

Learn transferable features for semantic image segmentation in the presence of label noise

First Author
Institution1
Institution1 address
firstauthor@i1.org

Second Author
Institution2
First line of institution2 address
secondauthor@i2.org

August 14, 2017

Abstract

The ABSTRACT HERE

1 Introduction

Why noisy labels? This paragraph should discuss the ubiquity of label noise and difficulties of collecting perfect annotations.

The recent success of deep neural networks benefits from the availability of large-scale supervised datasets such as [26, 5, 21, 17]. These datasets allow researchers to develop deep neural network models for object recognition[28], object detection[7], semantic image segmentation[18, 35] and other applications assuming the existence of perfect ground-truth segmentation. However, collecting well-annotated datasets on a large scale can be tremendously expensive and time-consuming in general. When annotators work on the tasks, it is natural for them to make mistakes as a result of lack of expertise, inherent ambiguity of tasks or unconscious bias. Enormous efforts were made in various manners to ensure the correctness of annotations as reported in, for example, [26, 17]. Saving efforts made for correctness will result in a noisy dataset but also potentially more annotated images. Trade-offs need to be made between the impact of by label noise and the gain of a larger dataset. For particular tasks, there exist freely available labels as alternatives to manual annotations. But these labels are often noisy owing to the way they were created. For example, one can use digital maps,

like OpenStreetMap, to segment aerial images. These segmentations constructed from maps suffer from the incomplete annotation as well as registration problems.[20] Besides, Pl@ntNet¹, a crowdsourcing platform, provide millions of images of plants and corresponding labels which may or may not be correct. Strong supervision is therefore required to correct the mistakes, for example double-checking the annotations over and over again and ensembling opinions from multiple annotators. It is therefore sometimes inevitable for deep neural network models to accept the existence of noisy annotations.

Why feature’s noise robustness? This paragraph should discuss 1. CNNs arch: hierarchical feature extraction and task-specific layers 2. CNNs based models benefit from transferability of hierarchical feature

Convolutional neural networks (CNNs) based models often contain two principal components: a stack of convolutional layers to extract hierarchical features and a few task-specific layers to fit the specific learning objectives.² The convolutional layers were proved “transferable” not only to another dataset[33] but also to another

¹<https://identify.plantnet-project.org/>

²M: As a non-expert it is rather unclear to me how to identify or design generic layers and task-specific ones, or what it actually means or why it is important to make this distinction. In addition, I wonder to what extent end-to-end training actually turns generic layers into task-specific ones. J: I think the idea to make this distinction here is to highlight that hierarchical features are reusable for many different applications and tasks but task-specific ones are not as they may have different number of neurons due to the different number of classes, or have even completely different architectures because the learning objectives are different for different app/tasks.

application[7, 18]. This feature transferability allows reusing the convolutional features for tasks and applications different from which they were originally trained with in a transfer learning scenario, and it helps to achieve better performance when there are only limited number of training samples[18].

This paragraph should introduce the idea of studying the impact of annotation noise on feature transferability. In cases where only a clean but small dataset is available, an extra noisy but large dataset in a similar domain might be helpful as it can be used to pre-train the convolutional features of the CNNs model. Previous studies[29, 23] have reported a negative impact of label noises on classification performance, but not yet on convolutional feature transferability. In general, optimal classification performance on test set often indicates that the extracted features are also optimal, whereas suboptimal classification performance does not necessarily reflect the convolutional features are also suboptimal, especially concerning feature transferability. Feature transferability describes the *generality of features*, i.e., the category-independence of features. Low-level features were proved to be less dependent to categories and thus more transferable to new tasks than high-level features. [33] We experimented in Section 3 that how much label noises interfere the transferability of convolutional features.

Narrow the problem of discussion down to Segmentation In this paper, we considered three types noises that happen to semantic image segmentation: mis-segmentation, mis-classification and inexhaustive segmentation, as described in details in Section 3.1. We chose to study segmentation errors because:

1. It is more difficult to correct noisy segmentation than to correct classification noises for object recognition;
2. Semantic segmentation can be treated as pixel-wise classification so that the existing noise-robust methods can be applied by assuming classifications for each pixels are made independently.

Explain why PU learning is relevant If we consider the inexhaustive segmentation issue only, i.e., the segmentation contains only a proportion of the target instances while leaving the rest unsegmented, the problem becomes similar to a so-called *positive and unlabelled learning* (PU learning) setup[15]. In the positive and unlabeled learn-

ing setup, the training dataset has two sets of examples: the *positive (P) set*, containing only positive examples, and the *unlabeled (U) set*, containing a mix of positive or negative examples. The segmented pixels in the presence of inexhaustive segmentation then form the P set and the unsegmented images construct the U set. Methods to learn with only positive examples and unlabeled examples become relevant because of the difficulty of distinguish object pixels and background pixels in the U set.

Table of contents

In the next section, we summarize related works in areas of transfer learning and learning with noisy labels for deep learning. In Section 3 we formulate the segmentation errors into three categories, mis-segmentation, mis-classification and incomplete annotation, and test feature transferability against them separately. In Section 4 we connect training with inexhaustive segmentations to positive and unlabeled (PU) learning. We experiment with synthesized noisy dataset in Section 5.1 to study whether the mis-segmentation, misclassification and inexhaustive segmentation noises have impacts on learning transferable features. Section 5.2 contains experiments for methods to learn with positive and unlabeled examples. Conclusions are summarized in section ??.

2 Related work

Transfer Learning *transfer learning* We sometimes have a learning task in one domain of interest, but we only have sufficient training data in another domain which does not share a feature space with the domain of interest. Transfer learning arises in this scenario to transfer knowledge from one domain to another and to improve the performance of learning by avoiding much expensive data-labeling efforts.[22] Weights of convolutional neural networks (CNNs) show outstanding transferability to another task. For example, weights trained on ImageNet images to perform image classification were shown successfully transferred to new categories and new learning problems[7, 18, 27]. Better performance were achieved for these tasks by using ImageNet pre-trained CNNs as initialization than training full model from scratch. Yosinski et al. discovered that feature transferability is negatively affected by the specialization of higher layer neurons and optimization difficulties caused by breaking co-

adapted neurons. Their experiments showed that low-level features, which are less dependent to particular categories, are more transferable than high-level features. **TODO:R Relations to our work.** We studied if feature transferability is negatively affected by the presence of label noises.

