

# Learn representations in the presence of segmentation label noises

Jihong Ju

Computer Vision Lab, TU Delft

Mekelweg 4, 2628 CD Delft

j.ju@student.tudelft.nl

August 26, 2017

## Abstract

Transferring pre-trained representations is widely adopted by convolutional neural networks for semantic segmentation because the training data is often available only on a small scale. In domains where directly adapting classification networks is challenging, we propose to train representations with segmentation datasets containing mislabeled objects and unsegmented objects. Our experiments demonstrate that both mislabeled segments and incomplete segmentation lower the fine-tuning performance of the learned representations. To get rid of the negative effect of objects label noises, we propose to assign objects of any categories a foreground label instead of the exact object categories. Learning representations by segmenting foreground and background turns out to improve the fine-tuning performance significantly when label noises are dominant in the pre-training data. In the existence of unsegmented objects, a sigmoid loss for the background class is proposed to achieve high recall while keeping the precision better than simply weighting the classes. The proposed class dependent, sigmoid (or softmax) loss achieves both better pre-training performance and better fine-tuning performance than the class-weighted loss in the presence of incomplete segmentation.

## 1 Introduction

The often limited availability of training samples motivates most state-of-the-art deep learning based segmentation models [26, 2, 14] to transfer convolutional neural

network (CNN) models [19, 35, 37, 15] trained on a subset of images from ImageNet. Compared to object recognition tasks, it is tougher to collect a dataset for semantic segmentation on a large scale. The difficulty of obtaining manual segmentations is natural because it costs much more efforts for people to segment than to classify an image. One of the largest segmentation datasets, Microsoft COCO2014 [24], contains 123,287 images of 80 object categories, smaller in the number of images than the IL-SRVC dataset by a factor of 10 approximately. Transferring weights from the pre-trained ImageNet models can provide a segmentation performance boost in the limitation of lacking training samples, as reported in [26] and adopted by [2, 14].

However, it can be challenging to employ the ImageNet CNN models directly for semantic segmentation. For example, the object recognition models pursue features invariance to better capture semantics regardless the variations in objects. The result translation invariant and resolution-reduced features reduce the localization accuracy which is not essential for object recognition but is critical for object segmentation. [41, 2] Besides, the ImageNet models were originally trained with natural images at relatively low resolution. However, images in the domain of interest may be (1) non-natural, such as aerial images and medical images, or (2) have different lighting conditions, or (3) have a higher resolution than the ImageNet ones. Therefore, it can be beneficial to retrain the pre-trained ImageNet segmentation datasets for fine-grained cues about boundaries in the domain.

The segmentation datasets for pre-training representations may contain label noises. The use of the crowd-

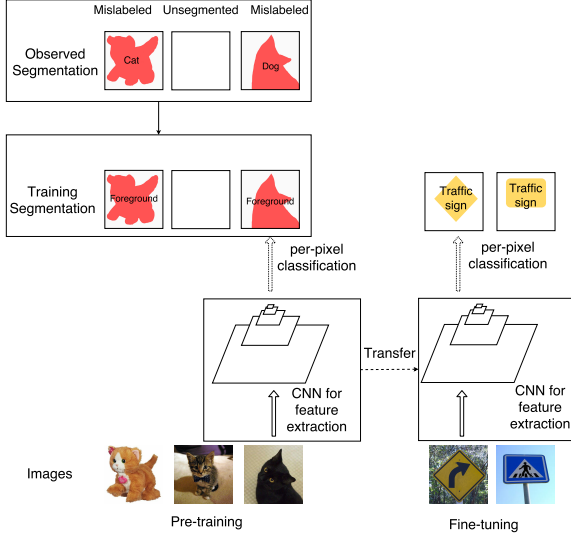


Figure 1: Learning representations with segmentation datasets containing label noises. For segmentations with mislabeled objects and unsegmented objects, we propose (1) to train with only segmentations instead of labeled segmentations, (2) to apply a sigmoid loss for the background class. The learned representations are then used as weights initialization and fine-tuned with a small set of true labels in the domain of interest.

sourcing platform like Mechanical Turk is common nowadays to collect annotations on a large-scale. It is natural for crowd-sourcing workers to make mistakes as a result of lack of expertise, inherent ambiguity of tasks or unconscious bias. Enormous efforts are required, according to [24, 10], to ensure the correctness of segmentations. If not requiring “gold standard” segmentations for training, the efforts saved for correctness can be made to segment more images for a larger dataset. Besides, some labels other than the manual ones may be freely available for particular tasks. These labels often contain structural noises depending on the way they were created. For example, digital maps, like OpenStreetMap, can be used to segment aerial images. Segmentations constructed from the maps could suffer from incompleteness as well as registration problems. [28]

Noises of different kinds can exist in segmentation la-

bels. In particular, we consider three types of label errors occurred to the whole objects instead of individual pixels: inexhaustive segmentation, objects mislabelling, and false segmentation. **Objects mislabelling** from one category to another exist occasionally even for well-annotated datasets. For example, the Microsoft COCO dataset [24] contains some misclassified bears and teddy bears even though annotators were asked to segment only one category at a time; **Inexhaustive segmentation** means that there exist objects left unsegmented. A typical scenario where incomplete segmentation emerges is to segment images containing massive amounts of objects of the same kind, e.g., a flock of sheep or a pile of products; **False positive segmentation** denotes that semantically meaningful objects from an undefined category are wrongly segmented, as objects of interest. It may occur due to the unclear definition of categories, visual similarities between objects, etc. With the simulated datasets described in Section 5.1, we report that objects mislabeling and inexhaustive segmentation both have a negative influence on the learned representations, whereas the false positive segmentation has little effects.

To overcome the negative influence of objects mislabelling, we propose to group all object categories into one foreground class and train representations by learning to segment foreground and background. Jain et al. [18] demonstrated a fully convolutional network trained on over one million images to for binary segmentation generalizes well to thousands of unseen object categories. A convolutional network probably learns generic knowledge about object boundaries if it can segment foreground and background for a wide range of categories sufficiently well. The generic knowledge is likely to transfer well to other categories that are unseen.

If we consider inexhaustive segmentation only, the problem becomes similar to a so-called *positive and unlabelled learning* (PU learning) setup [23]. In the positive and unlabeled learning setup, the training dataset has two sets of examples: a *positive (P) set*, containing only positive examples, and an *unlabeled (U) set*, containing a mix of positive or negative examples. The main characteristic of the U set is no easy way to generate reliable negative labels out of it. The set of background pixels mixed with unsegmented object pixels, in general, fulfills this property. In an incompletely segmented dataset, pixels of the segmented objects form the P set, and the rest pixels con-

struct the U set. Training with a segmentation dataset with incomplete segmentations is therefore similar to a learning problem with only positive examples and unlabeled examples. In this work, we treat the unlabeled set as a set of examples with noisy negative labels and propose to use the sigmoid loss for the negative class, following a branch of previous studies for PU learning as discussed in Section 2.

The main contributions of this work are:

1. We present that both objects mislabelling and incomplete segmentation have a negative influence on learning representations, while false positive segmentations do not.
2. We propose to learn representations by training foreground/background segmentations if the segmentation labels are noisy.
3. We employ a class-dependent sigmoid loss to balance precision and recall more effectively than weighting classes when training CNN models with positive and unlabeled examples.

The rest of this thesis is organized as follows: In the next section, we summarize related works. In Section 3 we formulate the problem of transferring representations learned with segmentations containing label noises, and simulation of the three types label noises. We introduce the class-dependent, sigmoid loss for the negative class for deep learning with positive and unlabeled examples in Section 4. Experiments in Section 5.1 are designed to investigate the influences of objects mislabeling, inexhaustive segmentations, and false positive segmentations independently, and validate whether our proposed methods can alleviate the negative influences. The proposed sigmoid loss is evaluated, compared to weighting classes, in simulated PU learning setups in Section 5.2. Discussions are presented in Section 6 and conclusions are summarized in Section 7.

## 2 Related works

**Transfer Learning** Weights of convolutional neural networks were proved “transferable” not only to another dataset [34, 40] but also to other applications [11, 26]. Transferring weights of pre-trained CNN models allows

improving the model performance without engaging in much efforts spent for data-labeling. [29] Yosinski et al. [40] discovered that the transferability of features is correlated with feature generality, i.e., how much the feature depends on a particular category. They also reported the weights from low-level layers of CNN models are well transferable to even completely dissimilar categories, for example, from natural objects to human-made objects. Because features are transferable regardless the exact categories they are trained with, we argue that binarizing or categorizing the pre-training classes can be expected to have no significant influence to the transferability of the result pre-trained models.

