

A survey on Semi-, Self- and Unsupervised Techniques in Image Classification

Similarities, Differences & Combinations

Lars Schmarje, Monty Santarossa, Simon-Martin Schröder, Reinhard Koch
 Multimedia Information Processing Group, Kiel University, Germany
 {las,msa,sms,rk}@informatik.uni-kiel.de

Abstract

While deep learning strategies achieve outstanding results in computer vision tasks, one issue remains. The current strategies rely heavily on a huge amount of labeled data. In many real-world problems it is not feasible to create such an amount of labeled training data. Therefore, researchers try to incorporate unlabeled data into the training process to reach equal results with fewer labels. Due to a lot of concurrent research, it is difficult to keep track of recent developments. In this survey we provide an overview of often used techniques and methods in image classification with fewer labels. We compare 21 methods. In our analysis we identify three major trends.

- 1. State-of-the-art methods are scalable to real world applications based on their accuracy.**
- 2. The degree of supervision which is needed to achieve comparable results to the usage of all labels is decreasing.**
- 3. All methods share common techniques while only few methods combine these techniques to achieve better performance.**

Based on all of these three trends we discover future research opportunities.

1. Introduction

Deep learning strategies achieve outstanding successes in computer vision tasks. They reach the best performance in a diverse range of tasks such as image classification, object detection or semantic segmentation.

The quality of a deep neural network is strongly influenced by the number of labeled / supervised images. ImageNet [26] is a huge labeled dataset which allows the training of networks with impressive performance.

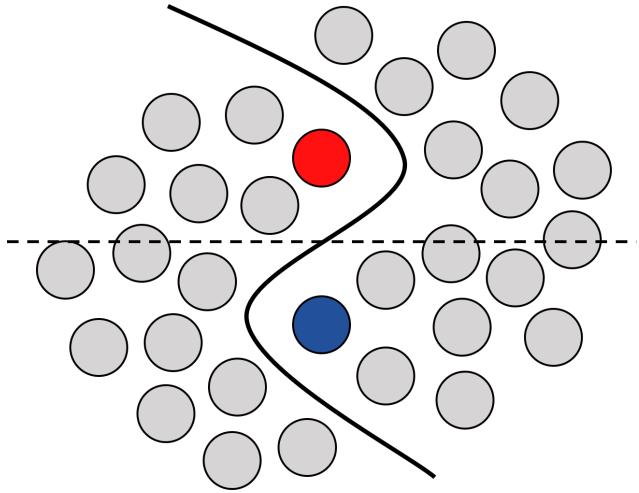


Figure 1: This image illustrates and simplifies the benefit of using unlabeled data during deep learning training. The red and dark blue circles represent labeled data points of different classes. The light grey circles represent unlabeled data points. If we have only a small number of labeled data available we can only make assumptions (dotted line) over the underlying true distribution (black line). This true distribution can only be determined if we also consider the unlabeled data points and clarify the decision boundary.

Recent research shows that even larger datasets than ImageNet can improve these results [31]. However, in many real world applications it is not possible to create labeled datasets with millions of images. A common strategy for dealing with this problem is transfer learning. This strategy improves results even on small and specialized datasets like medical imaging [40]. While this might be a practical workaround for some applications, the fundamental issue remains: Unlike hu-

mans, supervised learning needs enormous amounts of labeled data.

For a given problem we often have access to a large dataset of unlabeled data. Xie et al. were among the first to investigate unsupervised deep learning strategies to leverage this data [45]. Since then, the usage of unlabeled data has been researched in numerous ways and has created research fields like semi-supervised, self-supervised, weakly-supervised or metric learning [23]. The idea that unifies these approaches is that using unlabeled data is beneficial during the training process (see [Figure 1](#) for an illustration). It either makes the training with few labels more robust or in some rare cases even surpasses the supervised cases [21]. Due to this benefit, many researchers and companies work in the in the field of semi-, self- and unsupervised learning. The main goal is to close the gap between semi-supervised and supervised learning or even surpass these results. Considering presented methods like [49, 46] we believe that research is at the break point of achieving this goal. Hence, there is a lot of research ongoing in this field. This survey provides an overview to keep track of the major and recent developments in semi-, self- and unsupervised learning.