Unsupervised pre-training Apart from supervised pre-training, one can also obtain pre-trained features in an unsupervised or a semi-supervised way. The most common method is to train a generative model with either *auto-encoder* variants or *deep belief networks*. Vincent et al.[31] trained multiple levels of representation robust to the corrupted inputs with stacked denoising auto-encoders. Masci et al.[19] presented a stacked convolutional auto-encoder unsupervised pre-training for hierarchical feature extraction. Hinton et al.[11] proposed a greedy learning algorithm to train *deep belief nets* one layer at a time to train hierarchical features. Lee et al.[13] presented a *convolutional deep belief network*, to learn hierarchical convolutional representations. A few studies[4, 3, 1] highlighted the advantage of unsupervised pre-training compared to the random initialization, connecting unsupervised pre-training to a norm of regularization and a method that help disentangle the sample variations. However, better random initialization strategies, for example, xavier initialization[8] and its variants, have shortened the gap between unsupervised pre-training and random initialization. Using unsupervised pre-training or not now becomes a tradeoff between the time and resources invested and the performance gain. Unsupervised deep representation learning is in general not comparable to supervised representation learning especially when large scale dataset is available.

Deep Learning with Noisy Labels A few studies[29, 23] investigated the impact of label noise on classification performance with convolutional neural networks assuming the labels were randomly transited from one to another given the probabilities fall in a transition matrix. They found a significant decrease in classification performance along with the increase of false label proportion when the total number of examples is fixed. They then proposed methods to handle this label noise at random (NAR)[6] situation by either introducing a linear noise layer on top

of the output layer[29] or correcting the loss functions with an estimation of the noise transition matrix[23]. Xiao et al.[32] integrated a probabilistic graphic model to an end-to-end deep learning system to train predicting class labels, either correct or wrong, as well as to correct the wrong labels. Reed & Lee[24] proposed an empirical way of taking into account the *perceptual consistency* for large-scale object recognition and detection when incomplete and noisy labels exist by introducing a bootstrapping modification to the negative log-likelihood, in either a “Hard” or a “soft” favor.

Noise robustness In contrast to the works above, Rolnick et al.[25] argued that deep neural networks can learn robustly from the noisy dataset as long as an appropriate hyper parameters choice was made. They studied instead of replacing the correct labels with noisy labels but diluting correct labels with noisy labels to support their argument. They then concluded sufficiently large training set is of more importance than lower the level of noise. This work is closely related to our work in Section 3, except that we focus on the label noise robustness regarding the feature transferability instead of the classification performance. Additionally, most of these studies focus on the classification problems, whereas our work inclined more to the semantic segmentation problem.

Positive and Unlabeled Learning The previous studies about PU learning mainly focus on the binary classification for linear-separable problems[2, 14], whereas we showed in Section 4 that it is possible to train deep neural networks for multiple classes with only “positive” and unlabeled examples.

3 Noise robustness of feature transferability

From feature generality to feature robustness

Convolutional neural networks for images are believed to extract a rich hierarchy of image features. Some neurons capture concepts such as people and text while some units capture texture and material properties, such as dot arrays and specular reflections[7]. The conceptual features are specific for particular categories and the textural features can be shared by various categories. Further-

more, low-level features were found showing less specialization than the high-level features. By training a CNN with half of the ImageNet categories and transfer features to the other half, Yosinski et al.[33] found feature transferability decreases from the bottom layers to the top layers. The low level features, especially features of the first two layers, presented outstanding transferability between categories with far distance, for example natural categories and man-made classes as reported in [33]. Another evidence for this is that the first convolutional features for images are often observed as Gabor filters or color blobs even training with different datasets and different objectives[34, 13, 12, 27]. Given the evidence that low-level features can be category independent and shared by categories far from each other, we wondered if they are robust to learn with noisy training labels. More specifically speaking, the research question concerned in this section is:

1. Is transferability of convolution neural networks negatively influenced by the label noises?

Table of contents We formulated the learning task and annotation errors of interest in Section 3.1 and described the experiment setups in Section 5.1 to learn how transferable the features were when they were trained with a noisy dataset.

3.1 Problem Formulation

Semantic Segmentation *Segmentation is per-pixel classification* The goal of semantic segmentation tasks is to segment images into semantically meaningful partitions, a.k.a., *segments*. These segments are *target* if they depict instances of pre-defined object categories or *non-target* otherwise. A common interpretation for semantic image segmentation tasks with CNNs is the pixel-wise classification model: Given an image x , a segmenting model $f : R^{h \times w \times c} \rightarrow R^{h \times w}$ predicts a label for each pixel and predict a label map \hat{y} that has the same size as x . h, w are image height and weight respectively, and c is the number of image channels. The corresponding annotation y consist of labels for each pixel. Supposing there are K predefined categories, the label of pixel ij

$$y_{ij} = \begin{cases} k \in [1, K], & \text{for target pixels} \\ 0, & \text{for non-target pixels} \end{cases}$$

where $i \in [1, h], j \in [1, w]$.

Label corruption model *How noises synthesized* A straightforward way to model noisy labels is to corrupt true labels with a corruption model. The corruption model describes the probability of observed label map conditioning on the the true label map y , the image x and label errorness e :

$$p(\tilde{y}|x, y, e)$$

where the occurrence of errors e for pixels depends on the inputs x and true labels y . Such a noise model is called noisy not at random (NNAR) [6] because the noise depends on not only the true label y but also the inputs x . Given a corruption model and true annotations, one can synthesized the corresponding noisy annotations stochastically.

3.2 Noises in segmentation

Clarity for noise considered. In our works, we considered three types of errors for a semantic segmentation problem, *mis-segmentation*, *misclassification* and *inexhaustive segmentation*, and modified the NNAR model accordingly. Note that all these label errors apply to the whole segment instead of to individual pixels. That is pixels for the same segments will have the same true labels and observed labels. We excluded segmenting errors such as imprecise boundaries, oversegmenting or undersegmenting the objects from study because they are more difficulty to synthesize than the preceding classification errors. It would be of more value to use a real dataset for such errors than to synthesize.