Apart from supervised pre-training, one can also perform unsupervised learning to obtain pre-trained features typically with auto-encoders [39, 27], deep belief networks [16, 21]. Though a few studies [9, 8, 1] discussed the advantage of unsupervised pre-trained features compared to random weights initialization, the distance between the two has been shortened ever since the arises of modern initialization strategies, namely Xavier initialization [12] and its variants. We used random weights initialization as the lower baseline for pre-training with noisy labels. Features learned with supervision in the presence of label noises should at least outperform random weights because noisy information should be still better than no information.

**Deep Learning with Noisy Labels** The impact of randomly flipped labels on classification performance has been investigated by [36, 30] for convolutional neural networks. They both reported decreases in classification performance as the proportion of flipped labels increases for a fixed number of training samples. On the other hand, Rolnick et al. [32] argued that deep neural networks can learn robustly from noisy datasets as long as appropriate choices of hyperparameters were made. They studied the effects of diluting instead of replacing correct labels with noisy labels and concluded that larger training set is of more importance than lower the level of noise. None of these studies explored the influence of label noises on feature transferability. To the best of our knowledge, we are the first research to investigate feature robustness label noises.

Methods, including a linear noise layer on top of the

model output [36], loss correctness with an estimation of the noise transition matrix [30] were proposed to compensate the negative effect on classification performance introduced by flipped labels. To study if the proposed methods can alleviate the influence of noisy labels as expected, both works simulated the random flipped labels from clean labels. We also designed our experiments using stochastically simulated noisy segmentations from perfect segmentation.

**Positive and Unlabeled Learning** Traditional positive and unlabeled learning methods originally proposed for text classification [25, 23] do not extend well to deep learning models. These traditional methods often follow a two-step strategy: first identifying a set of reliable negative samples (RN set) from U set and then iteratively build a set of classifiers with RN set and P set, while updating the RN set with a selected classifier. It would take tremendously longer time to train a couple of deep learning models iteratively than to train a sequence of naïve Bayesian (NB) models or supported vector machines (SVMs). Therefore, it would be unrealistic to train CNN models in this manner.

Alternatively, one can treat all unlabeled examples as negative and simply reweigh positive and negative examples. [22] Under the assumption of randomly labeled positive examples, Elkan & Noto [7] demonstrated that a classifier trained on positive and unlabeled examples predicts probabilities that differ by only a constant factor from the true conditional probabilities of being positive. These two works considered only binary classification but it is possible to extend the weighted logistic loss and train deep neural networks for multiple classes with only positive and unlabeled examples.

The logistic loss grows to infinity as the confidence of wrong prediction increases to 1. This can be a problem of weighting the negative class: the superfluous penalty for confident, positive predictions, i.e., samples far from the decision boundary have a large influence on the final solution. [38] Du et al. [3] illustrated that logistic loss and hinge loss perform worse than ramp loss in the PU classification setting due to their superfluous penalty for confident predictions. se different losses for positive and negative data. The non-convex Ramp loss [4] and a convex double hinge-loss [3] were proposed separately to

learn from positive and unlabeled data by Du et al. But neither of the two losses are continuous, which is problematic for gradient based optimization.

To avoid over-punish confident predictions unnecessarily, Tax & Wang [38] uses class-dependent, non-convex loss to train classifiers to retrieve small set of relevant objects from objects of a large, dominant class; We proposed in Section 4 to use a sigmoid loss[38] for the negative class to alleviate superfluous punishment for confident, positive predictions. Additionally, Reed & Lee [31] proposed a bootstrapping loss to emphasize *perceptual consistency* in training when incomplete and noisy labels exist. It has a similar effect of reducing losses of confident predictions. We modified the hard bootstrapping loss to interpret the prior knowledge that positive labels are reliable.

## 3 Feature transferability with label noises

### 3.1 Problem Formulation

**Semantic Segmentation** A deep learning model for semantic segmentation normally consists of two principal functions: a CNN feature extractor  $g$  that extracts hierarchical feature maps  $F$  from images  $I$ , followed by a classifier  $h$  that generates pixel-by-pixel prediction to fit labels  $S$ . Together they form a segmentation model  $f$  to predict class probabilities for each of the pixels in a given image  $I$ :

$$f(I) = h(g(I))$$

Training is to find an optimal  $f$  from the space of functions which minimizes a loss function  $L$  that measures the distance between  $S$  and  $f(I)$ :

$$f^* = \underset{f}{\operatorname{argmin}} L(S, f(I))$$

**Feature Transferability** A pre-trained feature extractor,  $g_t$ , can be transferred as an initialization of the feature extractor and is expected to result in a better performed  $f^*$  on the fine-tuning test set than a random choice of initialization  $g_0$ . The corresponding fine-tuning performance improvement indicates the transferability of the pre-trained feature extractor. Ideally, a feature extractor

pre-trained with noisy labels would result in a fine-tuned model with equivalent performance as one pre-trained with true labels. The difference in the result fine-tuning performance improvement between the noisy and clean pre-trained feature extractor can tell the influence of label noises on feature transferability.

### 3.2 Label noises simulation

It is challenging to find a dataset with both clean and noisy labels available, so we simulate segmentation noises with correct labels. A straightforward way to simulate noisy labels is to corrupt true labels stochastically for each segment with a corruption model. The corruption model simply describes the probability of the observed labels given the true labels.

For segmentation problems, each pixel (or voxel for 3D segmentation) of a training image has a label assigned to one of the pre-defined categories. Supposing there are  $K$  pre-defined categories, the label of a pixel in the  $i$ -th row and  $j$ -th column from an image of size  $h \times w$ :

$$y_{ij} = \begin{cases} 1 \leq k \leq K, & \text{for foreground pixels} \\ 0, & \text{for background pixels} \end{cases}$$

where  $1 \leq i \leq h, 1 \leq j \leq w$  and  $i, j, k \in \mathbb{Z}^+$ .

**Inexhaustive segmentation** Given the true labels, which pixels belong to the same object is known. Inexhaustive segmentation is simulated by flipping all pixels for an object of the  $k$ -th category from  $k$  to 0 with probability  $p_{0k}$ . The probability for these pixels to have correct labels is then:  $1 - p_{0k}$ .

**Objects Mislabelling** To simulate objects mislabelling, we stochastically convert pixel labels of segment for an object of the  $k$ -th category from  $k$  to  $j$  with probability  $p_{jk}$ , where  $j, k \leq K$  and  $j, k \in \mathbb{Z}^+$ . Probabilities for all possible combinations of  $j$  and  $k$  form a *confusion matrix* in which probabilities in each column sum up to 1.

**False positive segmentations** The presence of false positive segmentations is in two steps: excluding classes except one from the foreground categories, forming the clean labels set, and assigning pixels of the excluded categories with foreground labels with a probability of  $p_{11}$ , constructing the noisy labels set.

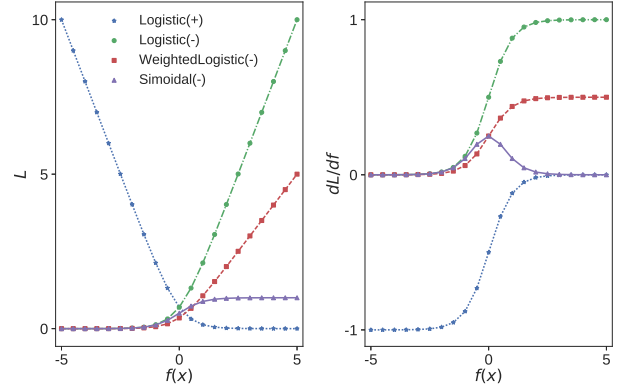


Figure 2: The differences in losses (top figure) and derivatives (bottom figure) with respect to the model output  $f$  between the weighted logistic loss and the sigmoid loss for the negative class. The + sign represents the loss of positive samples and the - sign stands for the loss of negative samples. The sigmoid loss of negative examples reaches a plateau and the derivative drops to zero in the very positive region. The weighted logistic loss for negative is a linearly scaled logistic loss.

## 4 Class-dependent sigmoid loss for PU Learning

**PU Learning** A training set for the positive and negative (PU) learning problems contains only a set of positive examples (P set) and a set of unlabeled samples (U set). Unlabeled examples can be either positive or negative, meaning that there is no reliable negative examples available. One straightforward way to generate negative examples for training is to treat all the unlabeled examples as a set of negative examples with noises. The problem then converts to learning with clean positive labels and noisy negative labels. The goal of solving a PU learning problem is to learn a classifier that predicts as many positives as possible while keeping the false positive rate low, regardless the influence of false negative labels. In other words, the purpose is to achieve high recall without sacrificing too much precision.

