rest of the algorithm in tact. In addition to cross entropy plus softmax, we replace cross entropy with Beta loss in two state of the art OSR techniques: G-Odin (Hsu et al., 2020) and CSI (Tack et al., 2020).

G-Odin is an OOD method that normalizes the final prediction layer of a CNN from  $f(x) = Wx + b$  to  $f(x) = \frac{Wx+b}{\|W\| \|x\|}$  for each class, which is normalized again by a single predictive logit  $g(x) = \operatorname{sigmoid}(W_g x + b_g)$ , and finally followed by a softmax activation. Although they also propose an image preprocessing strategy to improve OSR, we found that preprocessing produced consistently worse results than simply using  $f(x)$  as an outlier score.

CSI is a self-supervised method for learning a rich embedded representation of a dataset by contrasting a given sample with distributionally-shifted augmentations of itself. While their method generalizes to unsupervised learning, we use their supervised implementation, which fine tunes on the labels of a given dataset.

Other OSR methods from Section 4.3 do not have a standard classification loss that can be easily swapped with our Beta loss. For example, CVCap (Guo et al., 2021) used class labels by minimizing the KL-divergence between a class-conditional representation and a target fixed representation.

To determine the value of  $\lambda$  for the Beta loss versions of Softmax, G-Odin, and CSI, we perform a posthoc sensitivity analysis on  $\lambda$  over the CIFAR10 dataset. As with Sigmoid, we then select the  $\lambda$  value that produces the best AUC on CIFAR10 and apply it to the other datasets in our evaluation. Table 2 shows that the value of lambda that maximizes AUC for CIFAR10 is not always the same for different methods.

#### 4.5 EXPERIMENTAL SETUP

For all baseline methods, we use identical hyper-parameter settings and models to the original papers. Experiments described in Section 4.4 follow identical training procedures to the original papers for CSI and G-Odin and only replaced the cross entropy loss with Beta loss. For all of our other experiments, we use a DenseNet (Huang et al., 2017) and follow all the training procedures described by the original work. We use a batch size 64 for 300 epochs with a weight decay of 0.0001 for all weights except on the logits. Models are trained with Stochastic gradient descent, using a momentum of 0.9, learning rate 0.1, with a learning rate schedule at 50% and 75% epochs. Finally, all weights are initialized with He-initialization (He et al., 2015).

For our high resolution mars datasets, we use transfer learned models, which have been shown to elicit strong performance even on small datasets Huh et al. (2016). We use baseline publicly downloadable pre-trained models: a Resnet-50 (Targ et al., 2016) trained on ImageNet. This model is trained similarly to the DenseNet, but only needed 30 epochs to converge.

#### 4.6 DATASETS

**OSR** For our small-sized experiments, we use multiple common benchmark datasets, all of which are  $32 \times 32$  RGB images. For some experiments we use different splits of the training set for alternate analysis, following the setup described in (Neal et al., 2018). CIFAR6 randomly selects 6 classes as inliers to be trained on, with the other 4 considered unknown. CIFAR10+ and CIFAR50+ select the 4 non-animal classes as inlier, then uses the more diverse CIFAR100 to select outliers 10 and 50 animal classes, respectively, are randomly selected. We also use SVHN and TinyImageNet datasets, where 6 and 20 classes are selected from the 10 and 200 possible classes.

**OOD** We evaluate models trained on the entirety of both cifar10 and cifar100 datasets, and we compare these models against a variety of outlier datasets such as: cropped Imagenet, resized Imagenet, cropped LSUN, resized LSUN, iSUN and SVHN. While we report the averaged AUCs, full results, as well as results on the fixed versions of the resized datasets, are available in the appendix.

**Mars** We also use a recently released large-sized labeled dataset of Mars Surface images, from the Mars Science Laboratory (MSL) (Grotzinger et al., 2012) and Mars satellite images collected by the High Resolution Imaging Experiment (HiRISE) instrument onboard the Mars Reconnaissance Orbiter (McEwen et al., 2007). Both datasets (Wagstaff et al., 2021) are high-resolution at  $227 \times 227$ , with RGB and greyscale images respectively. The Rover dataset contains 2,900 images taken on the surface of the planet Mars, split into 19 classes. The Mars satellite dataset contains 10,815 images sorted 8 classes. We consulted a domain expert from NASA to split the dataset into meaningful openset tasks. For the Rover dataset, we set the “artifact” class (images of low quality or contain missing data) to be unknown and the rest to be inlier. We also split the dataset by “solday”, where only classes occurring during the first 100 martian days of data collection are used as inliers ( $K = 14$  classes), and the rest are used as outliers. Finally the orbital dataset sets “impact ejecta” and “spider” are used as outliers due to their low occurrence in the full dataset.

