

Learning Placeholders for Open-Set Recognition

Da-Wei Zhou

Han-Jia Ye[†]

De-Chuan Zhan

State Key Laboratory for Novel Software Technology, Nanjing University

{zhoudw, yehj}@lamda.nju.edu.cn, zhandc@nju.edu.cn

Abstract

Traditional classifiers are deployed under closed-set setting, with both training and test classes belong to the same set. However, real-world applications probably face the input of unknown categories, and the model will recognize them as known ones. Under such circumstances, open-set recognition is proposed to maintain classification performance on known classes and reject unknowns. The closed-set models make overconfident predictions over familiar known class instances, so that calibration and thresholding across categories become essential issues when extending to an open-set environment. To this end, we proposed to learn PlaceholderRs for Open-SEt Recognition (PROSER), which prepares for the unknown classes by allocating placeholders for both data and classifier. In detail, learning data placeholders tries to anticipate open-set class data, thus transforms closed-set training into open-set training. Besides, to learn the invariant information between target and non-target classes, we reserve classifier placeholders as the class-specific boundary between known and unknown. The proposed PROSER efficiently generates novel class by manifold mixup, and adaptively sets the value of reserved open-set classifier during training. Experiments on various datasets validate the effectiveness of our proposed method.

1. Introduction

Recent years have witnessed the rapid development of supervised learning, aiming to obtain the knowledge of finite known classes. During the testing process, the well-trained model matches an incoming instance to the class with the highest posterior probability. However, this closed-world assumption comes to an end when the test set includes unseen categories [16, 2, 42, 35]. Since it is impossible to cover all classes in the world as training set [19, 33], the model would treat all novel category instances as known ones. As a result, the performance decays, which is unbearable in real-world applications. Open-set recognition [27, 3, 38, 26] is thus

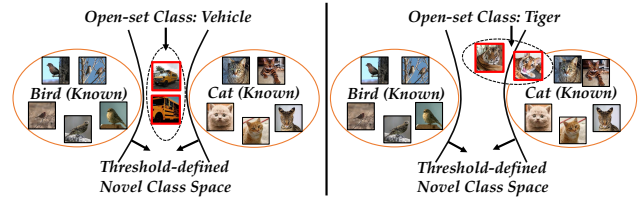


Figure 1. The drawback of threshold-based open-set recognition. It defines novel class space if confidence below a fixed threshold. Bird and cat are known classes, while vehicle and tiger are open-set classes in the left and right figures, respectively. Since the distributions of the vehicle and tiger categories differ, it is hard to rely on a single threshold to recognize unknown classes with diverse characteristics. The same threshold which well separates vehicles apart from known classes is not suitable for tigers.

proposed to conduct classification on known instances while at the same time detect those from unknown classes.

Facing the unknown input of novel categories, an intuitive way to separate known and unknown instances is to exert a threshold over the output probability [10]. It assumes the model produces a higher probability for known classes than unknowns. However, deep learning methods tend to overfit the training instances and produce overconfident predictions [27, 8, 9]. As a result, the model would output a high probability even for an unknown class instance, making the threshold hard to tune. Besides, the class-compositions are diverse, as shown in Figure 1. Since the semantic information of known classes differs in different tasks, it is hard to acquire an optimal threshold that suits all open-set tasks. Consequently, it is urgent to calibrate closed-set classifiers. Other methods try to foresee the distributions of novel classes and calibrate the output with open-set probability. [6] proposed G-OpenMax, which utilizes GAN to generate unknown samples for training novel classifier. [22] tried to generate images lying between decision boundaries as counterfactual instances. [26] combined self-supervision and augment input with generated open-set samples, which yields high disparity. These methods try to anticipate novel class distributions with generative models, and transform the closed-set training into open-set training.

The aim to boost open-set recognition can be summarized

[†]Correspondence to: Han-Jia Ye (yehj@lamda.nju.edu.cn)

as a calibration problem [8, 34]. Firstly, to make the closed-set model prepare for unknown classes, *data placeholders* of the novel class should be augmented and transform open-set into closed-set. Secondly, to better separate known and unknown instances, overconfident predictions should be calibrated by reserving *classifier placeholders* for novel classes.

Motivated by the problems above, we proposed to learn PlaceholderRs for Open-SEt Recognition (PROSER), aiming to calibrate open-set classifiers from two aspects. In detail, we augment the closed-set classifier with an extra *classifier placeholder*, which stands for the class-specific threshold between known and unknown. We reserve the placeholder for open-set classes to acquire the invariant information between target and non-target classes. Besides, to efficiently anticipate the distribution of novel classes, we consider generating *data placeholders*, which mimic open-set categories with a limited complexity cost. Consequently, we can transform closed-set classifiers into open-set ones, and *adaptively* predicts the class-specific threshold during testing. Experiments on various datasets validate the effectiveness of our proposed method on unknown detection and open-set recognition problems. Additionally, the visualization on decision boundaries indicates PROSER learns adaptive threshold for different class combinations.

In the following sections, we start with a brief review of related work, and then give the proposed PROSER and experiment results. After that, we conclude the paper.

2. Related Work

Open-set Recognition. There are two lines of work for open-set recognition, *i.e.*, discriminative models and generative models [7]. Discriminative models can be further divided into traditional machine learning-based methods and deep learning-based methods. There has been much progress in traditional methods. Based on SVM, [27] proposed 1-vs-Set machine, aiming to create a slab in the feature space. [2] extended the nearest class mean classifier [20], which calculates the distance between unknown and known class centers. [40] considered open-set recognition as a sparse representation learning problem, and tried to match the reconstruction error distributions with extreme value theory (EVT) [12]. In recent years, deep learning-based methods have attracted more attention due to the powerful representation ability. [3] first proposed to replace the softmax layer in the network with OpenMax, which calibrates the output probability with Weibull distribution. A similar work [28] replaced the softmax layer with one-vs-rest units. [36] utilized the latent representation for reconstruction, enabling robust unknown detection with known-class classification. These methods require a threshold to separate known and unknown, which face the challenge of threshold tuning.

Other works adopt generative models to anticipate the distribution of novel classes. G-OpenMax [6] followed the

guideline of OpenMax, which adopted a conditional generative network to synthesize unknown instances for network training. [22] proposed counterfactual images for open-set recognition (OSRCI). OSRCI generates instances lying in the decision boundary, and augments the original dataset with these generated instances. Class conditional auto-encoder (C2AE) [24] was proposed to tackle open-set recognition as a two-step problem, *i.e.*, closed-set training and open-set training. C2AE used class conditioned auto-encoders with novel training and testing methodology. Recently, [26] utilized self-supervision and augmented the input image to learn richer features to improve separation between classes. These methods work with extra generative models, and transform the closed-set classifier into open-set classifier with the generated novel instances.

Out-Of-Distribution Detection. The problem of out-of-distribution detection [10, 17, 1], anomaly detection [4, 18, 25] and novelty detection [11, 32, 41, 21, 43] are related topics to open-set recognition. They can be viewed as the unseen class detector in the open-set recognition problem. The essential difference between them and open-set recognition lies in that they are a binary classification problem. Since out-of-distribution detection methods are not designed for classifying known classes, they are not proper for open-set recognition.

3. From Closed-set to Open-set Recognition

In this section, we introduce the definition of closed-set classification and open-set recognition. Besides, we show how to train a closed-set classifier, and the limitations when extending it into open-set recognition.