**Weighted logistic loss** The mislabeled negative samples bias the classifier to have low recall. It is possible to balance precision and recall by simply weighing the positive and negative examples differently, namely, let the positive and negative examples have different rates of contribution to the total loss. Suppose a logistic loss is used, the corresponding weighted loss for a input-output pair  $(x, y)$  for a classifier determined by parameters  $\theta$  is:

$$l(x, y; \theta) = \begin{cases} -\alpha \log P(y = +1|x; \theta), & y = +1 \\ -\beta \log P(y = -1|x; \theta), & y = -1 \end{cases}$$

where  $\alpha$  and  $\beta$  are weights for positive and negative class respectively, and  $P(y|x; \theta) = \sigma(f(x; \theta))$  is the probabilistic predictions by model  $f(\cdot)$ , activated by the sigmoid function  $\sigma(\cdot)$ . This loss is referred to as the **weighted loss** in the rest of paper. Empirically, the choice of  $p, q$  can be made based on the highest precision and recall achieved on a validation set, or alternatively based on a class priors estimation[5].

**Sigmoid Loss for the negative class** As motivated in Section 2, we used a class-dependent loss to down-weight the loss contribution of very positive predictions negative labels and still making full use of the clean positive labels. The loss of positive examples is still a normal logistic loss and the loss of negative examples is replaced with a sigmoid loss [38]:

$$l(x, y; \theta) = \begin{cases} -\log P(y = +1|x; \theta), & y = +1 \\ 1 - P(y = -1|x; \theta), & y = -1 \end{cases}$$

We called this class-dependent loss **sigmoid loss** in a sense it uses a sigmoid function as the loss of negative samples. Using the logistic loss for the positive class and the sigmoid loss for the negative class is an interpretation of the prior knowledge that the positive labels are reliable whereas the negative labels are not.

Figure 2 shows the differences in losses and derivatives with respect to model output between weighted logistic loss and sigmoid loss. The main feature of sigmoid loss for negative examples is its small changes in the region of confident positive, compared to the weighted loss with  $\alpha = 1$  and  $\beta = 0.5$ . As a consequence, the corresponding derivative decreases to zero as the model prediction increases in the positive direction.

**Hard bootstrapping loss for the negative class** In addition to the proposed sigmoid loss, we also modify the hard bootstrapping loss by Reed et al. [31] for PU learning to set a benchmark. The modified class-dependent hard bootstrapping loss a pair of inputs and label  $(x, y)$  is:

$$l(x, y; \theta) = \begin{cases} -\log P(y = +1|x; \theta), & y = +1 \\ -\beta \log P(y = -1|x; \theta) - (1 - \beta) \log P(y = \hat{y}|x; \theta), & y = -1 \end{cases}$$

where  $\hat{y} = \operatorname{argmax}_{j \in \{-1, +1\}} P(y = j|x)$  is the class with the highest predicted probability and  $0 < \beta < 1$ . The first term of the objective is a weighted logistic loss and the second term can be considered as a regularization term to encourage consistent predictions. This loss is referred as **bootstrapping loss** for the rest of this paper.

**Extending PU learning from classification to segmentation** In Section 1, we argue that learning with unlabeled foreground pixels is similar to a PU learning setup. However, there is still differences between learning with unlabeled foreground pixels and learning with positive and unlabeled examples.

The first difference is that each example in the normal PU learning setup is independent of each other, whereas pixels in images are not. Assuming the probability of mislabeling foreground pixel as the background is independent of its neighbor pixels, the classification losses can be applied to segmentation problems by performing per-pixel classification problems.

Another difference between incomplete segmentation and a normal positive and unlabeled learning problem is that pixels for objects of various categories can be unlabeled. Supposing there are  $K$  categories of interest, varying from class 1 to class  $K$ , the class 0 is for unlabeled data which may or may not belong to the  $K$  defined categories. The sigmoid loss can be extended train deep learning models with unlabeled examples from various categories:

$$l(x, y; \theta) = \begin{cases} -\alpha_1 \log P(y = 1|x; \theta), & y = 1 \\ \vdots & \\ -\alpha_K \log P(y = K|x; \theta), & y = K \\ \beta(1 - P(y = 0|x; \theta)), & y = 0 \end{cases}$$

where  $\alpha_1, \dots, \alpha_K, \beta$  are weighting factors, and  $P(y|x; \theta)$  is the predicted probability for class  $y$ , i.e., model output activated by the softmax function. This loss is referred to as the **softmax loss** because it uses a softmax activation for the model output.

**Implementation details** We introduced the sigmoid loss only after training with a class-weighted cross entropy loss for a few epochs. The sigmoid loss of negative examples saturates for very positive outputs, meaning that the confident, positive prediction has little contribution to the weights update. The wrong confident predictions can introduce problems at the beginning of the training procedure when the confident predictions are likely to be made at random. Optimization would reach the plateau when the model made all positive predictions with high confidence. Besides, we also introduce the modified hard bootstrapping loss only after a few epochs trained with class-weighted loss because it also relies on a nonrandom model for sufficiently reliable prediction  $\hat{y}$ .

Another problem encountered in the PU learning setup is the class imbalance introduced by negatively labeled positive samples. A balanced dataset can become imbalanced in the presence of false negative labels, especially if only a small portion of positive samples are correctly labeled. We reweighed positive and negative samples based on their occurrences of the observed labels to alleviate the influence of imbalance for training. Note that the class-weighted logistic loss reweighed the classes in addition to this frequency balancing class weight.

## 5 Experiments

### 5.1 The influence of label noises in segmentation

To investigate the influence of the label noises, (1) inexact segmentations, (2) objects mislabelling and (3) false positive segmentations, on the learned representations, we set up experiments with simulated label noises from a well-annotated dataset, the PASCAL VOC2011 segmentation dataset [10].

**Dataset** In this experiment, fifteen out of twenty categories of the VOC2011 dataset were selected to form

a *pre-training dataset* and the other categories formed a *fine-tuning dataset*. The pre-training dataset was used to train a Fully Convolutional Network with AlexNet (FCN-AlexNet) model [26] for segmentation in the presence or absence of simulated label noises. Inexhaustive segmentations, objects mislabelling, and false segmentations were simulated independently with stochastic corruptions to the well-annotated pre-training dataset, followed the descriptions in Section 3.1. The fine-tuning dataset was used to fine-tune the weights of convolutional layers from the pre-trained FCN-AlexNet models. To avoid that the choice of pre-training and fine-tuning splitting for categories influence the results, the 20 categories of VOC2011 were divided equally into four folds. The training dataset was enriched with extra segmentations by Har-iharan et al. [13] To keep the segmentation task simple, we used only single-object images, resulting in totally 4000 training images for 20 categories available for pre-training, fine-tuning and evaluation. We subsampled the original images by four times to accelerate the training process.

**Experimental setup** The fine-tuned models were evaluated by mean intersection over union ratio (mean IU) achieved on the fine-tuning test set, referring to as the *fine-tuning performance*. Performance improvement of fine-tuning transferred models compared to a randomly initialized model indicates the transferability of pre-trained weights. The non-transferable layers of FCN-AlexNet were randomly initialized with Xavier Initialization. Experiments run for each fold independently, and the exact partitions of each fold are listed in Table 1. Fully Convolutional Networks with AlexNet was used for experiments because of its relatively small capacity and thus short training time. The pre-trained AlexNet model was used to set a baseline of performance, denoted as the ImageNet model. The ImageNet model and completely random weight initialization were considered as the upper baseline and lower baseline, respectively, for various pre-trained weights summarized in Table 1. The default hyperparameters of FCN-AlexNet in [26] were kept unchanged. The training process run 240,000 iterations for pre-training phase, and 12,000 iterations for fine-tuning phase. Snapshots for trained models were taken every 4,000 iterations. Each experiment was repeated three

times, mean and standard deviation were computed over the last five snapshots for all repetitions.

**Objects mislabelling** The negative influence of random object labels on feature transferability is demonstrated in Table 1. Compared to the model trained with true labels, both models trained with all random labels and half true half random labels are less transferable to the fine-tuning dataset. Transferring the AllRandomLabels model or the HalfRandomLabels model is no better than randomly initializing model weights. The mislabeled objects in segmentations negatively impact the transferability of CNN models in this simulated experiment setup.

We then binarize pre-training classes to foreground and background and achieve a fine-tuning performance better than training with the exact (half) randomized labels and equivalent to training with correct labels. Randomized object labels were mislabeled among foreground classes so that binarizing labels a foreground and background classes can in a sense correct the randomized labels. Compared to the precise but inaccurate noisy labels, binarized labels are accurate but imprecise. This observation indicates that inaccurate labels have a larger influence on the learned representations than imprecise labels.