Mis-segmentation Mis-segmentation denotes the wrongly segmented objects for categories that are semantically meaningful but are not predefined. For example, a toy dog can be misannotated as a dog, assuming that “dog” is predefined and “toy dog” is not. In the presence of mis-segmentation, pixel labels of segment S transit from 0 to k with probability

$$p(\tilde{y}_{ij} = k|x, y_{ij} = 0), ij \in S$$

The dependence of observed pixel labels \tilde{y}_{ij} on the original image x interpret the premise that mis-segmentation

would only happen to semantically meaningful segments in an image. It is natural to include this premise because semantically meaningless partitions of an image are less likely to be segmented by an annotator. However, it is difficult to estimate the above probability in practice because it is conditioning on the semantic meaning of x . Alternatively, we can synthesize mis-segmentation errors by selecting part of the categories as non-target categories so that instances of these categories should have zero labels for correct segmentations. We can then misannotate these non-target instances stochastically with a simplified probability $p(\tilde{y}_{ij} = k | y_{ij} = 0) = p_k$ without interpreting semantic meaning of x in probability. Note that p_k sums up to 1 for all classes, including class 0, $\sum_0^K p_k = 1$.

Misclassification Different from mis-segmentation, misclassification error means labels were misclassified between pre-defined categories. For example, cats may be misclassified as dogs occasionally if both “cat” and “dog” are target classes. In the presence of misclassification, pixel labels of segment S are transited from k to j stochastically with probability:

$$p(\tilde{y}_{ij} = j | y_{ij} = k) = p_{jk}, i, j \in S, k, j \in [1, K]$$

where $\sum_{j=1}^K p_{jk} = 1$. We assumed misclassification error is independent of the exact shape and appearance of the objects, i.e. information from x . This model is often called *noisy at random* [6]. This assumption does not hold in every cases of practice, for example, some instances can be more likely to be misclassified due to its ambiguity in shapes or appearances. But the difficulty of modeling the dependence of x leads to simply assuming an input independence. Given the class transition probabilities, one can easily synthesize noisy annotations including misclassification errors given a well-annotated segmentation dataset.

Inexhaustive segmentation Inexhaustive segmentation denotes that there exists unsegmented instances for pre-defined categories. Pixels of an unsegmented object S have labels flipped from k to 0 with probability:

$$p(\tilde{y}_{ij} = 0 | y_{ij} = k) = q_k, i, j \in S, k \in [1, K]$$

In words, an instance of category k is left unsegmented stochastically with probability q_k . The probability of cor-

rectly segmented in annotations is then:

$$p(\tilde{y}_{ij} = k | y_{ij} = k) = 1 - q_k, i, j \in S, k \in [1, K]$$

Again we assumed inexhaustive segmentations to be noisy at random (NAR).

4 Positive and Unlabeled Learning

This part should explain the necessity of Positive and Unlabeled Learning setup Experiments in Section 5.1 indicates that inexhaustive segmentation can have negative influences on feature transferability. Besides, including mis-segmented objects for training can aggravate the inexhaustive segmentation problem. For example, the existence of a mis-segmented toy dog does not mean that every toy dogs are mis-segmented. The other unsegmented toy dogs then become a source of inexhaustive segmentation and lead to worse fine-tuning performance as we discovered in Section 5. Method to compensate inexhaustive segmentation is therefore necessary to train better transferable representation.

Learning with inexhaustive segmentation fits a PU learning setup *This part should formulate the PU learning problem and discuss why PU Learning instead of semi-supervised learning is the way to go.*

To simplify the demonstration, let us consider a binary segmentation problem where we have *positive* and *negative* pixels denoting pixels corresponding to target and non-target segments. The goal of compensating incompleteness is to learn a model that predicts as many positive pixels as possible while keeping the false positive rate low. A pixel with a negative label can be either a truly negative pixel or a wrongly unsegmented positive pixel. In other words, there is no explicitly labeled negative examples in the training set. This naturally fits in a Positive and Unlabeled Learning (PU Learning) setup where the training dataset consists of only a set of positive examples (P set) and a set of unlabeled examples (U set), meaning that there are no reliable negative examples available. The traditional semi-supervised learning techniques are not applicable in such a setup as a result of the absence of negative training samples.

Class-weighted Logistic Loss *This part should discuss the difficulty of applying traditional PU learning methods to deep learning and then introduce the idea of re-weight positive and negative classes.* Traditional PU learning methods often follow a two-step strategy: first identifying a set of reliable negative samples (RN set) from U set and then iteratively build classifiers, either naive Bayesian (NB) or supported vector machine (SVM), with RN and P sets and update the RN set with expectation-maximization (EM) algorithm. However, this strategy requires training multiple classification models. It is unrealistic to train iteratively a sequence of deep learning models because it would take significantly longer time to train neural networks than to train a naive Bayesian (NB) model or supported vector machine (SVM).

Our original goal was to achieve high precision as well as high recall regardless the existence of false negative labels. A straightforward way to achieve this goal is to simply reweigh the positive and negative examples, namely let the positive and negative examples have different rates of contribution to the total loss. Suppose logistic loss is used, the corresponding losses for positive and negative samples are:

$$\begin{aligned} l_{\tilde{y}_i=+1} &= -\log p(y_i = +1|x_i) \\ l_{\tilde{y}_i=-1} &= -q \log p(y_i = -1|x_i), 0 < q < 1 \end{aligned}$$

where $p(y_i|x_i) = \sigma(f(x))$ denotes the probabilistic output of the model $f(\cdot)$ for the i -th example. Empirically, the choice of q can be made based on the most highest precision and recall achieved on validation set. Alternatively, one can also roughly assign $q = p(y = -1|\tilde{y} = -1)$. This turns out to be part of the backward corrected loss proposed in [23]:

$$\begin{aligned} l_{\tilde{y}_i=-1} &= -p(y_i = -1|\tilde{y}_i = -1) \log p(y_i = -1|x_i) \\ &\quad -p(y_i = +1|\tilde{y}_i = -1) \log p(y_i = +1|x_i) \end{aligned}$$

with $p(y_i = -1|\tilde{y}_i = -1) = q$ and $p(y_i = +1|\tilde{y}_i = -1) = 1 - q$.

Confidence of contrary prediction and probability of false negative label *The idea of alleviating punishment for confident predictions.* Losses above assume that the true label of an observed negative label is independent of the inputs x . However, the true label y is often dependent

on both x and \tilde{y} in practice. For instance, in segmentation problems, whether a pixel is left unsegmented correctly or not can be determined by the semantical meaning of this pixel and pixels around. It is nevertheless difficult to estimate $p(y|x, \tilde{y})$ directly due to the difficulty of exploring the joint distribution of x and \tilde{y} . Alternatively, one can use the probabilistic prediction $p(y|x; \theta)$ a classification model to provide extra information about the inputs. The predicted probabilities indicate how confident the current model is for the predictions made. Given a good enough model, more confident positive predictions for negatively labeled examples can have a higher probability of being wrongly annotated than the less confident ones.³ With a random model by which predictions are made randomly, the high and low confidences do not convey any information about the underlying inputs distribution. By contrast, confident predictions made by a trained model indicate that the corresponding samples are possibly close to some previous training samples with positive labels in the feature space. The natural trend of similar examples having the same labels supports a higher probability of being annotated wrongly for confident predictions which have opposite labels. This thought leads to the idea of not punishing confident contrary predictions more than unconfident ones as the traditional logistic/softmax loss will do.