3.1. Closed-Set Classification

Traditional closed-set classifiers are trained with $\mathcal{D}_{tr} = \{(\mathbf{x}_i, y_i)\}_{i=1}^L$ and tested with $\mathcal{D}_{te} = \{(\mathbf{x}_i, y_i)\}_{i=1}^M$, where $\mathbf{x}_i \in \mathbb{R}^D$ is a training instance, and $y_i \in Y = \{1, 2, \dots, K\}$ is the associated class label. In the closed-set assumption, \mathcal{D}_{tr} and \mathcal{D}_{te} are drawn from the same distribution \mathcal{D} . An algorithm should fit a model $f(x) : X \rightarrow Y$, which minimizes the expected risk:

$$f^* = \operatorname{argmin}_{f \in \mathcal{H}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_{te}} \mathbb{I}(y \neq f(\mathbf{x})), \quad (1)$$

where \mathcal{H} is the hypothesis space, $\mathbb{I}(\cdot)$ is the indicator function which outputs 1 if the expression holds and 0 otherwise. Assume the model f is composed of embedding function $\phi(\cdot) : \mathbb{R}^D \rightarrow \mathbb{R}^d$ and linear classifier $W \in \mathbb{R}^{d \times K}$, *i.e.*, $f(\mathbf{x}) = W^T \phi(\mathbf{x})$. We denote the k -th column of W as $\mathbf{w}_k \in \mathbb{R}^d$, *i.e.*, $W = [\mathbf{w}_1, \dots, \mathbf{w}_K]$, thus the output logits of class k is $\mathbf{w}_k^T \phi(\mathbf{x})$. Generally, it can be optimized with cross-entropy to gain the discrimination among known classes. A validation set \mathcal{D}_{val} drawn from closed set distribution \mathcal{D} can be used to measure the closed-set performance.

3.2. Open-Set Recognition

Facing the emergence of open-set classes, the model is still trained with $\mathcal{D}_{tr} = \{(\mathbf{x}_i, y_i)\}_{i=1}^L$. However, now $\hat{\mathcal{D}}_{te}$ is filled with instances of novel categories, *i.e.*, $\hat{\mathcal{D}}_{te} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, where $y_i \in \hat{Y} = \{1, \dots, K, K + 1\}$ is the associated class label. Note that class $K + 1$ is a group of novel categories, which may contain more than one class. Since there is no side-information in the training set, and we are unable to decompose the novel class group into sub-categories. An optimal open-set classifier minimizes the expected risk [38]:

$$\hat{f}^* = \operatorname{argmin}_{f \in \mathcal{H}} \mathbb{E}_{(\mathbf{x}, y) \sim \hat{\mathcal{D}}_{te}} \mathbb{I}(y \neq f(\mathbf{x})), \quad (2)$$

Since $\hat{\mathcal{D}}_{te}$ is composed of known classes and open-set class. The overall risk aims to classify known classes and meantime detect the unknown categories as class $K + 1$. Traditional classifiers predict the instance with highest posterior probability, *i.e.*, $\hat{y} = \operatorname{argmax}_{k=1, \dots, K} \mathbf{w}_k^\top \phi(\mathbf{x})$. However, since the model has not seen instances from open-set, it always predicts the lowest probability on class $K + 1$. As a result, directly employing closed-set classifiers into open-set recognition will predict all novel instances into known categories, yielding unsatisfactory performance in open-set recognition.

Consequently, an intuitive way to adopt closed-set classifier for open-set recognition is thresholding [10]. Taking the max output probability as confidence score, *i.e.*, $\text{conf} = \max_{k=1, \dots, K} \mathbf{w}_k^\top \phi(\mathbf{x})$. It assumes the model is more confident of closed-set instances than open-set. We can extend closed-set classifier by:

$$\hat{y} = \begin{cases} \operatorname{argmax}_{k=1, \dots, K} \mathbf{w}_k^\top \phi(\mathbf{x}) & \text{conf} > th \\ K + 1 & \text{otherwise} \end{cases}, \quad (3)$$

where th is the threshold. However, due to the overconfidence phenomena of deep neural networks, the output confidence of known and unknown is both close to 1 [3]. As a result, tuning a threshold that well separates known from unknown is hard and time-consuming. Furthermore, since the relationship between known and unknown may be different, the threshold in Eq. 3 relies on the essential similarity between known and unknown, which differs in diverse known-unknown compositions. To conclude, the closed-set classifier should be equipped with an extra calibration process to suit open-set recognition requirements.

4. Learning Placeholders for Open-Set Recognition

Facing the difficulty of closed-set classifier calibration, we need placeholders to prepare the closed-set model for novel classes. The key idea of PROSER is to design placeholders in two aspects, *i.e.*, data placeholders that anticipate

novel classes, and classifier placeholders that separate known from unknown. Learning data placeholders aims to mimic the emergence of novel classes, and transform closed-set training into open-set training. Reserving classifier placeholders for novel classes seeks to augment the closed-set classifier with dummy classifier, which adaptively outputs the class-specific threshold to separate known and unknown.

4.1. Learning Classifier Placeholders

To handle the diverse compositions of known-unknown categories, we need to extract invariant information from the target and non-target classes. Reserving classifier placeholders aims to arrange an extra dummy classifier and optimize it to represent the threshold between known and unknown. Assume we have the well-trained closed-set classifier; we first augment the output layer with an extra dummy classifier:

$$\hat{f}(\mathbf{x}) = [W^\top \phi(\mathbf{x}), \hat{\mathbf{w}}^\top \phi(\mathbf{x})]. \quad (4)$$

Note that the dummy classifier $\hat{\mathbf{w}}$ shares the same embedding $\phi(\cdot)$ with the closed-set classifier, and only create an extra linear layer $\hat{\mathbf{w}} \in \mathbb{R}^{d \times 1}$. These augmented logits will then pass the softmax layer to produce the posterior probability. The definition of dummy classifier is to well separate known and unknown, which is a dynamic threshold depend on input \mathbf{x} , acting as the classifier placeholder. To this aim, we fine-tune the model and make the dummy classifier to output the second-largest probability for known instances. Through this way, the invariant information between known class classifier and dummy classifier can be transferred into the detection process. Since the current output is augmented with dummy classifier, the classification loss can be expressed as:

$$l_1 = \sum_{(\mathbf{x}, y) \in \mathcal{D}_{tr}} \ell(\hat{f}(\mathbf{x}), y) + \beta \ell(\hat{f}(\mathbf{x}) \setminus y, K + 1), \quad (5)$$

where ℓ can be cross-entropy or other loss. The first item corresponds to optimizing the augmented output to match the ground truth, and maintaining performance in the closed-set. In the second term, $\hat{f}(\mathbf{x}) \setminus y$ means removing the probability of ground truth label, *i.e.*, set the predicted probability of ground-truth $\mathbf{w}_y^\top \phi(\mathbf{x})$ to 0. The second item matches the masked-probability with class $K + 1$, forcing the dummy classifier to output the second-largest probability. With the help of Eq. 5, the model learns to both correctly classify known instances and train the dummy classifier to place between the target and non-target classes. Note that loss is calculated with all training data \mathcal{D}_{tr} , which does not need novel class instances.