## 5 RESULTS AND DISCUSSION

OSR has a variety of metrics used for determining the effectiveness of different methods. The most popular being area under the receiver operating characteristic curve (AUC) (Boult et al., 2019). AUC provides a calibration-free measure and characterizes the performance for a given score by varying the discrimination threshold between inlier data and openset data. Other metrics exists such as true negative rate at 95% true positive rate (TNR@TPR95), which forces inlier recall to be directly at 95%. While having a high TNR is more challenging than a high AUC, we find the results to be highly correlated. Another metric is F-score, which combines outliers as a  $K + 1$  class and measures the accuracy. We do not find this metric to be meaningful as it requires a threshold for choosing outliers, and conflates the difficulties of the OSR task with strict classification. Therefore, we report only AUCs in our main text. TNRs and accuracy are reported in the supplementary material.

Table 3 shows the main results. All experiments are averaged over five random splits of known and unknown classes for both OSR (Neal et al., 2018) and OOD detection (Hendrycks & Gimpel, 2017). We draw attention to three elements. First, adding Beta Loss to each to a given method results in a very consistent improvement in performance for almost all tasks. Second, we achieve start-of-the-art results by using CSI with Beta loss for most OSR tasks, and we achieve start-of-the-art results by using Generalized Odin with Beta Loss on the standard OOD datasets. Lastly, using Beta Loss with the inverse logistic link function, in a one-vs-all classification setting, results in significant improvements over using standard binary cross entropy (Sigmoid  $\beta$  vs Sigmoid).

Table 4 shows the results of applying Beta Loss to the mars datasets. Here we find these datasets to be much more challenging. While no one method achieves the best results for all three mars-OSR tasks, adding Beta loss to a given method generally results in a performance increase, including

	Svhn	Cifar6	Cifar10+	Cifar50+	Tiny Imagenet	Cifar10 avg	Cifar100 avg
Softmax	0.886	0.677	0.816	0.805	0.577	-	-
Openmax	0.894	0.695	0.817	0.796	0.576	-	-
G-Openmax	0.896	0.675	0.827	0.819	0.58	-	-
OSRCI	0.91	0.699	0.838	0.827	0.586	-	-
CROSR	0.899	-	-	-	0.589	-	-
C2AE	0.892	0.711	0.81	0.803	0.581	-	-
GFROR	0.955	0.831	-	-	<b>0.657</b>	-	-
CGDL	0.896	0.681	0.794	0.794	0.653	-	-
RPL	0.931	0.784	0.885	0.881	0.624*	-	-
CVCap	0.956	0.835	0.888	0.889	0.558*	-	-
Focal*	0.915	0.830	0.891	0.883	0.636	0.949	0.877
Softmax*	0.931	0.838	0.924	0.917	0.646	0.971	0.890
Sigmoid*	0.890	0.806	0.857	0.855	0.612	0.934	0.730
G-Odin*	0.932	0.859	0.931	0.933	0.641	0.990	0.970
CSI*	<b>0.958</b>	0.857	0.967	0.963	0.623	0.977	0.896
Sigmoid $\beta$	0.927	0.843	0.939	0.933	0.625	0.989	0.751
Softmax $\beta$	0.935	0.837	0.942	0.940	0.646	0.991	0.971
G-Odin $\beta$	0.937	0.853	0.935	0.937	0.646	<b>0.992</b>	<b>0.973</b>
CSI $\beta$	<b>0.958</b>	<b>0.874</b>	<b>0.975</b>	<b>0.967</b>	0.628	0.988	0.923

Table 3: Full results of AUC averaged over 5 random splits. A Beta loss (indicated with a  $\beta$ ) shows a consistent improvement over log loss for nearly every dataset with pure classification loss. These results are mainly significant (see the supplement for full standard deviations, TNR at TPR95, and OOD results). The top third of the table is as reported in (Guo et al., 2021), including the empty cells. \* denotes our own runs of the experiment.

a substantial increase of about 0.1 AUC in the case of Softmax  $\beta$  on Artifact. Note that the  $\lambda$  parameter was not tuned on the mars datasets and we may experience an even larger improvement with tuning of  $\lambda$ .