Most investigated research topics share a variety of common ideas while differing in goal, application contexts and implementation details. This survey gives an overview in this wide range of research topics. The focus of this survey is on describing the similarities and differences between the methods. Moreover, we will look at combinations of different techniques.

While we look at a broad range of learning strategies, we compare these methods only based on the image classification task. The addressed audience of this survey consists of deep learning researchers or interested people with comparable preliminary knowledge who want to keep track of recent developments in the field of semi-, self- and unsupervised learning.

1.1. Related Work

In this subsection we give a quick overview about previous works and reference topics we will not address further in order to maintain the focus of this survey.

The research of semi- and unsupervised techniques in computer vision has a long history. There has been a variety of research and even surveys on this topic. Un-

supervised cluster algorithms were researched before the breakthrough of deep learning and are still widely used [30]. There are already extensive surveys that describe unsupervised and semi-supervised strategies without deep learning [47, 51]. We will focus only on techniques including deep neural networks.

Many newer surveys focus only on self-, semi- or unsupervised learning [33, 22, 44].

Min et al. wrote an overview about unsupervised deep learning strategies [33]. They presented the beginning in this field of research from a network architecture perspective. The authors looked at a broad range of architectures. We focus ourselves on only one architecture which Min et al. refer to as "Clustering deep neural network (CDNN)-based deep clustering" [33]. Even though the work was published in 2018, it already misses the recent development in deep learning of the last years. We look at these more recent developments and show the connections to other research fields that Min et al. didn't include.

Van Engelen and Hoos give a broad overview about general and recent semi-supervised methods [44]. While they cover some recent developments, the newest deep learning strategies are not covered. Furthermore, the authors do not explicitly compare the presented methods based on their structure or performance. We provide such a comparison and also include self- and unsupervised methods.

Jing and Tian concentrated their survey on recent developments in self-supervised learning [22]. Like us the authors provide an performance comparison and a taxonomy. They do not compare the methods based on their underlying techniques. Jing and Tian look at different tasks apart from classification but ignore semi- and unsupervised methods.

Qi and Luo are one of the few who look at self-, semi- and unsupervised learning in one survey [38]. However, they look at the different learning strategies separately and give comparison only inside the respective learning strategy. We distinguish between these strategies but we look also at the similarities between them. We show that bridging these gaps leads to new insights, improved performance and future research approaches.

Some surveys focus not on the general overviews about semi-, self- and unsupervised learning but on special details. In their survey Cheplygina et al.

present a variety of methods in the context of medical image analysis [6]. They include deep learning and older machine learning approaches but look at different strategies from a medical perspective. Mey and Loog focused on the underlying theoretical assumptions in semi-supervised learning [32]. We keep our survey limited to general image classification tasks and focus on their practical application.

Keeping the above mentioned limitations in mind the topic of self-, semi- and unsupervised learning still includes a broad range of research fields. In this survey we will focus on deep learning approaches for image classification. We will investigate the different learning strategies with a spotlight on loss functions. Therefore, topics like metric learning and general adversarial networks will be excluded.

2. Underlying Concepts

In this section we summarize general ideas about semi-, self- and unsupervised learning. We extend this summarization by our own definition and interpretation of certain terms. The focus lies on distinguishing the possible learning strategies and the most common methods to realize them. Throughout this survey we use the terms learning strategy, technique and method in a specific meaning. The *learning strategy* is the general type/approach of an algorithm. We call each individual algorithm proposed in a paper *method*. A method can be classified to a learning strategy and consists out of *techniques*. Techniques are the parts or ideas which make up the method/algorith.

2.1. Learning strategies

Terms like supervised, semi-supervised and self-supervised are often used in literature. A precise definition which clearly separates the terms is rarely given. In most cases a rough general consensus about the meaning is sufficient but we noticed a high variety of definitions in borderline cases. For the comparison of different methods we need a precise definition to distinguish between them. We will summarize the common consensus about the learning strategies and define how we view certain borderline cases. In general, we distinguish the methods based on the amount of used labeled data and at which stage of the training process supervision is introduced. Taken together, we call the semi-, self- and unsupervised (learning) strategies *re-*

duced supervised (learning) strategies. [Figure 2](#) illustrates the four presented deep learning strategies.