Effects of dummy classifiers: The explanation of Eq. 5 is shown in Figure 2(a). Eq. 5 creates a way to calibrate the closed-set classifier, where the first item aims to push

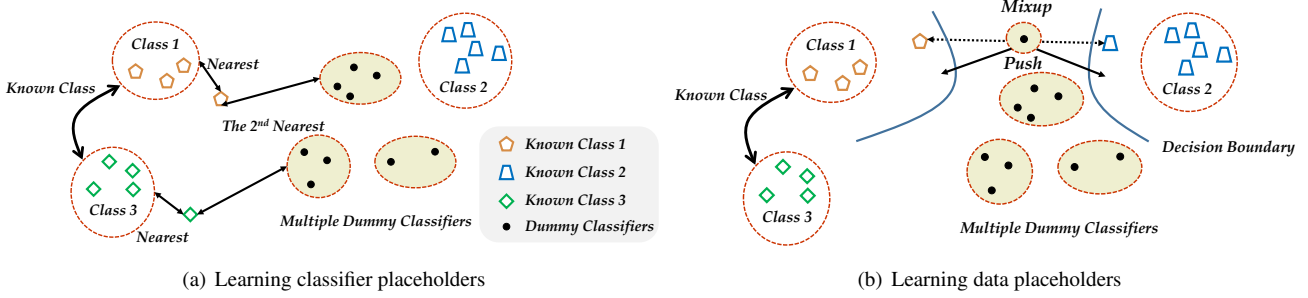


Figure 2. Illustration of proposed PROSER. The left figure corresponds to optimizing the classifier placeholders by Eq. 5 to output the second-largest probability, thus place them between target and non-target classes. The right figure corresponds to anticipating novel class patterns by manifold mixup as Eq. 7, which generates data placeholders near the manifold of the decision boundary, and pushing the decision boundary much tighter.

the instance towards its corresponding cluster to maintain correct classification. The second term seeks to relate the instance with the dummy classifier in the center space, and control the distance to the dummy classifier to be the second nearest among all class centers. As a result, it seeks a trade-off between correctly classifying closed-set instances and reserving novel classes’ probability as the classifier placeholder. In the training process, it learns to place between target class and non-target classes. When faced with novel classes, the prediction of dummy classifier would be high since all known classes are non-target ones. As a result, it acts as the instance-dependent threshold which can well fit every known class.

Learning multiple dummy classifiers: As discussed in Eq. 2, the $K + 1$ class may be composed of more than one novel class. We can learn more dummy classifiers by augmenting $\hat{W} \in \mathbb{R}^{d \times C}$, where C corresponds to the number of dummy classifiers. Under such circumstance, Eq. 4 turns into augment $f(x)$ with the highest dummy logit, *i.e.*, $\hat{f}(\mathbf{x}) = [W^\top \phi(\mathbf{x}), \max_{k=1, \dots, C} \hat{w}_k^\top \phi(\mathbf{x})]$, which only considers the nearest dummy classifier. With the help of multiple dummy classifiers, the model would adaptively choose the nearest dummy classifier to optimize. Learning multiple classifier placeholders boosts the variety of dummy classifiers from a united group into several scattered clusters. As a result, it leads to smoother decision boundaries, which in turn facilitates open-set recognition.

4.2. Learning Data Placeholders

The target of learning data placeholders is to change closed-set training into open-set training. The synthesized data placeholders should have two main characters, *i.e.*, the distribution of these instances seems novel, and the generating process should be quick. Traditional generative-based open-set models tend to utilize powerful generative models to mimic novel patterns [22, 26]. However, the distribution of natural images is hard to modeling [36]. To this end, we provide a simple yet effective way to anticipate novel class

instances without any extra time complexity.

Taking the above two characters into consideration, we mimic novel patterns with manifold mixup [31]. Assume the embedding module $\phi(\cdot)$ of the model can be decomposed by the middle hidden layer: $\phi(\mathbf{x}) = \phi_{post}(\phi_{pre}(\mathbf{x}))$, where ϕ_{pre} corresponds to the pre-layers before middle layer, which maps input into the hidden representation. Correspondingly, ϕ_{post} maps the hidden representation into the final embedding $\phi(\mathbf{x})$. We choose two instances from different class, and mix them up at the middle layer:

$$\tilde{\mathbf{x}}_{pre} = \lambda \phi_{pre}(\mathbf{x}_i) + (1 - \lambda) \phi_{pre}(\mathbf{x}_j), y_i \neq y_j, \quad (6)$$

where $\lambda \in [0, 1]$ is sampled from Beta distribution. The mixed $\tilde{\mathbf{x}}_{pre}$ will then pass the later layers, yielding $\phi_{post}(\tilde{\mathbf{x}}_{pre})$. Considering the interpolation between two different clusters are often regions of low-confidence predictions [31], *i.e.*, places of non-target classes. As a result, we can treat the embedding $\phi_{post}(\tilde{\mathbf{x}}_{pre})$ as the embedding of open-set classes, and train them as novel ones:

$$l_2 = \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}_{tr}} \ell([W, \hat{\mathbf{w}}]^\top \phi_{post}(\tilde{\mathbf{x}}_{pre}), K + 1) \quad (7)$$

$$\tilde{\mathbf{x}}_{pre} = \lambda \phi_{pre}(\mathbf{x}_i) + (1 - \lambda) \phi_{pre}(\mathbf{x}_j), y_i \neq y_j.$$

Note that we do not combine all possible pairs of $(\mathbf{x}_i, \mathbf{x}_j)$ in the whole dataset to form $\tilde{\mathbf{x}}_{pre}$. On the contrary, these combinations are produced within mini-batches [39, 31], *i.e.*, once we get the training batch of size B , another order of instances can be derived by shuffling this mini-batch. We then mask the pairs of the same class, and conduct manifold mixup with the pairs from different classes. As a result, the calculation complexity is of the same magnitude as vanilla training and would not cost extra time.

Effects of data placeholders: It is obvious that Eq. 7 does not consume extra time complexity, which generates novel instances lying between the manifold of decision boundaries. Besides, manifold mixup can better generate novel patterns in the improved embedding space leveraging interpolations

of deeper hidden representations, which better stands for the novel distribution. The visualization of data placeholders is shown in Figure 2(b), where the mixed instances push the decision boundary in the embedding space towards the composition class y_i and y_j . With the help of data placeholders, the embeddings of known classes would be much tighter, leaving more place for novel classes.

Discussion about vanilla mixup: Vanilla mixup [39, 30] aims to produce augmented instances with linear interpolations between two known ones in the input space: $\tilde{\mathbf{x}} = \lambda \mathbf{x}_i + (1 - \lambda) \mathbf{x}_j$. However, a fatal problem may occur that mixed $\tilde{\mathbf{x}}$ in the input space may situate near another class y_k between y_i and y_j , *i.e.*, $\lambda \mathbf{x}_i + (1 - \lambda) \mathbf{x}_j \approx \mathbf{x}_k$. Since we will treat mixed instances as class $K + 1$, optimizing such $\tilde{\mathbf{x}}$ harms the discriminability among closed-set classes and semantic information of original inputs. On the contrary, [31] proved that manifold mixup can move the decision boundary away from the data in all directions, resulting in a compact embedding space. Besides, since the features in the input space are fixed, the generated novel patterns by vanilla mixup are thus deterministic and cannot be optimized. By contrast, the generated novel patterns can be optimized with the embedding module $\phi(\cdot)$.

4.3. Calibration and Guideline for Implementation

Previously, we talk about the elements in PROSER, and how they are separately trained. We then discuss the combination of these two parts and extra calibration tricks to further improve the classifier.

Rethinking the overconfidence problem of closed-set classifier, we should make dummy classifier output the same magnitude logits as closed-set classifier. To fully explore the magnitude relationship between closed-set classifier $W^\top \phi(\mathbf{x})$ and dummy classifier $\hat{w}^\top \phi(\mathbf{x})$, we consider a further calibration step. With the help of validation set \mathcal{D}_{val} , we can calculate the highest logit of known classes and dummy classifier. To calibrate them into same magnitude, we calculate the difference of them: $\max_{k=1, \dots, K} \mathbf{w}_k^\top \phi(\mathbf{x}) - \max_{k=1, \dots, C} \hat{\mathbf{w}}_k^\top \phi(\mathbf{x})$, and divide it into several equal intervals as *bias*. The calibrated logits is $[W^\top \phi(\mathbf{x}), \hat{w}^\top \phi(\mathbf{x}) + bias]$, we obtain the best *bias* by ensuring 95% of validation data in \mathcal{D}_{val} to be recognized as known [29, 26]. As a result, the dummy logit will be further calibrated to the same magnitude as closed-set classifiers with bias tuning. Note that the bias tuning is conducted with the validation set \mathcal{D}_{val} drawn from the same distribution as \mathcal{D}_{tr} and with no access to the unknown dataset.