**Precise labels vs. Accurate labels** To validate that accurate, imprecise labels have less negative effect on feature transferability than inaccurate, precise labels. In figure 3, we decrease the preciseness of pre-training labels by grouping the pre-training classes into a various number of categories to demonstrate the corresponding change of feature transferability. Addition to the binary categorization, we also categorized the fifteen pre-training classes into four meaningful categories: person, animal, vehicle, indoor according to [10]. The result transferred and fine-tuned model is shown as the error bars on the solid line at categories=4. It has almost the same fine-tuning performance as the model trained with binarized labels and the model trained with precision labels (shown as error bar at categories=15).

As a comparison, the fifteen classes were randomly categorized into 4, 7, 11 categories and shown as isolated error bars in figure 3 at categories=4, 7, 11 respectively. Figure 3 reveals that reducing label preciseness by categorizing pre-training classes has little effect on the fine-

tuning performance of transferred models. Even categorizing classes at random without explicit meaning can pre-train weights better than random initialization (shown as error bar at categories=0).

In contrast to the significant influence of randomized precise labels, the imprecise binarized labels have little effect to the learned representations. Binarizing is helpful to learn better representations when mislabeled objects is dominant.

**Inexhaustive segmentation** Fine-tuning performance for pre-trained models with complete segmentations and incomplete segmentations are summarized in Table 1. Pre-training data with half of the objects unsegmented results in fine-tuned models with a worse mean IU (average across four folds) than pre-training with complete labels by 0.04, the same as no pre-training. This observation demonstrates that inexhaustive segmentations have a negative impact on representations transferability.

By applying the class-dependent sigmoid loss to pre-train models with half of the objects segmented, a fine-tuning performance comparable to the model pre-trained with complete segmentation is achieved. The class-dependent sigmoid loss compensates the negative influence of inexhaustive segmentations on the fine-tuning performance of the pre-trained model. Note that we train foreground/background segmentation instead of multi-class segmentation when applied the class-dependent sigmoid loss because binary segmentation can produce as good representation as multi-class segmentation.

**False positive segmentations** We observe a little influence of the incorrectly segmented objects from the non-predefined but meaningful categories, i.e., the false positive segmentations, on the learned representations. To simulate the false positive segmentations errors, we selected one category, either cat or dog depending on the folds, as the target category and all the other 14 categories in the pre-training dataset became non-target, as discussed in Section 3.1. In the presence of false segmentations, instances from non-target categories can be misannotated as the target category with a probability of  $p_{10} = 0.5$ . The result pre-trained models is named as the HalfFalsePositive model in Table 1. The noise-free counterpart of is the model pre-trained with segmentations of the selected



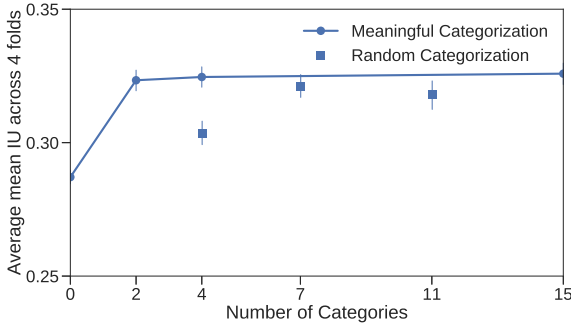


Figure 3: Test performance for the fine-tuned models pre-trained with varying categorization of pre-training classes. Zero categorize means no pre-trained weights used and the model was random initialized. Error bars located on lines denote the meaningful categorization, and isolated error bars denote random categorizations (RC) of the 15 classes. The displayed mean IU/mean accuracies and standard deviations were averaged over four folds listed in Table 1. The line shows that binarizing and categorizing classes meaningfully had little negative effect on feature transferability.

target category only, and the other 14 categories remained unsegmented, denoted as the NoFalsePositive model.

We observe from Table 1 that transferring the Half-FalsePositive model performs almost the same as the No-FalsePositive model. Therefore, false positive segmentation for objects from the 14 categories have a little negative impact on pre-training CNN models.

## 5.2 Learning with only positive and unlabeled samples

### 5.2.1 2D non-linear dataset

To investigate the decision boundaries led by the sigmoid loss for the negative class, we trained a two-layer multilayer perceptron, with six neurons per layer, using the normal logistic loss, the class-weighted logistic loss, and the class-dependent sigmoid loss respectively and visualize the optimal decision boundaries.

**Experimental setup** The training data contains four hundred samples per class drawn randomly from two interleaving half circles with noises added with a minor standard deviation, as shown in Figure 4. Half of the positive examples were assigned negative labels, resulting in a training data with reliable positive labels but noisy negative labels. The three different losses were trained with this noisy training data, and the result decision boundaries are drawn as white regions in Figure 4. The same multilayer perceptron classifier was also trained with true labels to present a baseline decision boundary. The weights for positive class and negative class in the weighted logistic loss were chose to be 1 and 0.5 respectively.

**Results** If trained with the sigmoid loss, the decision boundary is distant from the positive cluster with a relatively large margin as shown in Figure 4. By contrast, the weighted logistic loss results in a decision boundary still closed to the positive examples. For sigmoid loss, the mislabeled positive examples far away from the decision boundary do not contribute more loss than samples less distant from the decision boundary. As a consequence, the loss derivative with respect to the model weights is larger for uncertain predictions in overall than for confident predictions, illustrated by marker sizes in Figure 5. Therefore, training with the sigmoid loss emphasizes the positive predictions with low confidence instead of equally down-weighting all incorrect positive predictions. The emphasis of uncertain predictions by sigmoid loss leads to a margin from the decision boundary to both positive examples and negative examples.

### 5.2.2 CIFAR dataset

To compare the precision and recall achieved by the class-dependent softmax loss and the class-weighted cross entropy, we trained a CNN classifier to distinguish images of multiple relevant categories from non-relevant images when relevant images are partially labeled.

**Experimental setup** We trained an eight-layer CNN model to classify images into eleven classes: ten relevant classes from CIFAR10 and one non-relevant class for all categories from CIFAR100. Relevant images are partially labeled, and the rest forms an unlabeled (U) set

|                            | Initial Feature Extractor      | Fine-tuning mean IU per pretraining-finetuning fold |                                    |   |                                      | Average mean IU across four folds  |
|----------------------------|--------------------------------|---|------------------------------------|---|--------------------------------------|------------------------------------|
| Fine-tuning categories     |                                | aeroplane, bicycle, bird, boat, bottle              | bus, car, cat, chair, cow          | dining table, dog, horse, motorbike, person | potted plant, sheep, sofa, train, TV |                                    |
| Baseline                   | ImageNetModel<br>RandomWeights | $0.42 \pm 0.01$<br>$0.29 \pm 0.01$                  | $0.51 \pm 0.01$<br>$0.29 \pm 0.03$ | $0.49 \pm 0.01$<br>$0.27 \pm 0.01$          | $0.47 \pm 0.01$<br>$0.30 \pm 0.02$   | $0.47 \pm 0.01$<br>$0.29 \pm 0.02$ |
| Objects Mislabelling       | TrueLabels                     | $0.29 \pm 0.01$                                     | $0.36 \pm 0.01$                    | $0.29 \pm 0.01$                             | $0.37 \pm 0.01$                      | <b><math>0.33 \pm 0.01</math></b>  |
|                            | AllRandomLabels                | $0.29 \pm 0.01$                                     | $0.33 \pm 0.03$                    | $0.26 \pm 0.01$                             | $0.28 \pm 0.01$                      | $0.29 \pm 0.01$                    |
|                            | HalfRandomLabels               | $0.27 \pm 0.01$                                     | $0.33 \pm 0.02$                    | $0.25 \pm 0.01$                             | $0.29 \pm 0.01$                      | $0.29 \pm 0.01$                    |
|                            | BinarizedLabels                | $0.30 \pm 0.02$                                     | $0.35 \pm 0.01$                    | $0.29 \pm 0.02$                             | $0.35 \pm 0.03$                      | <b><math>0.32 \pm 0.02</math></b>  |
| Inexhaustive segmentation  | CompleteLabels                 | $0.29 \pm 0.01$                                     | $0.36 \pm 0.01$                    | $0.29 \pm 0.01$                             | $0.37 \pm 0.01$                      | <b><math>0.33 \pm 0.01</math></b>  |
|                            | HalfUnsegmented                | $0.26 \pm 0.01$                                     | $0.30 \pm 0.03$                    | $0.28 \pm 0.03$                             | $0.32 \pm 0.02$                      | $0.29 \pm 0.02$                    |
|                            | SigmoidLoss                    | $0.30 \pm 0.01$                                     | $0.37 \pm 0.01$                    | $0.31 \pm 0.02$                             | $0.34 \pm 0.02$                      | <b><math>0.33 \pm 0.02</math></b>  |
| False positive segmentaion | NoFalsePositive                | $0.26 \pm 0.01$                                     | $0.37 \pm 0.03$                    | $0.27 \pm 0.01$                             | $0.33 \pm 0.04$                      | <b><math>0.31 \pm 0.02</math></b>  |
|                            | HalfFalsePositive              | $0.27 \pm 0.01$                                     | $0.34 \pm 0.01$                    | $0.30 \pm 0.01$                             | $0.32 \pm 0.01$                      | <b><math>0.31 \pm 0.01</math></b>  |