Exponential Loss for unlabeled examples *This section should mention the class dependent losses and introduce ExponentialUnlabeledLoss* In a PU learning setup, the positive (P) set contains only reliable positive labels, whereas the unlabeled (U) set can be considered as noisy negative labels. The problem then converts to training with clean positive examples and noisy negative examples. We used a class-dependent loss to compensate the noisy negative labels while still making full use of the clean positive labels. The loss was made of a normal logistic loss for positive examples and an exponential loss [30] for examples with negative labels:

$$\begin{aligned} l_{\tilde{y}_i=+1} &= -\log p(y_i = +1|x_i) \\ l_{\tilde{y}_i=-1} &= 1 - p(y_i = -1|x_i) \end{aligned}$$

Figure 1 shows the weighted logistic loss with $q = 0.5$ and the exponential unlabeled loss respectively and their

³J: This argument needs more detailed explanation.

derivatives with respect to logits by varying the logit from negative to positive. The main feature of exponential loss is its relatively small changes in the region of confident positive for negatively labeled examples, compared to the logistic loss and class weighted logistic loss. As a consequence of this feature, the corresponding derivative decreases to zero as the prediction increases in the positive direction. This feature interprets the idea of not punishing positive prediction with confidence for negatively labeled samples. Another modified loss function, namely the Focal Loss[16], instantiates the same idea with a similar form of expression.

This section discuss the loss and derivative difference between the losses. Figure 2 shows the loss and derivate difference between logistic loss, class-weighted logistics loss and exponential negative (unlabeled) loss with a two-dimensional example. It demonstrates that, for exponential unlabeled loss, unlabeled positive examples farther from the decision boundary do not have larger loss contributions as ones closer to but still distant from the decision boundary. Confident positive predictions, shown as examples located on the positive side of the decision boundary and far from it, has little effect for updating model weights. The consequence of this effect is that positive examples push the decision boundary away from the positive cluster while negative examples closed to the decision boundary instead of those away from the decision boundary pull the decision boundary towards the positive cluster. This characterisic of exponential loss lead to a selectively counting weights update contributions of negative samples than simply lower the overall estimation for all negative samples while optimization. The exponential loss was introduced in [30] to get rid of the effect of outliers. In our case, the negative examples given confident predictions by classifier can be considered as outliers.

Minimum entropy regularization and bootstrapping objective An alternative way of encouraging confident positive predictions is to introduce minimum entropy regularization[9]. Reed et al.[24] proposed an emprical modification to softmax loss, a.k.a. the *bootstrapping loss*, which can be used to learn with unlabeled examples:

$$l_{\tilde{y}_i=-1} = -\beta \log p(y_i = -1|x_i) - (1 - \beta) \log p(y_i = \hat{y}_i|x_i)$$

where $\hat{y} = \operatorname{argmax}_{j \in \{-1, +1\}} p(y_i = j|x_i)$ is the model prediction and $0 < \beta < 1$. The first term of the objective

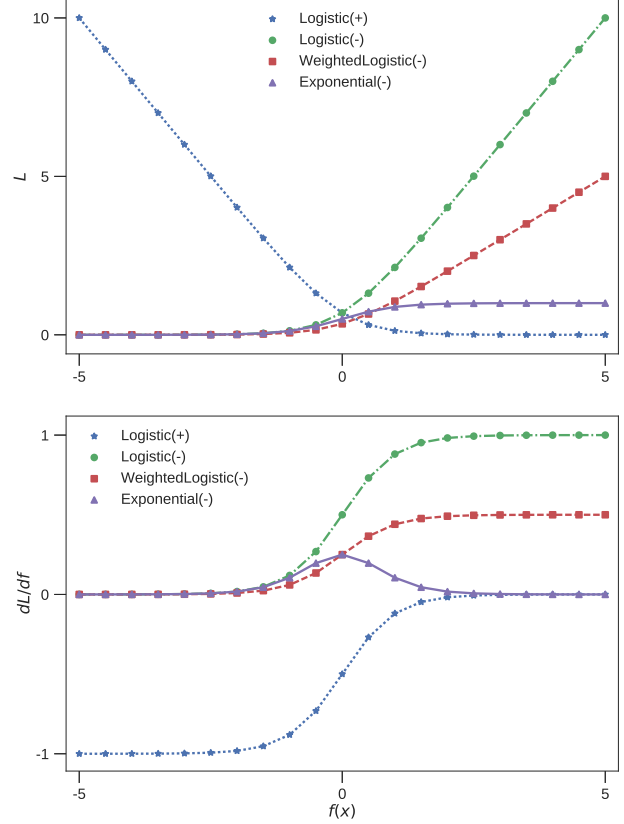


Figure 1: The Logistic Loss, Weighted Logistic Loss, Exponential Loss and their derivatives with respect to logits.

is a weighted logistic loss and the second term can be considered as a variation of minimum entropy regularization:

$$\begin{aligned} H &= - \sum_{j \in \{-1, +1\}} p(y_i = j|x_i) \log p(y_i = j|x_i) \\ &\sim - \sum_{j \in \{-1, +1\}} \delta(y_i - \hat{y}_i) \log p(y_i = j|x_i) \end{aligned}$$

which intuitively encourages the model to make confident predictions[9].