2.1.1 Supervised

Supervised learning is the most common strategy in image classification with deep neural networks. We have a set of images X and corresponding labels or classes Z . Let C be the number of classes and $f(x)$ the output of a certain neural network for $x \in X$. The goal is to minimize a loss function between the outputs and labels. A common loss function to measure the difference between $f(x)$ and the corresponding label z is cross-entropy.

$$\begin{aligned} CE(f(x), z) &= \sum_{c=1}^C P_{f(x)}(c) \log(P_z(c)) \\ &= H(P_z) + KL(P_z | P_{f(x)}) \end{aligned} \quad (1)$$

P is a probability distribution over all classes. H is the entropy of a probability distribution and KL is the Kullback-Leibler divergence. The distribution P can be approximated with the output of neural network $f(x)$ or the given label z . It is important to note that cross-entropy is the sum of entropy over z and a Kullback-Leibler divergence between $f(x)$ and z . In general the entropy $H(P_z)$ is zero due to one-hot encoded label z .

Transfer Learning

A limiting factor in supervised learning is the availability of labels. The creation of these labels can be expensive and therefore limits their number. One method to overcome this limitation is to use transfer learning. Transfer learning describes a two stage process of training a neural network. The first stage is to train with or without supervision on a large and generic dataset like ImageNet [26]. The second stage is using the trained weights and fine-tune them on the target dataset. A great variety of papers have shown that transfer learning can improve and stabilize the training even on small domain-specific datasets [40].

2.1.2 Unsupervised

In unsupervised learning we only have images X and no further labels. A variety of loss functions exist

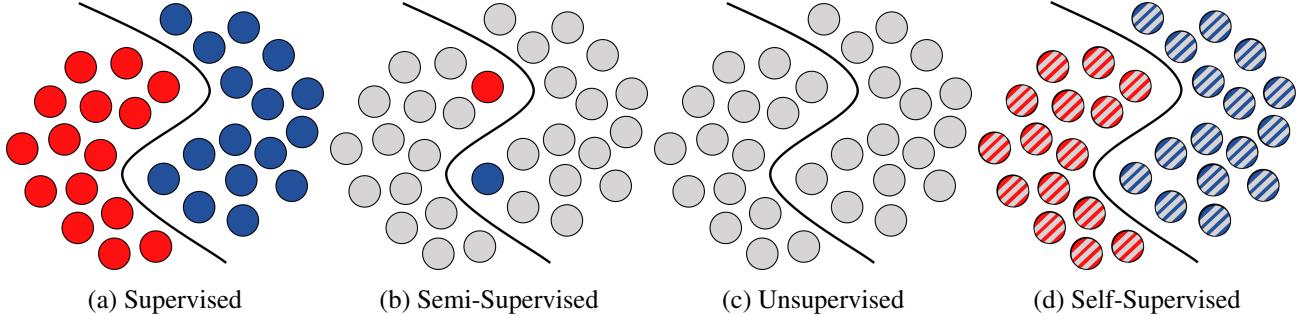


Figure 2: Illustrations of the four presented deep learning strategies - The red and dark blue circles represent labeled data points of different classes. The light grey circles represent unlabeled data points. The black lines define the underlying decision boundary between the classes. The striped circles represent datapoints which ignore and use the label information at different stages of the training process.

in unsupervised learning [5, 21, 45]. In most cases the problem is rephrased in such a way that all inputs for the loss can be generated, e.g. reconstruction loss in auto encoders [45]. Despite this automation or self-supervision we do call these methods unsupervised. Please see below for our interpretation of self-supervised learning.

2.1.3 Semi-Supervised

Semi-supervised learning is a mixture of unsupervised and supervised learning. We have labels Z for a set of images X_l like in supervised learning. The rest of the images X_u have no corresponding label. Due to this mixture, a semi-supervised loss can have a variety of shapes. A common way is to add a supervised and an unsupervised loss. In contrast to other learning strategies X_u and X_l are used in parallel.