Table 1: Segmentation performance for fine-tuned FCN-AlexNet models pre-trained on 15 categories from the Pascal VOC2011 dataset and fine-tuned on the other 5 categories. **ImageNetModel** represents the pre-trained ImageNet model; **RandomWeights** indicates that the randomly initialized weights; All the other extractors were pre-trained in the presence/absence of the corresponding label noises listed in the leftmost column. Half of the objects unsegmented (**HalfUnsegmented**) result in pre-trained models not better than random weight initialization. Introducing the sigmoid negative loss in the pre-training phase was able to improve the fine-tuning performance to be comparable to pre-trained model with complete segmentation (**CompleteLabels**). Using random labels (**RandomLabels**) decreased the fine-tuning performance of transferred models, compared to using true labels (**TrueLabels**). Binarizing the pre-training classes as foreground and background help overcome the negative effects of random labels. Applying the sigmoid loss to the negative class when pre-trained with inexhaustive segmentations achieves a comparable fine-tuning performance to pre-training with the complete segmentations. Including the false positive segmentations in pre-training (**HalfFalsePositive**) achieves the same fine-tuning performance as not including the false positive segmentations (**NoFalsePositive**), better than random initialization.

together with the non-relevant images. An eight layer CNN model was trained with the cross-entropy loss, the class-weighted cross-entropy loss and the class-depend sigmoid loss respectively in the simulated PU learning setup, where 50% of the positive examples were unlabeled. The CNN model was also trained with the modified hard bootstrapping loss introduced in Section 4 to set a benchmark for the state-of-the-art method to learn in the presence of label noises. Model performances were evaluated on a separate test set with true labels. The architecture of the CNN model can be found in Appendix C. Each model was trained from scratch with Adam optimizer and base learning rate 0.0001. Experiments were

repeated three times with random split of P set, and U set and standard deviations were around 0.01 if not explicitly mentioned.

**Results** Table 2 shows using the softmax loss for the non-relevant class achieves better recall than the class-weighted cross-entropy loss without lowering precision significantly. With 50% of the relevant examples correctly labeled and the rest assigned non-relevant labels, the normal cross-entropy loss leads to an imbalanced model with high precision but low recall, and therefore with a low f1-score. By reweighing the loss for the non-relevant class by a factor of 0.5, the model becomes balanced for pre-

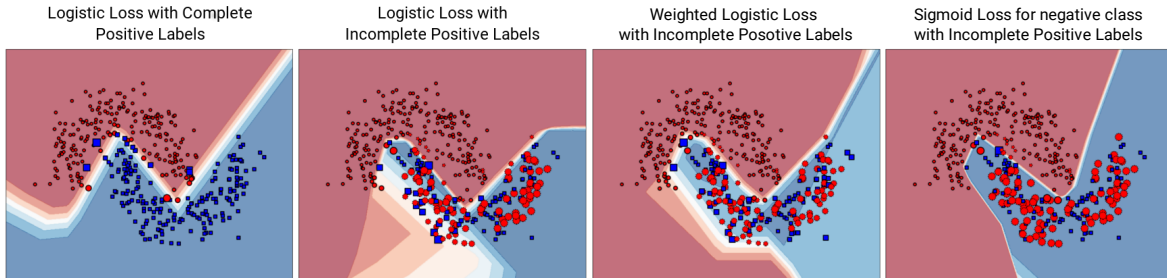


Figure 4: Decision boundaries of a 2-layer multilayer perceptron trained with different losses on a 2D moons dataset with the unlabeled positive. A **red circle** indicates an example labeled as positive and a **blue square** indicates the example has a negative label. The **background colors** represent the classifier prediction in the corresponding area: **red** for negative class, **blue** for positive class and **white** for the class transition areas, i.e., decision boundaries. The **markers sizes** demonstrates the training loss normalized per-class. Compared to the normal logistic loss and weighted logistic loss (positive:negative=1:0.5), the decision boundary optimized with the sigmoid loss has a larger margin from the positive and negative clusters. (Best viewed in color)

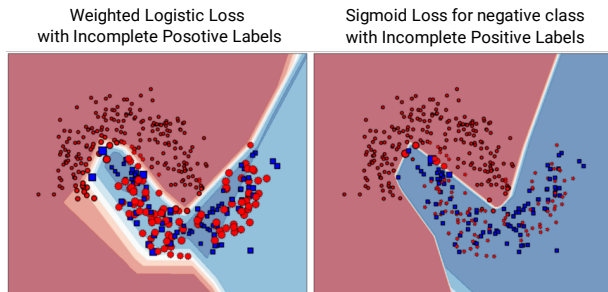


Figure 5: Derivatives w.r.t the last layer output for the two losses (normalized per class and shown as the marker size). The sigmoid loss has small derivatives for predictions farther from the decision boundary. (Best viewed in color)

cision and recall so that the result f1-score is improved significantly. Compared to the class-weighted cross entropy, the class-dependent softmax loss improves recall by 0.08 while reduces precision only by 0.01. The f1-score achieved by the class-dependent softmax loss is slightly better than the class-weighted loss, though not as good as training with clean labels with either 50% of the sample or the complete training set. The state-of-the-art benchmark method, the hard bootstrap loss, achieves the same

f1-score as the softmax loss but not as high recall. This result is as expected because the softmax loss and the hard bootstrap loss share an idea of encouraging confident predictions.

By varying the labeled percentage of the relevant images, Figure 6 demonstrates that the class-dependent softmax loss performs slightly better than the class weighted cross-entropy when the labeled percentage of relevant images is neither too high ( $> 0.8$ ) nor too low ( $< 0.2$ ).

### 5.2.3 PASCAL VOC2011 Segmentation

To compare the class dependent sigmoid loss with the normal logistic loss, and class-weighted logistic loss for training foreground/background segmentation with incomplete segmentations, we used again the PASCAL VOC2011 dataset with extra segmentation [13].

**Experimental setup** We simulated inexhaustive segmentations the same way as described in Section 5.1. The same AlexNet-FCN model was trained together with the different loss functions to predict binary segmentation, determining whether a pixel belongs to an object or not. The same hyperparameters for optimization were used as in Section 5.1. The trained models were evaluated with the test set of PASCAL VOC2011 segmentation dataset with binary segmentations.

| Annotation | Loss          | acc.        | mean prec. | mean rec.   | mean $F_1$  |
|------------|---------------|-------------|------------|-------------|-------------|
| R+N        | CrossEntropy  | 0.87        | 0.88       | 0.82        | 0.85        |
| 50%(R+N)   | CrossEntropy  | 0.83        | 0.84       | 0.78        | 0.80        |
| 50%R+U     | CrossEntropy  | 0.66        | 0.94       | 0.38        | 0.49        |
| 50%R+U     | ClassWeighted | 0.78        | 0.75       | 0.75        | 0.76        |
| 50%R+U     | SoftmaxLoss   | 0.79        | 0.74       | <b>0.83</b> | <b>0.78</b> |
| 50%R+U     | BootstrapHard | <b>0.80</b> | 0.76       | 0.81        | <b>0.78</b> |

Table 2: Comparing different losses for training a 2-layer multilayer perceptron to classify ten relevant classes and one non-relevant class with partially labeled relevant examples and unlabeled non-relevant examples. The trained classifiers are evaluated on a test set of true labels. For each of the relevant classes, precision, recall, and f1-score are measured with the one-vs-all strategy and averaged. **R+N** denotes model trained with the complete relevant labels (R set) and non-relevant labels (N set); **50%(R+N)** represents model trained with the half of the relevant labels and non-relevant labels respectively; **50%R+U** means the model is trained with half of the relevant samples, and the rest relevant samples are mixed with non-relevant samples (U set). Weighting the non-relevant class by a factor of 0.5 improves the mean f1-score from 0.49 to 0.76. Both the class-dependent softmax loss and the hard bootstrapping loss perform even better than simply weighing the classes, but not as good as training with a set of labeled negative examples. Using the softmax loss for non-relevant class achieves the highest mean recall, whereas the modified hard bootstrapping loss has a higher accuracy with 50% positive examples unlabeled.