Multiclass PU learning *This part should describe how to extend the losses for multiclass scenarios and the difficulty of having multiple positive classes. So far, we have*

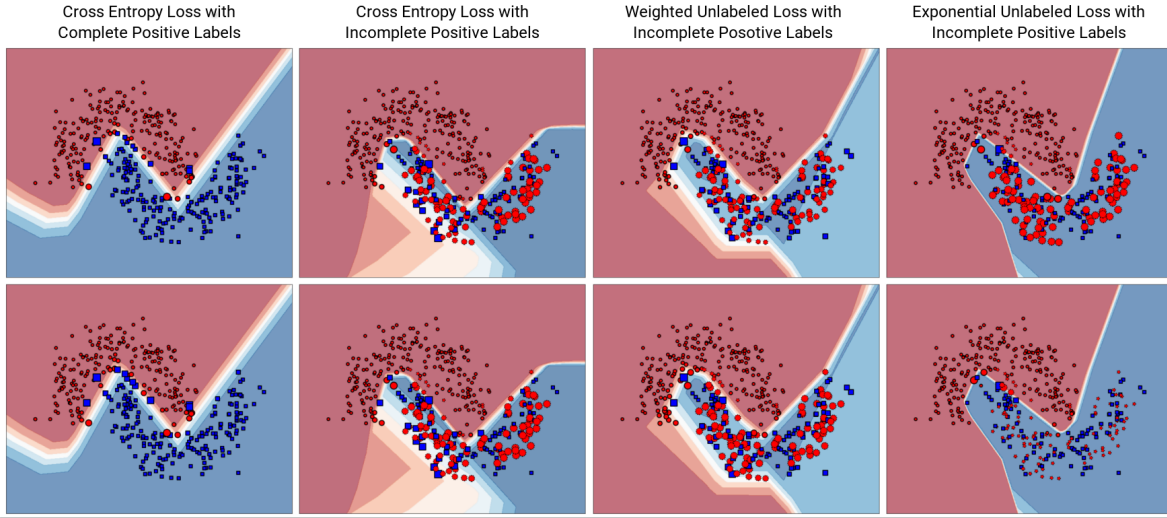


Figure 2: 2D moons dataset with non-linear separable decision boundary. Four hundreds samples per class were drawn randomly from two interleaving half circles with noises added with a minor standard deviation. A **red circle** indicates an example labelled as positive whilst a **blue square** indicates the example has a negative label. The **leftmost** figures have complete positive labels, meaning the positive and negative labels are all correct, whereas, in **the other figures** only half of the positives were correctly labelled and the rest were mixed with the negative samples. The **background colors** represent the probability for the area to be positive given by the classifier trained with the given samples and labels: **red** for high probability areas, **blue** for low probability areas and **white** for the class transition areas, i.e. decision boundaries. The **size of the markers** in the top row denotes the per-class normalized training losses and the **size of the markers** in the bottom row the per-class normalized derivatives w.r.t the output of the last layer for the trained Multilayer Perceptron (MLP) with the different losses.

been discussing modifications to the logistic loss for binary classification and segmentation. Similar modifications can be applied to softmax loss as well. In a multi-class positive and unlabeled learning setup, $K > 1$ positive classes are supposed to be distinguished from one negative class, whereas only part of the positive examples are labeled, and the others are not. Softmax loss can be still used for positive classes because positive labels are reliable as long as they are given. By contrast, the negatively labeled samples are a mix of samples assigned with true negative labels and false positive labels. The weighted unlabeled logistic, exponential unlabeled loss and bootstrapping can be used for negative by simply replacing logistic function with softmax function as the activation function for logits.

From classification to segmentation *This paragraph should highlight the problem of spatial independence assumption.* Classification for each pixel is made independently when segmentation task is considered as per-pixel classification. The objective modifications to alleviate the negative impact of unlabeled positive samples can be applied to classification for each pixel by assuming the probability of missing positive label for a pixel is independent of its neighbor pixels, which may not hold in practice. For example, we stochastically flipped pixel labels for the whole segment, instead of individual pixel label, to synthesize incomplete segmentation errors in Section 3. Even with annotation errors like under-segmented objects, the probability for a pixel to be unsegmented depends on its neighbor pixels. We made a strong assumption of spatial independence for false negative labels for inexhaustive segmentation noises. Future works maybe required to relax this assumption and make use of the spatial dependence of label noises to achieve higher mean intersection over union (mean IU) not only higher accuracy.

Implementation details *This paragraph should explain fade-in was introduced to avoid all-positive initial prediction;* The exponential loss for negative (unlabeled) examples saturates for very positive outputs, meaning that the confident positive prediction has little contribution to the total loss. This can introduce problems at the beginning of the training procedure when the confident predictions are likely to be made at random. Additionally, optimiza-

tion could reach the plateau where the model made all positive predictions with high confidence. Therefore, we introduced the exponential loss after training with the normal logistic/softmax loss for a few epochs. We applied a similar “fade-in” mechanism to the bootstrapping objective as well because it also requires a non-random model for sufficiently trustworthy prediction \hat{y} .

This part should explain the influence of the imbalanced problem and how to overcome. Another problem encountered in the PU learning setup is the class imbalance introduced by negatively labeled positive samples. Even 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 Synthesized mis-segmentation, misclassification and inexhaustive segmentations

Experiment setup In order to investigate the influence of mis-segmentation, misclassification and inexhaustive segmentation on feature transferability respectively, we set up experiments with a perfectly annotated dataset, the PASCAL VOC2011 dataset[5]. Fifteen out of twenty categories 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[18] for segmentation in the presence or absence of synthesized segmentation errors. The fine-tuning dataset was used to fine-tune the convolutional weights from the pre-trained FCN-AlexNet models. Fine-tuned models were then evaluated by mean intersection over union ratio (mean IU) achieved on the test set of fine-tuning dataset, referring to as the *fine-tuning performance*. Performance improvement of fine-tuning models compared to a randomly initialized model indicates the transferability of pre-trained weights.

Experiment details To avoid 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. Each fold was studied separately, and the exact partitions of each fold is listed in Table 1. The training dataset was enriched with extra segmentations by Hariharan et al.[10] 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 dataset and evaluation. In order to accelerate the training process, we subsampled the original images by four times. Fully Convolutional Networks with AlexNet was used for experiments because its relatively small capacity and thus short training time. The existence of an ImageNet model for AlexNet can be beneficial to set a performance reference. Only convolutional filters of AlexNet were transferred from the pre-training phase to fine-tuning phase because the transferability of convolutional weights were the focus of this work. The other layers were random initialized with Xavier Initialization. The ImageNet model and completely random weight initialization were considered as the upper bound and lower bound, respectively, for various pre-trained weights summarized in Table 1. The default hyperparameters of FCN-AlexNet in [18] were kept unchanged. Training 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.

Feature Transferability robustness to segmentation noises *What Table 1 tell us. How annotation errors were synthesized; How synthesizations are different from reality; Transferability of noisy models compared to clean models.* Mis-segmentation, misclassification and inexhaustive segmentation were synthesized separately with stochastic corruptions to the well-annotated pre-training dataset followed the descriptions in Section 3.1.