	Artifact	HiRISE	Solday
Sigmoid	0.604	0.627	0.529
Softmax	0.655	0.654	0.583
G-Odin	0.843	0.608	0.523
Sigmoid $\beta$	0.615	0.634	<b>0.596</b>
Softmax $\beta$	0.759	0.640	0.588
G-Odin $\beta$	<b>0.850</b>	<b>0.655</b>	0.532

Table 4: AUC of Mars tasks averaged over 5 trials. Beta loss generally provides a consistent improvement over the baseline loss functions.

## 6 CONCLUSION

We have introduced a new loss function for OSR called the Beta loss, which generalizes cross entropy loss. The Beta loss is a proper composite loss with an associated weight function over the  $P(Y|X)$  range that can be tailored to a specific OSR task. This additional flexibility produces consistent increases in AUC over cross entropy and these increases yield state of the art results relative to recent OSR algorithms.

For future work, we will explore other link functions aside from the logistic link function that still preserve the convexity of the proper composite loss function. It remains an open question if there are other weight and link function combinations that are more effective for OSR than Beta loss.

## REFERENCES

- Abhijit Bendale and Terrance E Boult. Towards open set deep networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1563–1572, 2016.
- T. E. Boult, S. Cruz, A. Dhamija, M. Gunther, J. Henrydoss, and W. Scheirer. Learning and the unknown: Surveying steps toward open world recognition. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):9801–9807, 2019.
- A. Buja, W. Stuetzle, and Y. Shen. Loss functions for binary class probability estimation and classification: Structure and applications. Technical report, University of Pennsylvania, Nov 2005.
- G. Chen, L. Qiao, Y. Shi, P. Peng, J. Li, T. Huang, S. Pu, and Y. Tian. Learning open set network with discriminative reciprocal points. In *Proc. of European Conference in Computer Vision (ECCV)*, pp. 507–522, 2020.
- T. G. Dietterich. Steps toward robust artificial intelligence. *AI Magazine*, 38(3):3–24, 2017.
- Z. Ge, S. Demyanov, and R. Garnavi. Generative openmax for multi-class open set classification. In *Proc. of British Machine Vision Conference (BMVC)*, 2017.
- T. Gneiting and A. E. Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378, March 2007.
- I. Goodfellow, J. Shelnus, and C. Szegedy. Explaining and harnessing adversarial examples. In *Proceedings of the International Conference on Learning Representations*, 2015.
- John P Grotzinger, Joy Crisp, Ashwin R Vasavada, Robert C Anderson, Charles J Baker, Robert Barry, David F Blake, Pamela Conrad, Kenneth S Edgett, Bobak Ferdowski, et al. Mars science laboratory mission and science investigation. *Space science reviews*, 170(1):5–56, 2012.
- Yunrui Guo, Guglielmo Camporese, Wenjing Yang, Alessandro Sperduti, and Lamberto Ballan. Conditional variational capsule network for open set recognition. *CoRR*, abs/2104.09159, 2021.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pp. 1026–1034, 2015.
- Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *Proceedings of International Conference on Learning Representations*, 2017.
- Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier exposure. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=HyxCxhRcY7>.
- Yen-Chang Hsu, Yilin Shen, Hongxia Jin, and Zsolt Kira. Generalized odin: Detecting out-of-distribution image without learning from out-of-distribution data. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4700–4708, 2017.
- Minyoung Huh, Pulkit Agrawal, and Alexei A Efros. What makes imagenet good for transfer learning? *arXiv preprint arXiv:1608.08614*, 2016.
- S. Kornblith, M. Norouzi, H. Lee, and G. Hinton. Similarity of neural network representations revisited. In *Proceedings of the 36th International Conference on Machine Learning, PMLR 97*, pp. 3519–3529, 2019.
- Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(2):318–327, 2020. doi: 10.1109/TPAMI.2018.2858826.