In the training process of PROSER, we first augment the dummy classifiers with pretrained closed-set classifiers. After that, we separate each batch into two equal parts, and separately calculate l_1 and l_2 over corresponding parts. The manifold mixup loss l_2 is calculated by shuffling the mini-batch and mixup instances of different classes, which takes

Algorithm 1 Training PROSER for open-set recognition

Input: Closed-set classifier: f ; Dummy classifier number: C ; Closed-set training set: $\mathcal{D}_{tr} = \{(\mathbf{x}_i, y_i)\}_{i=1}^L$;

Output: Open-set classifier: \hat{f} ; Calibration bias: *bias*;

- 1: Initialize $\hat{W} \in \mathbb{R}^{d \times C}$ as dummy classifier;
 - 2: Augment the output of f with dummy classifier as Eq. 4;
 - 3: **repeat**
 - 4: Get a batch of training instance $\{(\mathbf{x}_i, y_i)\}_{i=1}^B$;
 - 5: Separate the batch into two parts with equal size;
 - 6: Calculate the dummy loss on the first part $l_1 \leftarrow$ Eq. 5;
 - 7: Conduct manifold mixup on the second part, calculate mix loss $l_2 \leftarrow$ Eq. 7;
 - 8: Calculate the overall loss $l_{total} = l_1 + \gamma * l_2$;
 - 9: Obtain derivative and update the model;
 - 10: **until** reaches predefined epoches
 - 11: Calibrate the output of dummy classifier over validation set by ensuring 95% of validation data recognized as known;
-

no extra time. After the training process, we utilize the validation set to get the calibration bias. The guideline for implementation is shown in Algorithm 1. The optimal bias would be used to accumulate dummy classifier logits during open-set testing.

5. Experiment

In this section, we compare PROSER on benchmark datasets with state-of-the-art methods. We separately evaluate the performance of unknown detection and open-set recognition tasks. Ablations show the model robustly tackles datasets with different complexity, and visualizations indicate PROSER’s ability in adaptive threshold choosing.

5.1. Unknown Detection

Following the protocol defined in [22, 36, 26], we evaluate the performance of related methods. Several classes are sampled from a multi-class dataset and the others are viewed as open-set class [22]. We simulate the sampling process over five trials [22], and report the mean results. The commonly used metric *i.e.*, average area under the ROC curve (a.k.a AUC) evaluates the recognition performance. Since real-world scenarios are complex, where ratio of seen and unseen differs in diverse tasks, we utilize openness* [26, 22, 27, 7] to represent the complexity of the open-set task:

$$\text{Openness} = 1 - \sqrt{\frac{N_{train}}{N_{test}}}, \quad (8)$$

where N_{train} and N_{test} corresponds to the number of classes in training set and testing set. As we discussed in the preliminaries, $N_{train} = K$. We compare to other methods on the following benchmark datasets:

*There are two similar definitions of openness, and we follow [22, 26].

Table 1. Unknown detection performance in terms of the mean AUC. Results are averaged among five randomized trials. We report the full table with standard deviation in the supplementary.

Methods	SVHN	CIFAR10	CIFAR+10	CIFAR+50	Tiny-ImageNet
Softmax	88.6	67.7	81.6	80.5	57.7
OpenMax [3]	89.4	69.5	81.7	79.6	57.6
G-OpenMax [6]	89.6	67.5	82.7	81.9	58.0
OSRCI [22]	91.0	69.9	83.8	82.7	58.6
C2AE [24]	89.2	71.1	81.0	80.3	58.1
GFROSR [26]	93.5	83.1	91.5	91.3	64.7
PROSER	94.3	89.1	96.0	95.3	69.3

Table 2. Closed-set accuracy between the plain CNN (closed-set classifier) and PROSER. Although PROSER aims at learning placeholders for novel classes, there is no significant degradation in closed-set accuracy.

Methods	SVHN	CIFAR10	Tiny-ImageNet
Plain CNN	96.5	92.8	52.2
PROSER	96.4	92.6	52.1

- **SVHN [23]** and **CIFAR10 [13]**: There are total 10 classes in these datasets. SVHN contains street view house numbers, and CIFAR10 contains images of vehicles and animals. We randomly sample six classes to be known and the other four classes to be open-set classes. The openness of these tasks is 22.54%.
- **CIFAR+10 [22]** and **CIFAR+50**: To create a dataset with higher openness, four classes from CIFAR10 is sampled as known class. Besides, 10 and 50 classes are sampled from CIFAR100 as open-set classes, yielding CIFAR+10 and CIFAR+50. The openness for them are 46.55% and 72.78%, respectively.
- **Tiny-ImageNet [15]**: Tiny-ImageNet is a subset of ImageNet [5] with 200 classes. We sample 20 classes as known and the other 180 as open-set classes. The openness is 68.37% for this dataset.

We compare to the SOTA methods, *i.e.*, Softmax, OpenMax [3], G-OpenMax [6], OSRCI [22], C2AE [24], CROSR [36], GFROSR [26]: **Softmax**: utilizes the highest softmax probability as the confidence score for detection; **OpenMax**: replaces softmax layer with OpenMax and calibrates the confidence to predict novel class; **G-OpenMax**: trains conditional GAN to generate open-set instances and adopts OpenMax for recognition; **OSRCI**: generates instances lying near the decision boundary; **C2AE**: uses class conditioned auto-encoders with novel training and testing methodology; **CROSR**: utilizes the latent representation learning for reconstruction, which enables robust unknown

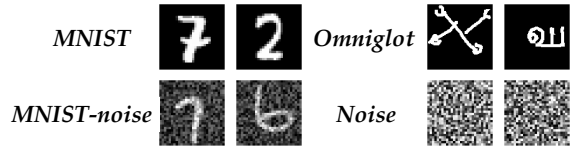


Figure 3. Dataset example of original MNIST, Omniglot, MNIST-noise and Noise.

detection; **GFROSR**: adopts self-supervision and augments the input image to learn richer features to improve separation between classes;

We adopt the same network backbone and dataset splits as GFROSR [26]. Table 1 shows the performance of unknown detection. We report the baseline performance from [26, 36, 22]. From Table 1, we can infer that the results on traditional digital number datasets SVHN are almost saturated, and little progress can be further achieved. However, our proposed PROSER improves the recognition ability on natural images by a substantial margin, *i.e.*, in the experiments on CIFAR10, we push forward 6% than SOTA method GFROSR. We also improve the performance over CIFAR+10, CIFAR+50, and Tiny-ImageNet by 4%. We also provide the closed set accuracy in Table 2, which indicates that learning PROSER does not sacrifice the model’s discriminative ability in the closed-set classification.

5.2. Open-Set Recognition

The ultimate goal of open-set recognition is to not only detect unseen classes, but also correctly classify known ones with superior performance. In this section, we validate the performance of our proposed PROSER with open-set recognition tasks. Following the protocol defined in [26], the models are trained by all training instances of the original dataset. While in the testing process, instances from another dataset are augmented to the original test set as open-set classes. In this setting, *macro-averaged F1-scores* over all known classes and the augmented novel class is used to measure the performance. Two benchmark datasets, *i.e.*, MNIST and CIFAR10 are used to simulate this setting.

Table 3. Open-set recognition results on CIFAR10 with various outliers added to the test set as unknowns. The performance is evaluated by macro F1 in 11 classes (10 known classes and unknown).