**Results** As shown in Table 3, the class dependent sigmoid loss achieves the highest mean accuracy, approximately 0.07 better than training with the normal logistic loss, and 0.04 better than the class-weighted loss. In contrast to the improvement of mean accuracy, improvement of mean IU for both the class-weighted loss and the class-dependent logistic loss are insignificant. The increase in the mean accuracy is caused by an increase in foreground accuracy and a decrease in background accuracy. The decrease in background accuracy interferes the improvement mean IU since the mean IU counts for both low false positive rate and low false negative rate.

Selective predictions made by the models trained with the sigmoid negative loss and the cross entropy loss were presented in Figure 7. For the two example images shown, the model trained with the cross entropy loss failed to

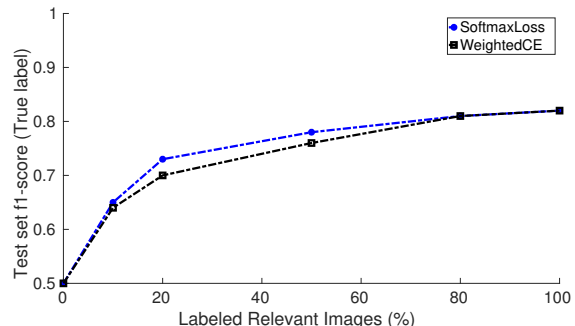


Figure 6: Comparing f1-score for the class-dependent softmax loss and the class-weighted cross entropy with varying percentage of relevant images labeled. The class-dependent softmax loss achieves better test f1-score than the class-weighted cross entropy when 20% and 50% percentage of relevant images are labeled, and the others are mixed with non-relevant images.

segment objects from images whereas sigmoid negative loss predicted segmentations on the position of the objects. The coarse outlines were mainly due to the limited compacity of the FCN-AlexNet model. The third column shows predictions given by model trained with complete training segmentation, and it did not produce more accurate outlines. There is no example observed correctly segmented by the model with the class-weighted logistic loss but not by the model with the class-dependent sigmoid loss.

## 6 Discussion

We investigate in this work the influence of label noises but not segmentation noises for a noisy segmentation dataset. In practice, noisy segmentation dataset may also contain segmentation errors such as imprecise boundaries, over-segmenting and under segmenting the objects. Future studies are necessary to learn representations in the presence of segmentation noises.

When we simulate the label noises in our experiments, we assume unsegmented objects and mislabeled objects occur independently to the shape and appearance of objects. In practice, some categories may have a higher probability to be mislabeled due to its ambiguity in shapes

| Annotation | Loss          | overall acc. | mean acc.   | f.w. IU | mean IU     |
|------------|---------------|--------------|-------------|---------|-------------|
| Complete   | LogisticLoss  | 0.90         | 0.85        | 0.82    | 0.75        |
| 50%Unseg.  | LogisticLoss  | 0.85         | 0.68        | 0.73    | 0.60        |
| 50%Unseg.  | ClassWeighted | 0.84         | 0.71        | 0.73    | <b>0.62</b> |
| 50%Unseg.  | SigmoidLoss   | 0.83         | <b>0.75</b> | 0.72    | <b>0.62</b> |

Table 3: Training foreground/background segmentation with different losses when 50% of the objects are unsegmented. The performances are achieved on the test set of PASCAL VOC2011 segmentation dataset. Both the class-dependent sigmoid loss and the class-weighted logistic loss perform better than the normal logistic loss when 50% objects unsegmented but not as good as the model trained with complete segmentations. The class-dependent sigmoid loss has a better mean accuracy than the class-weighted logistic loss and a similar mean IU as the class-weighted logistic loss.

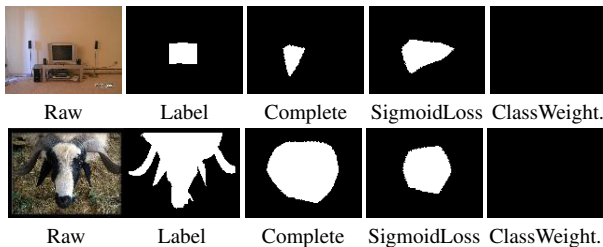


Figure 7: Example predictions made by models trained with the logistic loss and the class-dependent sigmoid loss. This figure presented two selective images for which the model trained with the logistic loss failed to segment objects, whereas the model trained with the class-dependent sigmoid negative loss succeed.

or appearances, such as bear and teddy bear. Similarly, ambiguous objects are more likely to be missing than easily recognizable objects. Experiments on a real dataset with both clean and noisy segmentation are valuable to validate that our findings.

We proposed to learn representations with binary segmentation when objects mislabelling dominates the pre-training dataset. Training binary segmentations can also be relevant when there are multiple segmentation datasets for pre-training, and category overlappings and affinities prevent from combining multiple datasets into one.

We argued to not over-punish confident, positive predictions when negative labels are noisy so that we come

up with the sigmoid loss for negative examples. But with the sigmoid loss, we are not able to determine what is the threshold of confident predictions. A generalized logistic function may be used to replace the normal logistic function as the activation function to achieve a more flexible S-shape and the tuning where the loss saturates. For example, a parametrized sigmoid loss for the negative class could be  $l_- = \alpha(\frac{1}{1+\exp(\beta z)})^\gamma$ , where  $z$  is the model output for the negative class,  $\alpha$  is the scale factor,  $\gamma$  affects where the loss starts,  $\beta$  determines where the loss saturates.

Besides, saturating loss for confident predictions has an effect of encouraging confident predictions. It is intuitively similar to the minimum entropy regularization for semi-supervised learning, which also encourages confident predictions. Both the class-dependent sigmoid loss for PU learning and the minimum regularization scheme favors low-density separations of input features. The hard bootstrapping loss by [31] has a “soft” alternative which is equivalent to the minimum entropy regularization. This similar in the sigmoid loss and the bootstrapping loss explains why the sigmoid loss has a similar performance as the hard bootstrapping loss.

Not over-punishing confident predictions have also its disadvantages: (1) If a classifier makes incorrect predictions with high confidence, it tends to keep being wrong for these examples and emphasize predictions by itself. (2) Punishing confident predictions more than uncertain predictions with the logistic loss is a design of choice for neural networks to optimize more effectively, whereas the sigmoid loss breaks it. These factors determine that the sigmoid loss often performs worse than the logistic loss when the dataset contains only correct labels or a few noisy labels. There is a trade-off to make between punishing and not punishing more for confident predictions, based on the prior knowledge: an estimation of the noisy negative labels percentage.

Lastly, applying the sigmoid loss to segmentation data with incomplete segmentations assumes that the pixels for objects are unlabeled independently. In practice, there is often a spatial dependence for noises in pixel labels which can be valuable to improve performance. For example, if one or a few pixels have a high probability of being foreground, their neighboring pixels are probably also foreground.

## 7 Conclusion

We investigate in this paper to learn representation by training with segmentation datasets containing label noises. Specifically, we report both objects mislabelling and inexhaustive segmentations negatively influence the transferability of learned representation. By contrast, false positive segmentations do not reduce the fine-tuned performance of learned representation as the other two types noises do. We present that binarizing classes as foreground and background slow the decrease of the fine-tuning performance of the learned representations due to the mislabeled objects. Incomplete segmentation causes the trained model to make predictions with a low recall. To overcome the influence of incomplete segmentations, we propose a class-dependent sigmoid loss to not over-punish the confident, positive predictions for samples assigned negative labels. Compared to simply reweighing classes differently, the proposed sigmoid loss for the negative class achieves higher recall while not sacrificing precision by much. Applying the sigmoid loss to the segmentation model pre-training improves both the pre-training and the fine-tuning performance for a pre-training dataset with simulated incomplete segmentations.