- W. Liu, X. Wang, J. Owens, and Y. Li. Energy-based out-of-distribution detection. In *Advances in Neural Information Processing Systems 33*, 2020.
- Alfred S McEwen, Eric M Eliason, James W Bergstrom, Nathan T Bridges, Candice J Hansen, W Alan Delamere, John A Grant, Virginia C Gulick, Kenneth E Herkenhoff, Laszlo Keszthelyi, et al. Mars reconnaissance orbiter’s high resolution imaging science experiment (hirise). *Journal of Geophysical Research: Planets*, 112(E5), 2007.
- Lawrence Neal, Matthew Olson, Xiaoli Fern, Weng-Keen Wong, and Fuxin Li. Open set learning with counterfactual images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 613–628, 2018.
- P. Oza and V. M. Patel. C2ae: Class conditioned autoencoder for open-set recognition. In *Proc. of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2302–2311, 2019.
- Pramuditha Perera, Vlad I. Morariu, Rajiv Jain, Varun Manjunatha, Curtis Wigington, Vicente Ordonez, and Vishal M. Patel. Generative-discriminative feature representations for open-set recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11814–11823, June 2020. doi: 10.1109/CVPR42600.2020.01183.
- M. D. Reid and R. C. Williamson. Composite binary losses. *Journal of Machine Learning Research*, 11:2387–2422, Dec 2010.
- L. J. Savage. Elicitation of personal probabilities and expectations. *J. of the American Statistical Association*, 66(336):783–801, 1971.
- W. J. Scheirer, A. Rocha, A. Sapkota, and T. E. Boult. Towards open set recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(7):1757–1772, 2013.
- M. J. Schervish. A general method for comparing probability assessors. *The Annals of Statistics*, 17(4):1856–1879, 1989.
- E. Shuford, A. Albert, and H. E. Massengill. Admissible probability measurement procedures. *Psychometrika*, 31(2):125–145, June 1966.
- X. Sun, Z. Yang, C. Zhang, K.-V. Ling, and G. Peng. Conditional gaussian distribution learning for open set recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 13480–13489, 2020.
- J. Tack, S. Mo, J. Jeong, and J. Shin. Csi: Novelty detection via contrastive learning on distributionally shifted instances. In *Advances in Neural Information Processing Systems 33*, 2020.
- Sasha Targ, Diogo Almeida, and Kevin Lyman. Resnet in resnet: Generalizing residual architectures. *arXiv preprint arXiv:1603.08029*, 2016.
- Kiri Wagstaff, Steven Lu, Emily Dunkel, Kevin Grimes, Brandon Zhao, Jesse Cai, Shoshanna B Cole, Gary Doran, Raymond Francis, Jake Lee, et al. Mars image content classification: Three years of nasa deployment and recent advances. *arXiv preprint arXiv:2102.05011*, 2021.
- Ryota Yoshihashi, Wen Shao, Rei Kawakami, Shaodi You, Makoto Iida, and Takeshi Naemura. Classification-reconstruction learning for open-set recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4016–4025, 2019.
- Zhongqi Yue, Tan Wang, Qianru Sun, Xian-Sheng Hua, and Hanwang Zhang. Counterfactual zero-shot and open-set visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 15404–15414, 2021.

## A FULL EXPERIMENTAL RESULTS

The complete results of all over experiments, averaged over the 5 splits can be found in Tables 6, 7, 8, 9 and 10. Tables 6 and 7 contain all the AUCs and TNR metrics, respectively, for OOD experiments, both for CIFAR10 and CIFAR100. Tables 8 and 9 contain all the AUCs and TNR metrics, respectively, for OSR experiments. Finally, all the average accuracies for both OSR experiments, CIFAR10 and CIFAR100 can be found in Table 10.

### A.1 LATENT SPACE COMPARISONS

	linear CKA	rbf kernel CKA
test	0.716	0.764
train	0.758	0.806

Table 5: Comparing the latent representation of a baseline sigmoid model with a beta-weighted sigmoid model  $\lambda = 1$  on the cifar10 dataset.