Mis-segmentation: little effect on weights transferability To synthesize mis-segmentation noises, 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 mis-segmentation noises, instances from non-target categories can be misannotated as the target category with probability of $p_1 = 1$ and

$p_1 = 0.5$ respectively. The two choices of probability led to two different pre-training sets and thus two different pre-trained models, naming the AllMisSegmented model and the HalfMisSegmented model, in Table 1 respectively. The noise-free counterpart of mis-segmentation is an dataset containing segmentations of the selected target category only whilst the other 14 categories remained unsegmented. Model trained with this noise-free dataset was denoted as NoMisSegmented in Table 1.

Table 1 shows that all three models achieved better fine-tuning performances than random initialization. The dataset mis-segmented all non-target instances produced a model with even slightly better fine-tuning performance than the dataset segmented only the target category. However, it is less likely to happen in practice that annotators will mis-segment every single non-target instance in the dataset. Mis-segmentations often occur in annotations occasionally. Therefore, we also trained models with a training set containing half of the non-target objects to test if inexhaustively mis-segmenting non-target objects decreased the fine-tuning performance. The results show that the HalfMisSegmented model had a slightly worse fine-tuning performance than the AllMisSegmented model but was still comparable to the NoMisSegmented model. Based on these observations, we concluded that mis-segmenting semantically meaningful objects could have little impact when they are used for pre-training transferable weights.

Misclassification: negatively affect weights transferability Misclassification errors were also synthesized from the well-annotated pre-training dataset as described in Section 3.1. The noisy dataset containing labels for target segments stochastically transited to a random class with probability $p_{jk} = \frac{1}{20}$. The resulted trained model was denoted as the AllRandomLabels model in Table 1. Similarly, if a random half of the segmented objects were assigned random labels, the resulting pre-trained model is called the HalfRandomLabels model. The noise-free counterpart of these two noisy models was the model trained with true labels, denoting as the TrueLabels model.

Compared to the TrueLabels model, the noisy models trained with samples containing random labels, no matter if all labels were random or if only half of the labels were random, led to worse fine-tuning performances. Fine-tuned performances of the AllRandomLabels model and

Initial Representation	mean IU (aeroplane, bicycle, bird, boat, bottle)	mean IU (bus, car, cat, chair, cow)	mean IU (dining table, dog, horse, motorbike, person)	mean IU (potted plant, sheep, sofa, train, TV)	avg mean IU and avg std.
ImageNetModel	0.42 ± 0.01	0.51 ± 0.01	0.49 ± 0.01	0.47 ± 0.01	0.47 ± 0.01
RandomWeights	0.29 ± 0.01	0.29 ± 0.03	0.27 ± 0.01	0.30 ± 0.02	0.29 ± 0.02
NoMisSegmented	0.26 ± 0.01	0.37 ± 0.03	0.27 ± 0.01	0.33 ± 0.04	0.31 ± 0.02
AllMisSegmented	0.30 ± 0.02	0.35 ± 0.01	0.29 ± 0.02	0.35 ± 0.03	0.32 ± 0.02
HalfMisSegmented	0.27 ± 0.01	0.34 ± 0.01	0.30 ± 0.01	0.32 ± 0.01	0.31 ± 0.01
TrueLabels	0.29 ± 0.01	0.36 ± 0.01	0.29 ± 0.01	0.37 ± 0.01	0.33 ± 0.01
AllRandomLabels	0.29 ± 0.01	0.33 ± 0.03	0.26 ± 0.01	0.28 ± 0.01	0.29 ± 0.01
HalfRandomLabels	0.27 ± 0.01	0.33 ± 0.02	0.25 ± 0.01	0.29 ± 0.01	0.29 ± 0.01
InexhaustiveSegmented	0.26 ± 0.01	0.30 ± 0.3	0.28 ± 0.03	0.32 ± 0.02	0.29 ± 0.02

Table 1: Performances of fine-tuned FCN-AlexNet models with different representation initializations. **ImageNetModel** represents the pre-trained ImageNet model; **RandomWeights** indicates that the weights were randomly initialized; All the other weights were first trained with the pre-training dataset in the presence or the absence of different types of label noises. Each experiment was repeated three times, the mean and the standard deviation were computed over the last five snapshots for all repetitions.

the HalfRandomLabels model were no better than randomly initializing model weights, indicating poor weights transferabilities of a trained model in the presence of random labels to segmentations. In other words, misclassification noises in segmentation can impact the transferability of convolutional weights negatively.

Inexhaustive Segmentation: negatively affect weights transferability Inexhaustive segmentations in the training dataset were synthesized by randomly converting labels of segmented objects to 0 with probability $q_k = 0.5$. Similar as misclassification errors, inexhaustive segmentation can have negative impact on weights transferability. Pre-trained model trained from a dataset with 50% percentage of the instances unsegmented produced a fine-tuned model with an average mean IU 0.04 worse than the model pre-trained with true labels and it was almost the same as training a model with random weight initialization.

Categorizing classes for pre-training *This paragraph should discuss the potential benefit of categorizing classes. It had equivalent fine-tuning performance as using true labels; It achieved better performance than training to segment the exact classes but with misclassification labels.* Table 1 also showed that the fine-tuning performance of weights (AllMisSegmented) trained with bi-

narized labels was better than weights (HalfRandomLabels) pre-trained with random class labels. This observation can be relevant for learning transferable weights in the presence of misclassification because one can train binary segmentation if the pre-training dataset is dominated by misclassification errors. Besides, weights pre-trained by binary segmentation, segmenting object pixels from the background, can have a comparable fine-tuning performance to weights pre-trained by multi-class segmentation. In our experiment, the number of samples for each class in the pre-training dataset was limited. Most of the classes had only around one hundred images, and the dining table had only 20 images. The limited number of class samples can increase the difficulty to segment individual classes and may explain why binarized labels could achieve comparable fine-tuning performance to true labels. We then studied the influence of varying number of categories for pre-training. We categorized the fifteen classes in the pre-training set into person, animal, vehicle, indoor according to [5] and trained, shown as the error bars on lines in figure 3 at categories=4. The fifteen classes were also randomly categorized into 4, 7, 11 categories and shown as separate error bars in figure 3 at categories=4, 7, 11 respectively. Figure 3 shows that categorizing classes into categories had little effect to the fine-tuning performance of trained weights. Even categorizing classes into random categories without explicit meaning could pre-train weights better than random initialization. Binarizing or categorizing classes into higher hierarchy can be beneficial when the main learning objective is to train transferable convolutional weights, and the training dataset is corrupted by noisy labels.