## References

- [1] Yoshua Bengio. Deep learning of representations for unsupervised and transfer learning. In *Proceedings of ICML Workshop on Unsupervised and Transfer Learning*, pages 17–36, 2012.
- [2] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *arXiv preprint arXiv:1606.00915*, 2016.
- [3] Marthinus Du Plessis, Gang Niu, and Masashi Sugiyama. Convex formulation for learning from positive and unlabeled data. In *International Conference on Machine Learning*, pages 1386–1394, 2015.
- [4] Marthinus C du Plessis, Gang Niu, and Masashi Sugiyama. Analysis of learning from positive and unlabeled data. In *Advances in neural information processing systems*, pages 703–711, 2014.
- [5] Marthinus Christoffel Du Plessis and Masashi Sugiyama. Class prior estimation from positive and unlabeled data. *IEICE TRANSACTIONS on Information and Systems*, 97(5):1358–1362, 2014.
- [6] Vincent Dumoulin and Francesco Visin. A guide to convolution arithmetic for deep learning. *arXiv preprint arXiv:1603.07285*, 2016.
- [7] Charles Elkan and Keith Noto. Learning classifiers from only positive and unlabeled data. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 213–220. ACM, 2008.
- [8] Dumitru Erhan, Yoshua Bengio, Aaron Courville, Pierre-Antoine Manzagol, Pascal Vincent, and Samy Bengio. Why does unsupervised pre-training help deep learning? *Journal of Machine Learning Research*, 11(Feb):625–660, 2010.
- [9] Dumitru Erhan, Pierre-Antoine Manzagol, Yoshua Bengio, Samy Bengio, and Pascal Vincent. The difficulty of training deep architectures and the effect of unsupervised pre-training. In *Artificial Intelligence and Statistics*, pages 153–160, 2009.
- [10] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015.
- [11] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014.
- [12] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 249–256, 2010.
- [13] Bharath Hariharan, Pablo Arbeláez, Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. Semantic contours from inverse detectors. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, pages 991–998. IEEE, 2011.
- [14] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. *arXiv preprint arXiv:1703.06870*, 2017.
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

- [16] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006.
- [17] David H Hubel and Torsten N Wiesel. Receptive fields, binocular interaction and functional architecture in the cat’s visual cortex. *The Journal of physiology*, 160(1):106–154, 1962.
- [18] Suyog Dutt Jain, Bo Xiong, and Kristen Grauman. Pixel objectness. *arXiv preprint arXiv:1701.05349*, 2017.
- [19] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [20] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [21] Honglak Lee, Roger Grosse, Rajesh Ranganath, and Andrew Y Ng. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In *Proceedings of the 26th annual international conference on machine learning*, pages 609–616. ACM, 2009.
- [22] Wee Sun Lee and Bing Liu. Learning with positive and unlabeled examples using weighted logistic regression. In *ICML*, volume 3, pages 448–455, 2003.
- [23] Xiao-Li Li and Bing Liu. Learning from positive and unlabeled examples with different data distributions. *Machine Learning: ECML 2005*, pages 218–229, 2005.
- [24] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [25] Bing Liu, Yang Dai, Xiaoli Li, Wee Sun Lee, and Philip S Yu. Building text classifiers using positive and unlabeled examples. In *Data Mining, 2003. ICDM 2003. Third IEEE International Conference on*, pages 179–186. IEEE, 2003.
- [26] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015.
- [27] Jonathan Masci, Ueli Meier, Dan Cireşan, and Jürgen Schmidhuber. Stacked convolutional auto-encoders for hierarchical feature extraction. *Artificial Neural Networks and Machine Learning–ICANN 2011*, pages 52–59, 2011.
- [28] Volodymyr Mnih and Geoffrey E Hinton. Learning to label aerial images from noisy data. In *Proceedings of the 29th International Conference on Machine Learning (ICML-12)*, pages 567–574, 2012.
- [29] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2010.
- [30] Giorgio Patrini, Alessandro Rozza, Aditya Menon, Richard Nock, and Lizhen Qu. Making neural networks robust to label noise: a loss correction approach. *arXiv preprint arXiv:1609.03683*, 2016.
- [31] Scott Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and Andrew Rabinovich. Training deep neural networks on noisy labels with bootstrapping. *arXiv preprint arXiv:1412.6596*, 2014.
- [32] David Rolnick, Andreas Veit, Serge Belongie, and Nir Shavit. Deep learning is robust to massive label noise. *arXiv preprint arXiv:1705.10694*, 2017.
- [33] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.
- [34] Hoo-Chang Shin, Holger R Roth, Mingchen Gao, Le Lu, Ziyue Xu, Isabella Nogues, Jianhua Yao, Daniel Mollura, and Ronald M Summers. Deep convolutional neural networks for computer-aided detection: Cnn architectures, dataset characteristics and transfer learning. *IEEE transactions on medical imaging*, 35(5):1285–1298, 2016.
- [35] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [36] Sainbayar Sukhbaatar, Joan Bruna, Manohar Paluri, Lubomir Bourdev, and Rob Fergus. Training convolutional networks with noisy labels. *arXiv preprint arXiv:1406.2080*, 2014.
- [37] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.
- [38] David MJ Tax and Feng Wang. Class-dependent, non-convex losses to optimize precision. In *Pattern Recognition (ICPR), 2016 23rd International Conference on*, pages 3314–3319. IEEE, 2016.

- [39] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, and Pierre-Antoine Manzagol. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of Machine Learning Research*, 11(Dec):3371–3408, 2010.
- [40] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In *Advances in neural information processing systems*, pages 3320–3328, 2014.
- [41] Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Dalong Du, Chang Huang, and Philip HS Torr. Conditional random fields as recurrent neural networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1529–1537, 2015.



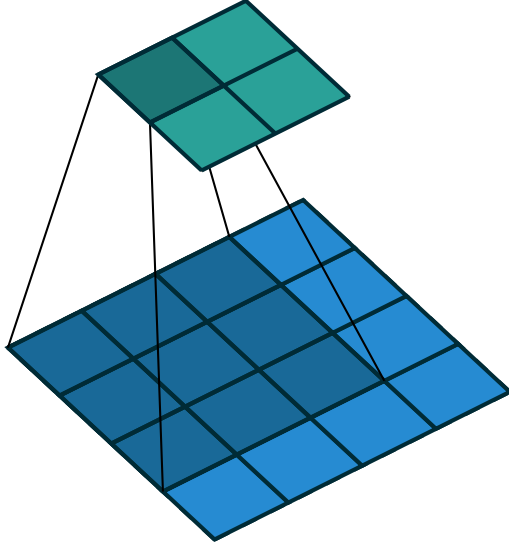


Figure 8: A basic convolution operation. A 3x3 convolutional filter convolves with a 3x3 window sliding over the image (bottom). The output of convolutions at each sliding position form a feature map (top). This figure was drawn by Dumoulin and Visin [6]

## A Convolutional Networks for Semantic Segmentation

### A.1 Convolutional Neural Networks

The main components of a typical convolutional neural network (CNN) are several layers of convolutions and sub-sampling, followed by a few fully-connected layers.

An example CNN model, LeNet-5 (1998) [20], is shown in Figure 9. The first convolutional layer of LeNet contains 6 convolutional kernels of size 5x5 and each convolutional kernel convolve with small windows sliding over the images and produce a feature map of size 28x28. Each output in the produced feature map is corresponding to a small sub-region of the visual field (the image), called a *receptive field*. A following max pooling layer subsamples the feature maps by a factor of two by extracting the maximum values for every two adjacent pixels literally

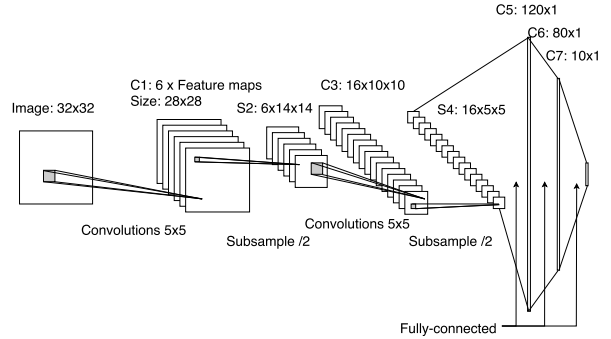


Figure 9: An example convolutional neural network, LeNet-5 [20]

and vertically. The result feature map S2 has a shape of 14 by 14 and a receptive field of 6 by 6. Another sequence of convolutional and pooling layers generate feature maps of size 5x5 with receptive field 16x16. Neurons in the last three layers of LeNet are fully connected to the layer before and the layer after if exists, creating the final prediction for 10 classes.

Features produced by CNN models have a rich hierarchy varying from local to global, from simple to complex. The bottom layers in the convolutional layer stack have smaller receptive fields and the top layers have larger receptive fields. A small receptive field means that the filter have access to information only in a local sub-region of the image while a large receptive field can convey more global information.

The various pattern responses from local to global, from simple to complex for stacked convolutional layers is a reflect of emulating animals visual cortex. In cat's visual cortex [17], two basic cell types of visual cortex have been identified: Simple cells respond maximally to specific edge-like patterns within their receptive field. Complex cells have larger receptive fields and are locally invariant to the exact position of the pattern. The shallower convolutional layers play a similar functionality as simple cells while the deeper layers maps are similar to complex cells.