We were also curious as to how much the latent space changes when the loss function changes from binary cross entropy to the Beta loss, with both techniques using a sigmoid function as the inverse link function. To answer this question, we applied the Centered Kernel Alignment (CKA) similarity index (Kornblith et al., 2019) to the latent representations obtained from the two techniques of baseline sigmoid and sigmoid with Beta loss. For each, we got the feature representation of all 50,000 data instances chosen from the training set and another 10,000 instances from the test set. Table 5 shows the results of CKA applied between both trainings and test datasets. Linear CKA produced similarity indices of 0.716 (test) and 0.758 (training) while RBF kernel CKA produced similarity indices of 0.764 (test) and 0.806 (training). These results indicate that while most of the latent representations are similar, the relatively smaller amount of latent space changes caused by training with Beta loss do have an impact on OSR performance.

indata	method	Imagenet	Imagenet resize	Imagenet resize fixed	LSUN	LSUN resize	LSUN resize fixed	iSUN	svhn
cifar10	focal	0.949±.013	0.945±.015	0.907±.004	0.948±.014	0.971±.007	0.898±.012	0.960±.012	0.922±.016
cifar10	sigmoid	0.936±.014	0.936±.018	0.875±.006	0.932±.007	0.962±.014	0.868±.010	0.950±.016	0.887±.041
cifar10	softmax	0.976±.007	0.973±.010	0.930±.002	0.979±.006	0.981±.008	0.911±.009	0.979±.010	0.937±.052
cifar10	godin	0.985±.002	0.985±.003	0.944±.002	0.989±.001	0.991±.002	0.930±.007	0.990±.002	0.998±.001
cifar10	CSI	0.983±.013	0.969±.030	0.939±.039	0.981±.008	0.981±.014	0.936±.030	0.978±.005	0.973±.025
cifar10	sigmoid $\beta$	0.988±.002	0.986±.003	0.923±.004	0.991±.001	0.990±.001	0.915±.010	0.990±.001	0.990±.007
cifar10	godin $\beta$	0.989±.004	0.989±.004	0.944±.001	0.990±.001	0.993±.001	0.936±.001	0.993±.001	0.999±.000
cifar10	softmax $\beta$	0.982±.003	0.978±.005	0.926±.001	0.989±.001	0.981±.004	0.902±.002	0.981±.004	0.998±.001
cifar10	CSI $\beta$	0.993±.002	0.982±.009	0.942±.036	0.996±.002	0.986±.005	0.943±.028	0.986±.005	0.984±.014
cifar100	focal	0.876±.020	0.856±.033	0.763±.004	0.857±.011	0.888±.028	0.702±.003	0.871±.027	0.916±.021
cifar100	sigmoid	0.738±.049	0.698±.066	0.725±.004	0.751±.015	0.733±.062	0.671±.006	0.706±.064	0.756±.040
cifar100	softmax	0.874±.027	0.840±.054	0.777±.007	0.939±.005	0.873±.035	0.713±.006	0.852±.044	0.963±.013
cifar100	godin	0.968±.003	0.972±.004	0.795±.005	0.955±.007	0.969±.006	0.714±.010	0.969±.005	0.987±.001
cifar100	CSI	0.953±.004	0.818±.044	0.802±.001	0.962±.004	0.859±.037	0.746±.011	0.873±.031	0.912±.014
cifar100	sigmoid $\beta$	0.673±.082	0.658±.060	0.667±.006	0.870±.029	0.699±.060	0.605±.004	0.680±.062	0.928±.027
cifar100	softmax $\beta$	0.858±.017	0.798±.025	0.773±.004	0.941±.005	0.839±.024	0.699±.007	0.810±.030	0.971±.007
cifar100	godin $\beta$	0.974±.006	0.977±.007	0.795±.002	0.948±.007	0.979±.011	0.700±.013	0.977±.009	0.986±.002
cifar100	CSI $\beta$	0.968±.004	0.865±.028	0.811±.001	0.987±.003	0.885±.025	0.753±.011	0.881±.029	0.951±.008

Table 6: Full AUC results on cifar10 and cifar100 OOD detection datasets averaged over 5 splits with standard deviations.