2.1.4 Self-supervised

Self-supervised uses a pretext task to learn representations on unlabeled data. The pretext task is unsupervised but the learned representations are often not directly usable for image classification and have to be fine-tuned. Therefore, self-supervised learning can be interpreted either as an unsupervised, a semi-supervised or a strategy of its own. We see self-supervised learning as a special strategy. In the following, we will explain how we arrive at such a conclusion. The strategy cannot be called unsupervised if we need to use any labels during the fine-tuning. There is also a clear difference to semi-supervised methods.

The labels are not used simultaneously with unlabeled data because the pretext task is unsupervised and only the fine-tuning uses labels. For us this separation of the usage of labeled data into two different subtasks characterizes a strategy on its own.

2.2 Techniques

Different techniques can be used to train models in reduced supervised cases. In this section we present a selection of techniques that are used in multiple methods in the literature.

2.2.1 Consistency regularization

A major line of research uses consistency regularization. In a semi-supervised learning process these regularizations are used as an additional loss to a supervised loss on the unsupervised part of the data. This constraint leads to improved results due to the ability of taking unlabeled data into account for defining the decision boundaries [42, 28, 49]. Some self- or unsupervised methods take this approach even further by using only this consistency regularization for the training [21, 2].

Virtual Adversarial Training (VAT)

VAT [34] tries to make predictions invariant to small transformations by minimizing the distance between an image and a transformed version of the image. Miyato et al. showed how a transformation can be chosen and approximated in an adversarial way. This adversarial transformation maximizes the distance between

an image and a transformed version of it over all possible transformations. [Figure 3](#) illustrates the concept of VAT. The loss is defined as

$$\begin{aligned} VAT(f(x)) &= D(P_{f(x)}, P_{f(x+r_{adv})}) \\ r_{adv} &= \operatorname{argmax}_{r: \|r\| \leq \epsilon} D(P_{f(x)}, P_{f(x+r)}) \end{aligned} \quad (2)$$

In this equation x is an image out of the dataset X and $f(x)$ is the output for a given neural network. P is the probability distribution over these outputs and D is a non-negative function that measures the distance. Two examples of used distance measures are cross-entropy [34] and Kullback-Leiber divergence [49, 46].

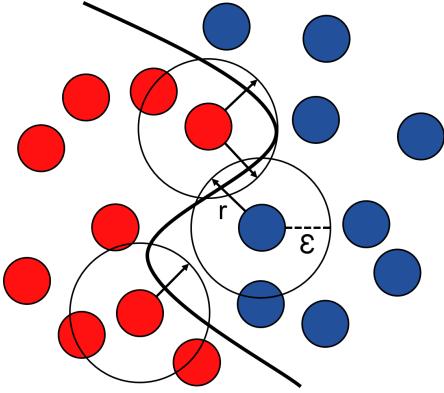


Figure 3: Illustration of the VAT concept - The blue and red circles represent two different classes. The line is the decision boundary between these classes. The ϵ spheres around the circles define the area of possible transformations. The arrows represent the adversarial change r which push the decision boundary away from any data point.

Mutual Information (MI)

MI is defined for two probability distributions as the Kullback Leiber (KL) divergence between the joint distribution and the marginal distributions [8]. This measure is used as a loss function instead of CE in several methods [19, 21, 2]. The benefits are described below. For images x, y , certain neural network outputs $f(x), f(y)$ and the corresponding probability distributions $P_{f(x)}, P_{f(y)}$, we can maximize the mutual infor-

mation by minimizing the following:

$$\begin{aligned} -I(P_{f(x)}, P_{f(y)}) &= -KL(P_{(f(x), f(y))} | P_{f(x)} * P_{f(y)}) \\ &= -H(P_{f(x)}) + H(P_{f(x)} | P_{f(y)}) \end{aligned} \quad (3)$$

An alternative representation of mutual information is the separation in entropy $H(P_{f(x)})$ and conditional entropy $H(P_{f(x)} | P_{f(y)})$.