5.2 PU Learning for classification and for segmentation

In order to compare exponential unlabeled loss with class weighted logistic/softmax loss, we synthesized positive and unlabeled learning setups for classification and segmentation respectively.

In the classification setup, we combined the CIFAR10 dataset and CIFAR100 dataset, using CIFAR10 as the positive (P) set and CIFAR100 as the negative (N) set. The learning objective is to classify images into eleven classes: ten positive classes from CIFAR10 and a neg-

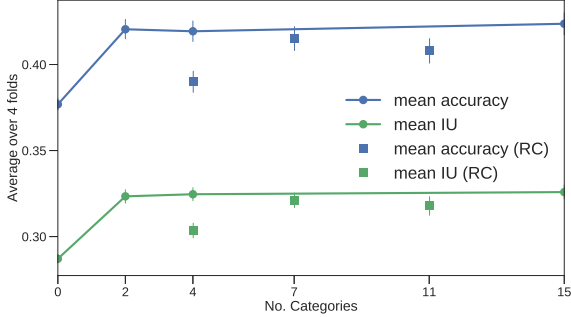


Figure 3: Test performance for fine-tuned models initialized with weights pre-trained with categorized 15 classes. Separate error bars located aside from the lines denote random categorizations (RC) of the 15 classes. The displayed mean IU/mean accuracies and standard deviations were averaged over four folds listed in Table 1. Experiments were repeated three times.

ative class for images from CIFAR100. Note that there is no category overlap between CIFAR10 dataset and CIFAR100 dataset. To synthesize a positive and unlabeled (PU) learning setup, we selected only part of positive images from CIFAR10 to be correctly labeled and the rest of the CIFAR10 images were mixed with CIFAR100 images, forming the unlabeled (U) set. Models were then trained with the labeled part of P set and U set. An eight layer VGG net was used together with different choices of losses. The architecture of this VGG8 model can be found in Appendix TODO:Jihong. Each model was trained from scratch with Adam optimizer base and learning rate 0.0001. Model performances were evaluated on a separate test set of combined CIFAR10 and CIFAR100 with true labels.

Table 2 summarizes the test precisions and recalls for training with different losses in the PU setup, compared with training with complete positive labels. With a training set containing 50% labeled CIFAR10 images and the rest unlabeled, the normal cross-entropy loss led to an imbalanced model with high precision but low recall, and therefore with a low f1-score. By reweighing the negative loss by a factor of 0.5, we were able to balance precision and recall and improve the resulting f1-score. Compared to the negative weighted loss, exponential loss and (hard)

Annotation	Loss	acc.	mean prec.	mean rec.	mean F_1
Complete	CrossEntropyU.	0.87	0.88	0.82	0.85
50%(P+N)	CrossEntropyU.	0.83	0.84	0.78	0.80
50%P+U	CrossEntropyU.	0.66	0.94	0.39	0.50
50%P+U	WeightedU.	0.78	0.75	0.75	0.76
50%P+U	ExponentialU.	0.81	0.85 \pm 0.03	0.72 \pm 0.03	0.77
50%P+U	BootstrapHard	0.80	0.76	0.81	0.78

Table 2: Accuracy, mean precision, mean recall and mean f1-score on test set of the CIFAR dataset with true labels. The complete dataset contains images from CIFAR10 as the **positive** (P) set and images from CIFAR100 as the **negative** (N) set. The unlabeled positive examples from P set construct the **unlabeled** (U) set, together with N set. Precision and recall were averaged over ten positive classes. 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.

bootstrapping loss were able to achieve slightly better f1-scores. TODO: explain Training f1 0.83 vs 0.81

As a reference for the performances, we trained a classifier with 50% of the positive samples and the same percentage of true negative samples. We referred this setup as positive and negative (PN) setup. The total number of training sample in PN setup is smaller because the rest unlabeled positive and negative samples were excluded from training. In Figure 4 we varied the percentage of labeled positive images, and compared the three different losses in the PU setup with a normal cross-entropy loss in the PN setup. In any of the labeled percentages for positive images, training with positive and negative examples can achieve higher f1-scores than any of the models trained with the same amounts of positive images and unlabeled images. The performance difference between learning with PN and learning with PU increases as the number of labeled positive images decreases. This result was expected because PN setup delivers extra information about which images in the unlabeled set are negative. The PU setup is therefore only relevant when it is difficult to annotate negative examples from the unlabeled data. And segmentation problem in the presence of inexhaustive segmentation can be such an example.

In the segmentation setup, we used again the PASCAL VOC2011 dataset with extra segmentation[10]. We synthesized inexhaustive segmentations the same way as described in Section 5.1. The same AlexNet-FCN model

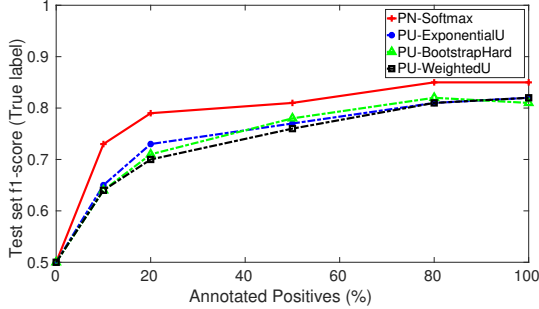


Figure 4: Varying percentage of labeled positive images. **P+N** represents training with percentage of images with reliable positive and negative labels and **P+U** stands for training with the positive (P) and unlabeled (U) sets.

were trained together with the different loss functions for class 0 to predict binary segmentation, determining whether a pixel is correspondent to an object or not. Only single-object images were used for training and testing in order to avoid the influence of two adjacent objects joining as one object because of binary segmentation. The same hyperparameters for optimization were used as in Section 5.1. The trained models were evaluated with the test set of PASACAL VOC2011 segmentation dataset with binary segmentations.

As shown in Table 3, the exponential unlabeled loss achieved the highest mean accuracy and a slightly lower overall accuracy. In contrast to the improvement of mean accuracy, mean IU for models trained with either weighted unlabeled loss or exponential unlabeled loss did not show significant improvement to the normal cross entropy loss.

Selective predictions for models trained with exponential unlabeled (ExpU.) loss and normal cross entropy (CrossEnt.) loss were presented in Figure 5. For these two example images, the model trained with cross entropy loss failed to segment objects from images whereas exponential unlabeled loss segmented on the position of the objects though with coarse outlines. The third column shows predictions given by model trained with complete training segmentation, and it did not give more accurate outlines. The coarse results were mainly due to the limited compacity of AlexNet model.