Methods	ImageNet-crop	ImageNet-resize	LSUN-crop	LSUN-resize
Softmax	63.9	65.3	64.2	64.7
OpenMax [3]	66.0	68.4	65.7	66.8
OSRCI [22]	63.6	63.5	65.0	64.8
LadderNet+Softmax [36]	64.0	64.6	64.4	64.7
LadderNet+OpenMax [36]	65.3	67.0	65.2	65.9
DHRNet+Softmax [36]	64.5	64.9	65.0	64.9
DHRNet+OpenMax [36]	65.5	67.5	65.6	66.4
CROSR [36]	72.1	73.5	72.0	74.9
GFROSR [26]	75.7	79.2	75.1	80.5
PROSER	84.9	82.4	86.7	85.6

Table 4. Open-set recognition on MNIST with various outliers added to the test set as unknowns. We report macro F1 in 11 classes (0–9 and unknown). The results are cited from [29].

Methods	Omniglot	MNIST-noise	Noise
Softmax	59.5	64.1	82.9
OpenMax	68.0	72.0	82.6
CROSR	79.3	82.7	82.6
PROSER	86.2	87.4	88.2

We first consider MNIST, which is formed with digital numbers between 0 – 9. According to [36], we choose the open-set classes from Omniglot [14], MNIST-noise, and Noise. Omniglot is a dataset of alphabet characters, and Noise is randomly generated by sampling each pixel between [0, 1] from a uniform distribution. Superimposing MNIST images over Noise yields the MNIST-noise dataset, which is similar to the original MNIST. We show the dataset example in Figure 3. We make the known-unknown ratio 1:1 by setting the number of test examples from open-set to 10,000, the same as the number of MNIST test set. We evaluate the performance with macro F1 scores between ten known classes and one unknown class. The results of open-set recognition are shown in Table 4. It indicates that open-set recognition on Omniglot and MNIST-noise are more challenging among the three datasets, since they have higher openness than Noise. However, our proposed PROSER handles these scenarios with the best performance.

We also conduct experiments with CIFAR10, which is a dataset of vehicles and animals. We choose open-set classes from ImageNet [5] and LSUN [37] according to [17]. LSUN has a testing set of 10,000 images of 10 different scenes. Since the image size of novel classes does not match CIFAR10, we design ImageNet-crop, ImageNet-resize, LSUN-crop, and LSUN-resize to align input size. For ‘crop’ datasets, we crop the original image into 32*32, and we

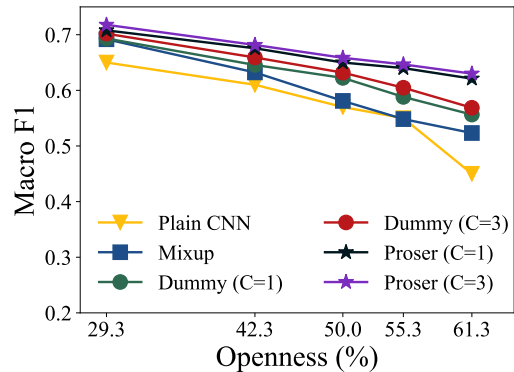


Figure 4. Macro F1 against varying openness with different base-lines for ablation analysis.

resize the image into 32*32 for ‘resize’ datasets. Like the experiments with MNIST, we set the size of novel instances to 10,000, yielding known-unknown ratio 1:1. We evaluate the performance with macro F1 scores between ten known classes and one unknown class. The results concerning CIFAR10 are reported in Table 3. We can infer from Table 3 that PROSER can handle open-set classes from diverse inputs and achieve better performance than SOTA methods.

5.3. Ablation study

In this section, we conduct an ablation study and analyze each part’s contribution with the CIFAR100 dataset. CIFAR100 contains 100 classes, and we change the composition of known and unknown by varying openness in Eq. 8. We randomly choose 15 out of 100 classes as known, and switch the number of unknown classes between {15, 30, 45, 60, 85}. The openness of these datasets is within the limits of 29.29% and 61.28%. We evaluate the performance by macro F1-scores over 15 known classes and unknown. The results are shown in Figure 4.

In the figure, ‘Plain CNN’ stands for thresholding softmax probabilities as Eq. 3. ‘Mixup’ stands for learning

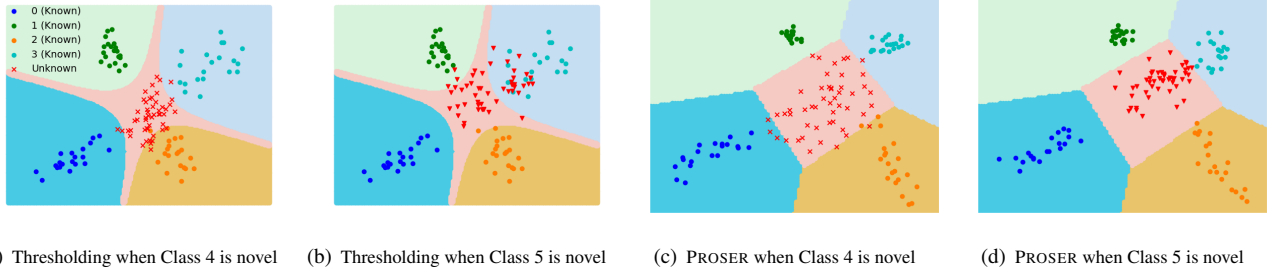


Figure 5. Visualization on open-set recognition of MNIST, known classes are visualized by dots with different colors, and the unknown classes are denoted by red cross/triangle. The shadow region represents the classification boundary of the model. Thresholding methods define unknown space if confidence is below some threshold. As shown in (a) and (b), the unknown class (class 4 and class 5) may come from diverse distributions and hard to be detected with a uniform threshold. It indicates a suitable threshold to detect open-set class 4 does not suit novelty from class 5. While (c) and (d) show PROSER can bear the distribution change of novel classes with the help of adaptive threshold choosing. Best viewed in color.

data placeholders only, and ‘Dummy’ stands for learning classifier placeholders only. The parameter C corresponds to the dummy classifier number. The trend of all methods with openness increasing declines, since the task becomes more complex. We can infer from the figure that adopting manifold mixup and dummy classifiers both improves the open-set recognition ability of plain CNN. The results indicate that data placeholders and classifier placeholders can both help calibration. Besides, learning more than one dummy classifier improves performance than one, indicating the diversity of dummy classifier matters. Furthermore, combining these two parts can further improve the performance, which outperforms the others. Ablations validate that placeholders help to calibrate the open-set classifier.

5.4. Visualization of Decision Boundaries

In this part, we visualize the learned decision boundaries on the MNIST dataset. Instances are shown in 2D by learning embedding module $\phi(\cdot) : \mathbb{R}^D \rightarrow \mathbb{R}^2$, *i.e.*, a linear layer attached to the CNNs as embedding. We plot the comparison between traditional threshold-based methods and our PROSER in Figure 5. In each figure, four known classes (in dots) and one open-set class (in red cross/triangle) are visualized. A threshold-based method detects open-set class 4 in Figure 5(a). The same decision boundary (threshold) is adopted in 5(b), where known classes are the same while class 5 is the novel class. We can infer from Figure 5(a) and 5(b) that the distribution of open-set classes differs, where a well-defined threshold for novel class 4 is not suitable for novel class 5. As a result, it is impossible to acquire an optimal threshold that well separates all kinds of unknown classes from known ones.

For comparison, we show the decision boundaries of PROSER in Figure 5(c), 5(d). With the help of data placeholders and classifier placeholders, PROSER is able to acquire the invariant information between target and non-target classes.

As a result, it learns to adaptively output the instance-specific threshold, and bears the distribution change of novel classes in the novel space. Figure 5 validates the effectiveness of PROSER facing various kind of open-set classes.

6. Conclusion

In real-world applications, instances from unseen novel classes may be fed to closed-set classifiers and be misclassified as known ones. Open-set recognition aims to simultaneously classify known classes and detect unknown ones. However, there are two main challenges in open-set recognition, *i.e.*, how to anticipating novel patterns and how to compensate for the overconfidence phenomena. In this paper, we propose PROSER to calibrate the closed-set classifiers in two aspects. On the one hand, PROSER efficiently mimics the distribution of novel classes as data placeholders, and transforms closed-set training into open-set training. On the other hand, we augment the closed-set classifier with classifier placeholders, which adaptively separates the known form unknown, and stands for the class-specific threshold. The proposed PROSER efficiently generates novel class by manifold mixup, and adaptively sets the value of reserved open-set classifier. How to extend open-set recognition into stream data scenarios, and utilize the detected novel patterns are interesting future works.