The main benefit of CNN compared to a standard multilayer neural network (multilayer perceptron) is that 1. take advantage of the 2D structure of an input image 2. it is easier to optimize because of spatial weights shar-

ing and local connectivity pattern of convolutional layers. Convolutional neurons and maximum pooling, translation invariance as well as scaling invariance and distortion invariance to some extent are achievable for convolutional neural networks. [20] Different from the traditional handcrafted features, learnable convolutional features normally generalize well and can achieve better performance for dataset with a complex input distribution. [19] By increasing the number of convolution layers and number of filters in each layer, one can create CNN models with high capacity, meaning a large space of representable functions. This can be beneficial for datasets of immense complexity, for example, ILSVRC [33], Microsoft COCO [24], as long as there are sufficient training samples with an appropriate optimization strategy.

## A.2 Semantic image segmentation

Semantic image segmentation is to segment images into semantically meaningful partitions, a.k.a., *segments*. It can be operated as classifying pixels into the corresponding pre-defined categories.

CNN models on object classification tasks can be adapted to perform semantic image segmentation tasks. [26] One of the primary challenges of applying CNN model to segmentation tasks is how to combine global information and local information to solve semantics and localization altogether. In contrast to object classification tasks, which normally only need global information to resolve semantics, segmentation tasks also require local information to resolve locations.

Long et al. [26] proposed a so-called skip architecture in the Fully convolutional networks (FCN) to aggregate information from the local low-level features in the hierarchy with global information from the high-level features. As we discussed in the previous session, convolutional layers can extract hierarchical features, varying from low-level to high-level encode information from local to global. The low-level features are fine, presenting appearances and the high-level features are coarse, revealing semantics. By combining them together, it becomes possible to create accurate and detailed segmentation.

Convolutional layers in FCN for feature extractions (solid arrows in Figure 10) can be transferred from ImageNet models.

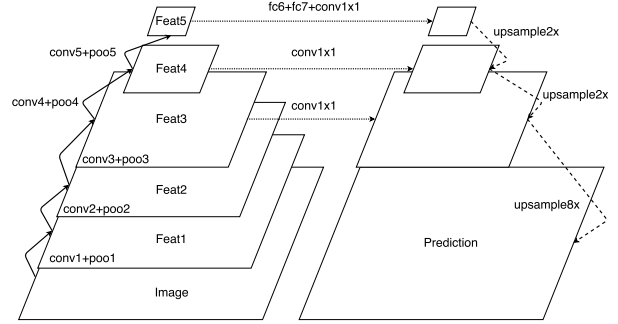


Figure 10: Fully convolutional network (FCN) by Long et al. (2015) [26]. Features of different resolutions are stacked in a feature pyramid on the left-hand side, with the image at the bottom of the pyramid. Predictions of different resolutions are piled in a prediction pyramid on the right-hand side. Each solid arrow denotes a few convolutional layers followed by a max pooling layer; Dotted arrow represents convolutional layers with kernel size one by one; Dashed arrows are up-sampling layers or transposed convolutional layers. From top to bottom, each level of prediction is upsampled and merged with the prediction under it. The bottom prediction is output as the final prediction, which has the same size as the image.

## B Cost function and Optimization

### B.1 Cross entropy loss

The cross entropy loss (a.k.a. softmax loss) is one of the most commonly used cost function for convolutional neural networks in classification problems. Let  $x^{(i)}$  be an input example from totally  $m$  examples,  $y^{(i)} \in 0, \dots, K$  be the corresponding label, and  $\theta$  be parameters of model  $f(\cdot)$ . The cross entropy loss is defined as:

$$J(\theta) = - \sum_{i=1}^m \sum_{k=0}^K 1\{y^{(i)} = k\} \log P(y^{(i)} = k | x^{(i)}; \theta)$$

In the equation above,  $1\{\cdot\}$  is the “indicator function” defined as:

$$1\{\text{statement}\} = \begin{cases} 1, & \text{statement is true} \\ 0, & \text{otherwise} \end{cases}$$

$P(y^{(i)} = k|x^{(i)}; \theta) = \sigma(f(x; \theta))_k$  is the likelihood of  $y^{(i)}$  being  $k$ , predicted by model  $f(\cdot)$ , where  $\sigma(\cdot)_k$  is the softmax function that applies to model output for the  $k$ -th class.

Model outputs  $f(x^{(i)}; \theta)$  is a vector of  $k$  elements with values varying from negative infinity to positive infinity. Each element of the output vector is corresponding to one class out of  $K$  classes. A larger output value for one class,  $k$ , than another,  $j$ , means that the example  $x^{(i)}$  is more likely to be class  $k$  than class  $j$ . The softmax function ensures that the model outputs are normalized to a region between 0 and 1, and sum up to 1 for all classes so that the result outputs fullfils a probability distribution over  $K$  different possible outcomes.

The cross entropy loss is a form of negative log-likelihood. The loss is closed to zero if the predicted probability of  $y^{(i)}$  is large, and takes a large positive value if the probability is small. Minimizing the negative log likelihood of the correct class can be interpreted as performing Maximum Likelihood Estimation (MLE), a commonly optimization.

## B.2 Gradient based optimization and Back-propagation

The model is optimized by solving the optimal  $\theta$  that minimizes the loss function. It is impossible to solve  $\theta$  for a non-linear model analytically so that a gradient-based optimization can be used as an efficient alternative.

The derivative of the cross entropy loss with respect to the  $k$ -th parameter of the last layer  $\theta_k^{(L)}$  is:

$$\nabla_{\theta_k^{(L)}} J(\theta) = - \sum_{i=1}^m [z^{(i)} (1\{y^{(i)} = k\} - P(y^{(i)} = k|x^{(i)}; \theta))]$$

where the superscription ( $L$ ) of  $\theta$  denotes the layer number of the last layer, and  $z^{(i)}$  is the output of the last layer for the  $i$ -th example.

Weights of the last layer in the  $t + 1$ -th iteration is updated by:

$$\theta_{t+1}^{(L)} = \theta_t^{(L)} - \alpha \nabla_{(\theta^{(L)})} J(\theta)$$

where  $\alpha$  is the learning rate determining how quickly the weights are updated.

Gradients in the layers  $l < L$  are calculated via a so-called back propagation of errors. The error of  $l$ -th layer

| layer name | output size    | 8-layer                                  |
|------------|----------------|--|
| conv1      | $16 \times 16$ | $3 \times 3, 32, \text{LeakyReLU}(0.2)$  |
|            |                | $3 \times 3, 32, \text{LeakyReLU}(0.2)$  |
|            |                | $2 \times 2$ max pool, dropout(0.2)      |
| conv2      | $8 \times 8$   | $3 \times 3, 64, \text{LeakyReLU}(0.2)$  |
|            |                | $3 \times 3, 64, \text{LeakyReLU}(0.2)$  |
|            |                | $2 \times 2$ max pool, dropout(0.2)      |
| conv3      | $4 \times 4$   | $3 \times 3, 128, \text{LeakyReLU}(0.2)$ |
|            |                | $3 \times 3, 128, \text{LeakyReLU}(0.2)$ |
|            |                | $2 \times 2$ max pool, dropout(0.2)      |
| fc         | $1 \times 1$   | flatten, 512-d fc, ReLU, dropout(0.5)    |
|            |                | 11-d fc, softmax                         |
| Parameters |                | 1,341,739                                |

Table 4: 8-layer Convolutional Neural Networks used in the classification of the CIFAR dataset.

is propagated the layer after  $l + 1$ :

$$\delta^{(l)} = \left( (\theta^{(l)})^T \delta^{(l+1)} \right) \bullet f'(z^{(l)})$$

where  $f'(z^{(l)})$  is the derivative of the activation function. Gradients with respect to weights for the  $l$ -th layer is:

$$\nabla_{\theta^{(l)}} J(\theta) = \delta^{(l+1)} (a^{(l)})^T$$

The weights update for the  $l$ -th layer in the  $t + 1$  is computed similarly as the last layer by:

$$\theta_{t+1}^{(l)} = \theta_t^{(l)} - \alpha \nabla_{(\theta^{(l)})} J(\theta)$$

## C Additional information

### C.1 An 8-layer Convolutional neural network

The architecture of the 8-layer convolutional neural network used in the classification of CIFAR dataset is presented in Table 4.

### C.2 Evaluation metrics

$$\text{accuracy} = \frac{\text{true pos.} + \text{true neg.}}{\text{true pos.} + \text{false pos.} + \text{true neg.} + \text{false neg.}}$$

$$\text{precision} = \frac{\text{true pos.}}{\text{true pos.} + \text{false pos.}}$$

$$\text{recall} = \frac{\text{true pos.}}{\text{true pos.} + \text{false neg.}}$$

$$F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

$$\text{IU} = \frac{\text{true pos.}}{\text{true pos.} + \text{false pos.} + \text{false neg.}}$$