Annotation	Loss	overall acc.	mean acc.	f.w. IU	mean IU
Complete	CrossEnt.U	0.90	0.85	0.82	0.75
50%Unseg.	CrossEnt.U	0.85	0.68	0.73	0.60
50%Unseg.	WeightedU	0.84	0.71	0.73	0.62
50%Unseg.	ExponentialU	0.83	0.75	0.72	0.62

Table 3: Best binary segmentation performance achieved on the test set of PASCAL VOC2011 segmentation dataset in the presence of inexhaustive segmentation. Class weight 0.7:1.75 was used to balance the sample frequency difference of the two classes and negative loss were further weighted by a factor of 0.5 for weighted unlabeled loss. Mean accuracy is equivalent to mean recall over classes. Mean IU is the average intersection over union ratio (IU) over two classes and f.w. IU is the frequency weighted average of IU over the two classes. Experiments were repeated twice and standard deviations were approximately 0.01.

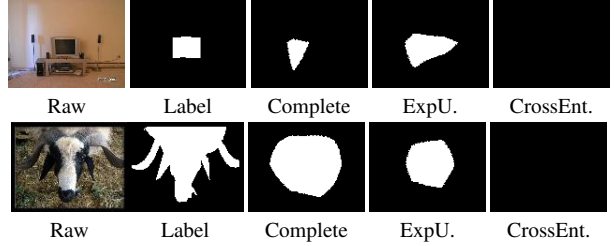


Figure 5: Selective predictions for models in Table 3.

6 Conclusion

- Feature transferability is robust to mis-segmentation but not to misclassification and inexhaustive segmentation.
- Misclassification noises can be alleviated by binarizing/categorizing classes.
- Inexhaustive segmentation can be translated as learning with only positive and unlabeled examples.
- We proposed a class dependent loss to not over-punish confident positive predictions in the presence of noisy negative labels, and it showed slightly better results than class-weighted loss.

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] 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.
- [3] 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.
- [4] 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.
- [5] 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.
- [6] Benoît Frénay and Michel Verleysen. Classification in the presence of label noise: a survey. *IEEE transactions on neural networks and learning systems*, 25(5):845–869, 2014.
- [7] 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.
- [8] 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.
- [9] Yves Grandvalet and Yoshua Bengio. Semi-supervised learning by entropy minimization. In *Advances in neural information processing systems*, pages 529–536, 2005.
- [10] 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.
- [11] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006.
- [12] 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.
- [13] 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.
- [14] Wee Sun Lee and Bing Liu. Learning with positive and unlabeled examples using weighted logistic regression. In *ICML*, volume 3, pages 448–455, 2003.
- [15] Xiao-Li Li and Bing Liu. Learning from positive and unlabeled examples with different data distributions. *Machine Learning: ECML 2005*, pages 218–229, 2005.
- [16] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. *arXiv preprint arXiv:1708.02002*, 2017.
- [17] 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.
- [18] 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.
- [19] 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.
- [20] 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.
- [21] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic segmentation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 891–898, 2014.
- [22] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2010.

- [23] 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.
- [24] 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.
- [25] David Rolnick, Andreas Veit, Serge Belongie, and Nir Shavit. Deep learning is robust to massive label noise. *arXiv preprint arXiv:1705.10694*, 2017.
- [26] 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.
- [27] 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.
- [28] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [29] Sainbayar Sukhbaatar, Joan Bruna, Manohar Paluri, Lubomir Bourdev, and Rob Fergus. Training convolutional networks with noisy labels. *arXiv preprint arXiv:1406.2080*, 2014.
- [30] 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.
- [31] 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.
- [32] Tong Xiao, Tian Xia, Yi Yang, Chang Huang, and Xiaogang Wang. Learning from massive noisy labeled data for image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2691–2699, 2015.
- [33] 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.
- [34] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer, 2014.
- [35] 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.

A Convolutional Networks for Semantic Segmentation

B Deep Learning with Label Noise

NNAR, NAR, NCAR

C Evaluation metrics

(overall) accuracy

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

precision

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

recall

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

f1-score

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

intersection over union (IU)

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

D Additional Results

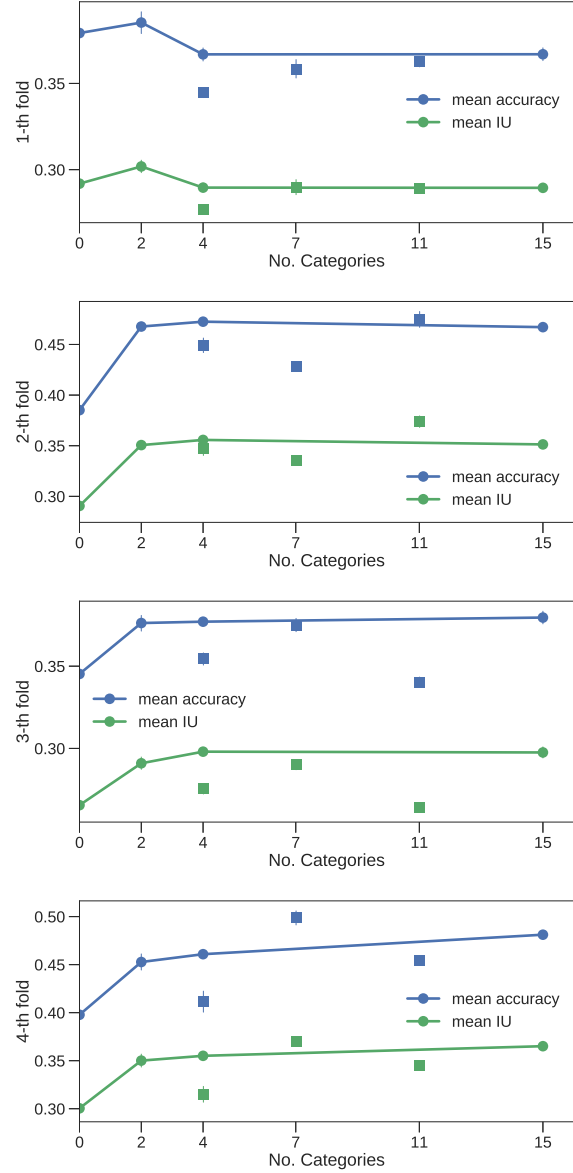


Figure 6: The influence of number of pre-training categories on performance of fine-tuned models for each fold, addition to Figure 3.