indata	method	Imagenet	Imagenet resize	Imagenet resize fixed	LSUN	LSUN resize	LSUN resize fixed	iSUN	svhn
cifar100	focal	0.704±.099	0.680±.099	0.564±.017	0.726±.072	0.834±.048	0.487±.070	0.776±.066	0.511±.118
cifar100	sigmoid	0.735±.102	0.733±.118	0.441±.008	0.735±.089	0.830±.105	0.425±.017	0.789±.114	0.431±.086
cifar100	softmax	0.876±.039	0.856±.062	0.612±.018	0.891±.039	0.898±.051	0.510±.045	0.884±.064	0.712±.217
cifar100	godin	0.922±.016	0.925±.020	0.697±.012	0.946±.009	0.958±.012	0.616±.037	0.952±.014	0.994±.002
cifar100	CSI	0.915±.083	0.855±.123	0.676±.116	0.872±.075	0.901±.070	0.666±.092	0.927±.249	0.849±.144
cifar100	sigmoid $\beta$	0.940±.010	0.929±.016	0.590±.024	0.950±.006	0.951±.006	0.530±.050	0.951±.007	0.950±.039
cifar100	softmax $\beta$	0.901±.021	0.876±.033	0.591±.009	0.943±.008	0.898±.026	0.476±.017	0.894±.026	0.990±.008
cifar100	godin $\beta$	0.952±.024	0.954±.027	0.701±.005	0.952±.006	0.977±.005	0.650±.005	0.973±.005	0.997±.001
cifar100	CSI $\beta$	0.974±.009	0.914±.047	0.693±.136	0.992±.004	0.936±.025	0.681±.086	0.935±.022	0.918±.082
cifar100	focal	0.399±.052	0.339±.079	0.208±.010	0.336±.026	0.407±.099	0.116±.011	0.350±.083	0.564±.151
cifar100	sigmoid	0.295±.075	0.240±.082	0.194±.009	0.292±.050	0.264±.089	0.105±.006	0.224±.084	0.272±.096
cifar100	softmax	0.436±.088	0.356±.128	0.231±.010	0.715±.021	0.407±.082	0.118±.004	0.370±.083	0.805±.058
cifar100	godin	0.812±.018	0.839±.022	0.265±.011	0.756±.035	0.818±.039	0.102±.012	0.821±.032	0.934±.009
cifar100	CSI	0.713±.020	0.285±.076	0.298±.007	0.768±.026	0.397±.080	0.142±.020	0.389±.072	0.521±.066
cifar100	sigmoid $\beta$	0.172±.075	0.113±.046	0.126±.013	0.562±.052	0.150±.031	0.068±.006	0.117±.024	0.691±.092
cifar100	softmax $\beta$	0.368±.072	0.216±.057	0.227±.008	0.714±.020	0.273±.053	0.105±.004	0.227±.053	0.848±.041
cifar100	godin $\beta$	0.851±.035	0.872±.042	0.265±.005	0.732±.032	0.886±.065	0.092±.014	0.881±.049	0.927±.013
cifar100	CSI $\beta$	0.822±.030	0.320±.084	0.304±.007	0.935±.017	0.412±.074	0.167±.019	0.409±.075	0.687±.058

Table 7: Full TNR at TPR 95% results on cifar10 and cifar100 OOD detection datasets averaged over 5 splits with standard deviations.

method	svhn	cifar6	cifar10+	cifar50+	TinyImagenet
focal	0.915 $\pm$ .019	0.830 $\pm$ .028	0.891 $\pm$ .012	0.883 $\pm$ .006	0.636 $\pm$ .018
sigmoid	0.890 $\pm$ .012	0.806 $\pm$ .030	0.857 $\pm$ .028	0.855 $\pm$ .022	0.612 $\pm$ .021
softmax	0.941 $\pm$ .009	0.848 $\pm$ .018	0.924 $\pm$ .006	0.917 $\pm$ .008	0.646 $\pm$ .024
godin	0.932 $\pm$ .009	0.859 $\pm$ .027	0.931 $\pm$ .015	0.933 $\pm$ .003	0.651 $\pm$ .016
CSI	0.958 $\pm$ .006	0.857 $\pm$ .036	0.967 $\pm$ .007	0.963 $\pm$ .001	0.623 $\pm$ .018
sigmoid $\beta$	0.927 $\pm$ .006	0.843 $\pm$ .024	0.939 $\pm$ .016	0.933 $\pm$ .007	0.625 $\pm$ .022
softmax $\beta$	0.935 $\pm$ .008	0.837 $\pm$ .020	0.942 $\pm$ .018	0.940 $\pm$ .006	0.646 $\pm$ .025
godin $\beta$	0.928 $\pm$ .010	0.859 $\pm$ .022	0.936 $\pm$ .009	0.936 $\pm$ .006	0.646 $\pm$ .020
CSI $\beta$	0.958 $\pm$ .008	0.874 $\pm$ .021	0.975 $\pm$ .007	0.967 $\pm$ .001	0.628 $\pm$ .019