Acknowledgments

This research was supported by National Key R&D Program of China (2020AAA0109401), NSFC (61773198, 61921006, 62006112), NSFC-NRF Joint Research Project under Grant 61861146001, Nanjing University Innovation Program for Ph.D. candidate (CXYJ21-53), Collaborative Innovation Center of Novel Software Technology and Industrialization, NSF of Jiangsu Province (BK20200313).

References

- [1] Haoyue Bai, Rui Sun, Lanqing Hong, Fengwei Zhou, Nanyang Ye, Han-Jia Ye, S-H Gary Chan, and Zhenguo Li. Decaug: Out-of-distribution generalization via decomposed feature representation and semantic augmentation. *arXiv preprint arXiv:2012.09382*, 2020. [2](#)
- [2] Abhijit Bendale and Terrance Boulton. Towards open world recognition. In *CVPR*, pages 1893–1902, 2015. [1](#), [2](#)
- [3] Abhijit Bendale and Terrance E Boulton. Towards open set deep networks. In *CVPR*, pages 1563–1572, 2016. [1](#), [2](#), [3](#), [6](#), [7](#)
- [4] Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. *ACM computing surveys (CSUR)*, 41(3):1–58, 2009. [2](#)
- [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009. [6](#), [7](#)
- [6] Zongyuan Ge, Sergey Demianov, Zetao Chen, and Rahil Garnavi. Generative openmax for multi-class open set classification. In *BMVC*, 2017. [1](#), [2](#), [6](#), [3](#)
- [7] Chuanxing Geng, Sheng-jun Huang, and Songcan Chen. Recent advances in open set recognition: A survey. *TPAMI*, page in press, 2020. [2](#), [5](#)
- [8] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *ICML*, pages 1321–1330, 2017. [1](#), [2](#)
- [9] Matthias Hein, Maksym Andriushchenko, and Julian Bitterwolf. Why relu networks yield high-confidence predictions far away from the training data and how to mitigate the problem. In *CVPR*, pages 41–50, 2019. [1](#)
- [10] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *ICLR*, 2017. [1](#), [2](#), [3](#)
- [11] Victoria Hodge and Jim Austin. A survey of outlier detection methodologies. *Artificial intelligence review*, 22(2):85–126, 2004. [2](#)
- [12] Samuel Kotz and Saralees Nadarajah. *Extreme value distributions: theory and applications*. World Scientific, 2000. [2](#)
- [13] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, 2009. [6](#)
- [14] Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015. [7](#)
- [15] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7, 2015. [6](#)
- [16] Fayin Li and Harry Wechsler. Open set face recognition using transduction. *TPAMI*, 27(11):1686–1697, 2005. [1](#)
- [17] Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. *ICLR*, 2018. [2](#), [7](#)
- [18] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. Isolation forest. In *ICDM*, pages 413–422, 2008. [2](#)
- [19] John McCarthy and Patrick J Hayes. Some philosophical problems from the standpoint of artificial intelligence. In *Readings in artificial intelligence*, pages 431–450. Elsevier, 1981. [1](#)
- [20] Thomas Mensink, Jakob Verbeek, Florent Perronnin, and Gabriela Csurka. Distance-based image classification: Generalizing to new classes at near-zero cost. *TPAMI*, 35(11):2624–2637, 2013. [2](#)
- [21] Xin Mu, Kai Ming Ting, and Zhi-Hua Zhou. Classification under streaming emerging new classes: A solution using completely-random trees. *TKDE*, 29(8):1605–1618, 2017. [2](#)
- [22] Lawrence Neal, Matthew Olson, Xiaoli Fern, Weng-Keen Wong, and Fuxin Li. Open set learning with counterfactual images. In *ECCV*, pages 613–628, 2018. [1](#), [2](#), [4](#), [5](#), [6](#), [7](#), [3](#)
- [23] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011. [6](#)
- [24] Poojan Oza and Vishal M Patel. C2ae: Class conditioned auto-encoder for open-set recognition. In *CVPR*, pages 2307–2316, 2019. [2](#), [6](#), [3](#)
- [25] Guansong Pang, Chunhua Shen, Longbing Cao, and Anton van den Hengel. Deep learning for anomaly detection: A review. *arXiv preprint arXiv:2007.02500*, 2020. [2](#)
- [26] 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 *CVPR*, pages 11814–11823, 2020. [1](#), [2](#), [4](#), [5](#), [6](#), [7](#), [3](#)
- [27] Walter J Scheirer, Anderson de Rezende Rocha, Archana Sapkota, and Terrance E Boulton. Toward open set recognition. *TPAMI*, 35(7):1757–1772, 2012. [1](#), [2](#), [5](#)
- [28] Lei Shu, Hu Xu, and Bing Liu. Doc: Deep open classification of text documents. In *EMNLP*, pages 2911–2916, 2017. [2](#)
- [29] Xin Sun, Zhenning Yang, Chi Zhang, Keck-Voon Ling, and Guohao Peng. Conditional gaussian distribution learning for open set recognition. In *CVPR*, pages 13480–13489, 2020. [5](#), [7](#)
- [30] Yuji Tokozume, Yoshitaka Ushiku, and Tatsuya Harada. Between-class learning for image classification. In *CVPR*, pages 5486–5494, 2018. [5](#), [2](#)
- [31] Vikas Verma, Alex Lamb, Christopher Beckham, Amir Najafi, Ioannis Mitliagkas, David Lopez-Paz, and Yoshua Bengio. Manifold mixup: Better representations by interpolating hidden states. In *ICML*, pages 6438–6447, 2019. [4](#), [5](#), [1](#), [2](#)
- [32] Xiu-Shen Wei, Han-Jia Ye, Xin Mu, Jianxin Wu, Chunhua Shen, and Zhi-Hua Zhou. Multiple instance learning with emerging novel class. *TKDE*, 2019. [2](#)
- [33] Yang Yang, Da-Wei Zhou, De-Chuan Zhan, Hui Xiong, and Yuan Jiang. Adaptive deep models for incremental learning: Considering capacity scalability and sustainability. In *KDD*, pages 74–82, 2019. [1](#)
- [34] Han-Jia Ye, Hexiang Hu, De-Chuan Zhan, and Fei Sha. Learning adaptive classifiers synthesis for generalized few-shot learning. *arXiv preprint arXiv:1906.02944*, 2019. [2](#)
- [35] Han-Jia Ye, De-Chuan Zhan, Yuan Jiang, and Zhi-Hua Zhou. Rectify heterogeneous models with semantic mapping. In *ICML*, pages 5630–5639, 2018. [1](#)
- [36] Ryota Yoshihashi, Wen Shao, Rei Kawakami, Shaodi You, Makoto Iida, and Takeshi Naemura. Classification-reconstruction learning for open-set recognition. In *CVPR*, pages 4016–4025, 2019. [2](#), [4](#), [5](#), [6](#), [7](#), [1](#), [3](#)

- [37] Fisher Yu, Yinda Zhang, Shuran Song, Ari Seff, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *CoRR*, abs/1506.03365, 2015. [7](#)
- [38] Yang Yu, Wei-Yang Qu, Nan Li, and Zimin Guo. Open-category classification by adversarial sample generation. In *IJCAI*, pages 3357–3363, 2017. [1](#), [3](#)
- [39] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *ICLR*, 2018. [4](#), [5](#), [2](#)
- [40] He Zhang and Vishal M Patel. Sparse representation-based open set recognition. *TPAMI*, 39(8):1690–1696, 2016. [2](#)
- [41] Da-Wei Zhou, Yang Yang, and De-Chuan Zhan. Detecting sequentially novel classes with stable generalization ability. In *PAKDD*, page To appear, 2021. [2](#)
- [42] Zhi-Hua Zhou. Learnware: on the future of machine learning. *Frontiers Comput. Sci.*, 10(4):589–590, 2016. [1](#)
- [43] Yue Zhu, Kai Ming Ting, and Zhi-Hua Zhou. Multi-label learning with emerging new labels. *TKDE*, 30(10):1901–1914, 2018. [2](#)

Supplementary Material

I. Additional Experimental Results

This section introduces some additional experiment results and then gives the implementation details.