Table 8: Full AUC results on OSR tasks averaged over 5 splits with standard deviations.

method	svhn	cifar6	cifar10+	cifar50+	tinyimagenet
focal	0.630 $\pm$ .067	0.327 $\pm$ .047	0.477 $\pm$ .047	0.452 $\pm$ .046	0.089 $\pm$ .010
sigmoid	0.635 $\pm$ .020	0.322 $\pm$ .048	0.455 $\pm$ .044	0.446 $\pm$ .036	0.074 $\pm$ .021
softmax	0.706 $\pm$ .029	0.371 $\pm$ .055	0.576 $\pm$ .026	0.558 $\pm$ .027	0.099 $\pm$ .022
godin	0.700 $\pm$ .023	0.366 $\pm$ .071	0.636 $\pm$ .073	0.655 $\pm$ .016	0.102 $\pm$ .012
CSI	0.760 $\pm$ .048	0.423 $\pm$ .064	0.822 $\pm$ .040	0.802 $\pm$ .008	0.086 $\pm$ .008
sigmoid $\beta$	0.645 $\pm$ .032	0.348 $\pm$ .037	0.698 $\pm$ .054	0.667 $\pm$ .031	0.088 $\pm$ .006
softmax $\beta$	0.672 $\pm$ .034	0.328 $\pm$ .056	0.699 $\pm$ .065	0.690 $\pm$ .022	0.086 $\pm$ .011
godin $\beta$	0.704 $\pm$ .026	0.379 $\pm$ .046	0.646 $\pm$ .055	0.656 $\pm$ .028	0.096 $\pm$ .015
CSI $\beta$	0.741 $\pm$ .067	0.466 $\pm$ .048	0.868 $\pm$ .039	0.844 $\pm$ .006	0.083 $\pm$ .008

Table 9: Full TNR at TPR 95% results on OSR tasks averaged over 5 splits with standard deviations.

method	svhn	cifar6	cifar10+ cifar50+	Tiny Imagenet	cifar10 all-classes	cifar100
focal	0.972 $\pm$ .003	0.959 $\pm$ .010	0.968 $\pm$ .001	0.507 $\pm$ .043	0.949 $\pm$ .002	0.763 $\pm$ .003
sigmoid	0.974 $\pm$ .003	0.964 $\pm$ .009	0.971 $\pm$ .003	0.473 $\pm$ .038	0.952 $\pm$ .001	0.743 $\pm$ .003
softmax	0.972 $\pm$ .002	0.962 $\pm$ .009	0.969 $\pm$ .001	0.535 $\pm$ .046	0.950 $\pm$ .001	0.766 $\pm$ .002
godin	0.971 $\pm$ .004	0.962 $\pm$ .009	0.968 $\pm$ .001	0.516 $\pm$ .040	0.950 $\pm$ .001	0.758 $\pm$ .002
CSI	0.973 $\pm$ .005	0.958 $\pm$ .004	0.975 $\pm$ .001	0.500 $\pm$ .002	0.936 $\pm$ .008	0.787 $\pm$ .002
sigmoid $\beta$	0.969 $\pm$ .004	0.953 $\pm$ .010	0.962 $\pm$ .001	0.512 $\pm$ .040	0.940 $\pm$ .001	0.556 $\pm$ .010
softmax $\beta$	0.969 $\pm$ .003	0.952 $\pm$ .010	0.964 $\pm$ .001	0.516 $\pm$ .037	0.940 $\pm$ .001	0.755 $\pm$ .004
godin $\beta$	0.972 $\pm$ .004	0.960 $\pm$ .011	0.969 $\pm$ .001	0.529 $\pm$ .039	0.953 $\pm$ .000	0.761 $\pm$ .000
CSI $\beta$	0.970 $\pm$ .003	0.957 $\pm$ .006	0.976 $\pm$ .001	0.498 $\pm$ .001	0.932 $\pm$ .002	0.780 $\pm$ .002

Table 10: Full accuracy results of our experiments on all OSR datasets the with standard deviations.