I.1. Sensitivity about Hyper-Parameters

In this section, we conduct experiments to explore the influence of hyper-parameters with the CIFAR100 dataset. The implementation details are the same as the ablations of the main paper. We choose 15 classes out of 100 as known ones, and another 15 out of 85 are open-set categories, making known-unknown ratio 1 : 1. The performance is measured with macro F1 over 15 known classes and unknown.

In the main paper, the overall loss is described as:

$$l_{total} = l_1 + \gamma * l_2, \quad (1)$$

where $l_1 = \sum_{(\mathbf{x}, y) \in \mathcal{D}_{tr}} \ell(\hat{f}(\mathbf{x}), y) + \beta \ell(\hat{f}(\mathbf{x}) \setminus y, K + 1)$ is the classifier placeholder loss, and $l_2 = \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}_{tr}} \ell([W, \hat{\mathbf{w}}]^\top \phi_{post}(\tilde{\mathbf{x}}_{pre}), K + 1)$ is the data placeholder loss, yielding three different parts. In this section, we report the effects of these hyper-parameters β and γ as well as dummy classifier number C in Figure 1. Considering the trade-off parameters are adopted to balance the loss of each parts, we tune them in range of $\{0, 10^{-2}, 10^{-1}, 10^0, 10^1\}$. As a result, we can get 25 results, corresponding to each combination from the set.

Note that we have provided an ablation study in the main paper, where ‘Mixup’ stands for only equipping the plain CNN with data placeholders, *i.e.*, $\gamma \rightarrow \infty$. Correspondingly, ‘Dummy’ stands for only training dummy classifiers, *i.e.*, $\gamma = 0$. The results in Figure 1(a) are consistent with the former conclusions that only employ part of the placeholder, *i.e.*, data or classifier, is not enough to produce the best performance. Additionally, we can observe that $\beta = 1, \gamma = 0.1$ leads to the best performance of the current task, which means a combination of data and classifier placeholders can jointly improve the model’s performance. This also guides the hyper-parameters setup in other tasks. We also show the influence of classifier placeholder number C in Figure 1(b). Since we arrange 15 classes as the open-set class, we tune it in the range of $\{1, 3, \dots, 15\}$. The results validate that multiple dummy classifiers increase the diversity of classifier placeholders, and can match novel patterns with the nearest classifier. However, learning too many dummy classifiers does not help open-set recognition, which shows a decline when $C > 11$.

I.2. Running Time Comparison

[36] point out that generative-based methods need more time in model training and instance generating. As a result, it takes more time to implement these methods in real-world

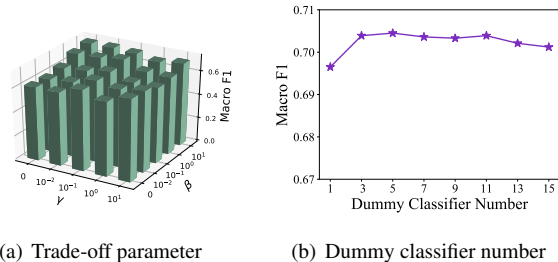


Figure 1. Sensitivity of hyper-parameters on CIFAR100 dataset. We report macro F1 between 15 known classes and unknown.

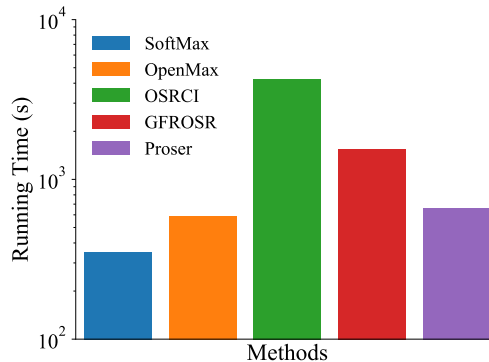


Figure 2. Running time comparison for different methods on MNIST dataset. The y-axis is in the logarithmic scale.

applications. We conduct experiments on MNIST, and compare to [22, 26, 3] as well as softmax in terms of the training time. The results are reported in Figure 2. The running time of generative methods includes training generative models and novel instance generation.

From Figure 2, we can tell that PROSER is of the same order of magnitude with Softmax and OpenMax [3]. In comparison, generative-based methods OSRCI [22] and GFROSR [26] consume much more time than PROSER. OSRCI needs to generate counterfactual images as novelty, and augment the initial dataset with these generated open-set instances, which consumes the most time. GFROSR trains an extra generative model to reconstruct images, and the reconstructed images are fed into the classification model, which consumes the second-most time. Since we only generate open-set classes with manifold mixup, the mixup loss is based on mixed embedding, which means we do not need more steps in training novel patterns. As a result, PROSER can generate novel classes and train models efficiently.

I.3. Manifold Mixup VS Vanilla Mixup

In the main paper, we discuss the pros and cons of manifold mixup [31] and the reason we do not adopt vanilla

Algorithm 2 Manifold mixup for data placeholders

Input: Embedding module $\phi(\cdot)$, which can be decomposed of pre-embedding $\phi_{pre}(\cdot)$ and post-embedding $\phi_{post}(\cdot)$;

Closed-set mini-batch: $\mathcal{D}_{tr} = \{(\mathbf{x}_i, y_i)\}_{i=1}^B$;

Output: Updated classifier \hat{f} ;

- 1: Calculate the pre-embeddings of this mini-batch $\phi_{pre}(\mathcal{D}_{tr})$;
- 2: Shuffle the mini-batch with random order, and get shuffled (embedding,label) pair $\mathcal{D}_{shuffle} = \{(\phi_{pre}(\hat{\mathbf{x}}_i), \hat{y}_i)\}_{i=1}^B$;
- 3: **for** $i = 1, \dots, B$ **do**
- 4: Mask the pairs of the same class, *i.e.*, $y_i = \hat{y}_i$;
- 5: **end for**
- 6: Sample λ from Beta distribution;
- 7: Calculate the manifold mixup pre-embeddings $\tilde{\mathbf{x}}_{pre}$ with unmasked pairs, *i.e.*, data placeholders;
- 8: Calculate the post-embeddings of $\tilde{\mathbf{x}}_{pre}$, *i.e.*, $\phi_{post}(\tilde{\mathbf{x}}_{pre})$;
- 9: Calculate the manifold mixup loss \leftarrow Eq. 7;

Table 1. Ablation study over manifold mixup and vanilla mixup, the configurations are the same as in Figure 1.

α	0.4	1	1.5	2
Vanilla Mixup	62.3	64.4	63.9	64.1
Manifold Mixup	63.0	68.6	70.0	70.7

mixup in the input space. In this section, we give the detailed pseudo code for manifold mixup for data placeholders and the performance comparison between manifold mixup and vanilla mixup [39, 30].

The guideline of generating data placeholders with manifold mixup is shown in Algorithm 2. Comparing to vanilla training, where mini-batch instances are fed into the model to forward passing and conduct back-propagation, our proposed method consumes the same complexity of forwarding and backward passing. Line 1 forward the mini-batch with the pre-embedding module, and get the middle representation of the original batch. These middle-representations are then shuffled with random order to form $\mathcal{D}_{shuffle}$, as shown in Line 2. To avoid mixing two instances from the same class, we mask the pairs of the same class with Line 4, and then conduct mixup to generate data placeholders in Line 7.

Note that the size of mixed embeddings $\tilde{\mathbf{x}}_{pre}$ should be no more than B , since we only combine unmasked instances. These data placeholders are then fed into the post-embedding module to get the ultimate representation in Line 8. As a result, the total forward and backward consumption is no more than vanilla training.

We also test the performance comparison between manifold mixup and vanilla mixup, *i.e.*, replace the ϕ_{pre} into

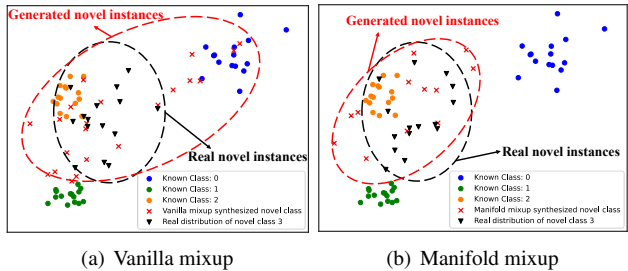


Figure 3. Visualization of generated novel instances and real distribution of novel instances. Left: novel instances generated by vanilla mixup. Right: novel instances generated by manifold mixup. Red crosses stand for generated novel instances, colored dots stand for known class instances, and black crosses stand for real distribution of novel class. The generated space of vanilla mixup covers known space, which may hurt the learning of embedding space.

identity mapping and mixup in the input space. As we stated in the main paper, we argue that manifold mixup is optimizable, and its benefits help the process of novelty generation. [31] proved that manifold mixup can move the decision boundary away from the data in all directions, which results in a compact embedding space. With the help of these data placeholders, the embeddings of known classes would be much tighter, thus leaving more place for embeddings of unknown classes. Moreover, we conduct ablations to validate the effectiveness of manifold mixup in Table 1. The experiment configuration is the same as Figure 1, and we tune different α in the Beta distribution to choose the best hyper-parameter α . Since the mixup percentage λ is influenced by α , we show the kernel density estimation with different α in Figure 4. We can infer that $\alpha < 1$ tends to sample λ near 0 or 1, while $\alpha > 1$ tends to sample λ near 0.5. Correspondingly, $\alpha = 1$ would result in a uniform distribution.

Let us return to our intuition where we try to synthesis novel patterns with known instances. Considering the places beside one class should not be a new class, while the middle point between two classes is often low-confidence area, an intuitive way we seek to mimic novel classes is to use $\alpha > 1$, as shown in Figure 4(c). To our relief, the results in Table 1 also validate the assumption that $\alpha = 2$ leads to the best performance. The results also indicate that compared to vanilla mixup, manifold mixup leads to a more compact embedding space, which boosts open-set recognition. The results also guide the setup of α in all experiments. We adopt $\alpha = 2$ in all experiments in the main paper without tuning the best task-specific value.

We also show the embedding of generated novel instances in Figure 3. We conduct experiments under the same setting as the visualization experiment part in the main paper, and show the embedding of generated instances by vanilla mixup and manifold mixup. The β parameter is the same

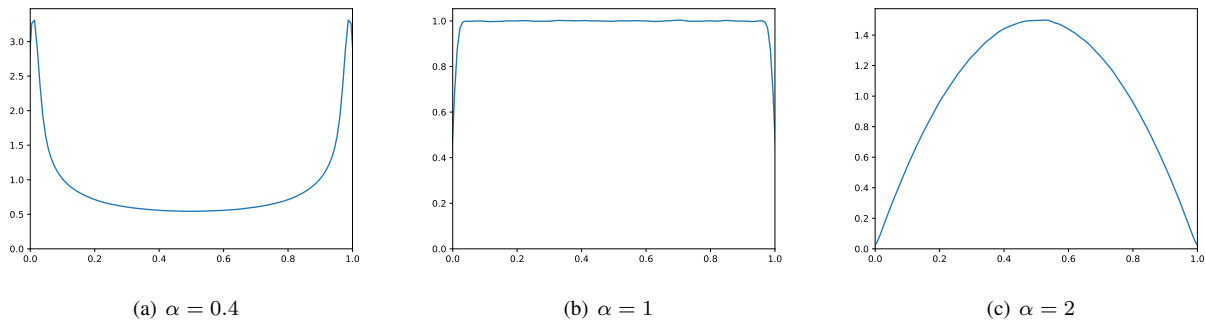


Figure 4. Kernel density estimation of Beta distribution when α varies.

Table 2. Unknown detection performance in terms of AUC (mean \pm std). Results are averaged among five randomized trials. N. R. means the original work did not provide a particular value.

Methods	SVHN	CIFAR10	CIFAR+10	CIFAR+50	Tiny-ImageNet
Softmax	88.6 \pm 1.4	67.7 \pm 3.8	81.6 \pm N. R.	80.5 \pm N. R.	57.7 \pm N. R.
OpenMax [3]	89.4 \pm 1.3	69.5 \pm 4.4	81.7 \pm N. R.	79.6 \pm N. R.	57.6 \pm N. R.
G-OpenMax [6]	89.6 \pm 1.7	67.5 \pm 4.4	82.7 \pm N. R.	81.9 \pm N. R.	58.0 \pm N. R.
OSRCI [22]	91.0 \pm 1.0	69.9 \pm 3.8	83.8 \pm N. R.	82.7 \pm N. R.	58.6 \pm N. R.
C2AE [24]	89.2 \pm 1.3	71.1 \pm 0.8	81.0 \pm 0.5	80.3 \pm 0.0	58.1 \pm 1.9
CROSR [36]	89.9 \pm 1.8	N. R.	N. R.	N. R.	58.9 \pm N. R.
GFROSR [26]	93.5 \pm 1.8	83.1 \pm 3.9	91.5 \pm 0.2	91.3 \pm 0.2	64.7 \pm 1.2
PROSER	94.3 \pm 0.6	89.1 \pm 1.6	96.0 \pm 0.4	95.3 \pm 0.3	69.3 \pm 0.5

between these two methods, and we can infer from Figure 3 that manifold mixup generates instances more similarly than the vanilla method. Besides, the generated space of vanilla mixup is much more than novel space, which also involves known space. As a result, utilizing vanilla mixup may destroy the learned embedding space to some extent.

II. Experiment Implementation

II.1. Unknown Detection Results

In the main paper, we report the averaged AUC of unknown detection tasks. We simulate the sampling process over five trials [22] to report the mean and standard deviation, and the full results are shown in Table 2. We report the baseline performance from [26, 36, 22]. Note that N.R. in the table means that the original paper did not report the standard deviation.

II.2. Implementation Details of PROSER.

We employ the same backbone architecture as [26, 22]. PROSER is trained with SGD with momentum of 0.9, and the initial learning rate is set to 0.001 in the experiment. We fix the batch size to 128 for all datasets. As we discussed in the hyper-parameter part, we set $\beta = 1, \gamma = 0.1$, and the number of classifier placeholders is set to 5 for all datasets. The calibration *bias* is obtained by ensuring 95% validation



Figure 5. Dataset example of CIFAR10 open-set recognition, each line stands for the open-set class in main paper.

data be recognized as known. The α parameter in Beta distribution is set to 2 for all datasets. We conduct the experiment on Nvidia RTX 2080-Ti GPU with Pytorch 1.6.0.

II.3. Dataset Configuration.

We also show the dataset example of CIFAR10 open-set recognition tasks in the main paper in Figure 5. The ‘crop’ datasets is part of the original picture, and ‘resize’ datasets is the full original picture resized into 32*32 pixel. As a result, detecting outliers from ‘crop’ datasets is easier than that of ‘resize’ datasets. This is consistent with macro-F1 results in the main paper.