Ji et al. describe the benefits of using MI over CE in unsupervised cases [21]. One major benefit is the inherent property to avoid degeneration due to the separation in entropy and conditional entropy. MI balances the effects of maximizing the entropy with a uniform distribution for $P_{f(x)}$ and minimizing the conditional entropy by equalizing $P_{f(x)}$ and $P_{f(y)}$. Both cases are undesirable for the output of a neural network.

Entropy Minimization (EntMin)

Grandvalet and Bengio proposed to sharpen the output predictions in semi-supervised learning by minimizing entropy [15]. They minimized the entropy $H(P_{f(x)})$ for all probability distributions $P_{f(x)}$ based on a certain neural output $f(x)$ for an image x . This minimization only sharpens the predictions of a neural network and cannot be used on its own.

Mean Squared Error (MSE)

A common distance measure between two neural network outputs $f(x), f(y)$ for images x, y is MSE. Instead of measuring the difference based on probability theory it uses the euclidean distance of the output vectors

$$MSE(f(x), f(y)) = \|f(x) - f(y)\|_2^2 \quad (4)$$

The minimization of this measure can contract two outputs to each other.

2.2.2 Overclustering

Normally, if we have k classes in a supervised case we use also k clusters in an unsupervised case. Research showed that it can be beneficial to use more clusters than actual classes k exist [4, 21]. We call this idea *overclustering*.

Overclustering can be beneficial in reduced supervised

cases due to the effect that neural networks can decide ‘on their own’ how to split the data. This separation can be helpful in noisy data or with intermediate classes that were sorted into adjacent classes randomly.

2.2.3 Pseudo-Labels

A simple approach for estimating labels of unknown data are Pseudo-Labels [29]. Lee proposed to predict classification for unseen data with a neural network and **use the predictions as labels**. What sounds at first like a self-fulfilling assumption works reasonably well in real world image classification tasks. Several modern methods are based on the same core idea of creating labels by predicting them on their own [42, 3].

3. Methods

In the following, we give a short overview over all methods in this survey in an alphabetical order and separated according to their learning strategy. Due to the fact that they may reference each other you may have to jump to the corresponding entry if you would like to know more. This list does not claim to be complete. We included methods which were referenced often in related work, which are comparable to the other methods and which are complementary to presented methods.

3.1. Semi-Supervised

Fast-Stochastic Weight Averaging (fast-SWA)

In contrast to other semi-supervised methods Athiwaratkun et al. do not change the loss but the optimization algorithm [1]. They analysed the learning process based on ideas and concepts of SWA [20], π -model [28] and Mean Teacher [42]. Athiwaratkun et al. show that averaging and cycling learning rates are beneficial in semi-supervised learning by stabilizing the training. They call their improved version of SWA fast-SWA due to a faster convergence and lower performance variance [1]. The architecture and loss is either copied from π -model [28] or Mean Teacher [42].

Mean Teacher

With Mean Teacher Tarvainen & Valpola present a student-teacher-approach for semi-supervised learning

[42]. They develop their approach based on the π -model and Temporal Ensembling [28]. Therefore, they also use MSE as a consistency loss between two predictions but create these predictions differently. They argue that Temporal Ensembling incorporates new information too slowly into predictions. The reason for this is **that the exponential moving average (EMA) is only updated once per epoch**. Therefore, they propose to use a teacher based on average weights of a student in each update step. For their model Tarvainen & Valpola show that the KL-divergence is an inferior consistency loss in comparison to MSE. An illustration of this method is given in [Figure 4](#).

MixMatch

MixMatch [3] uses a combination of a supervised and an unsupervised loss. Berthelot et al. use CE as the supervised loss and MSE between predictions and generated Pseudo-Labels as their unsupervised loss. These Pseudo-Labels are created from previous predictions of augmented images. They propose a novel sharpening method over multiple predictions to improve the quality of the Pseudo-Labels. Furthermore, they extend the algorithm mixup [50] to semi-supervised learning by incorporating the generated labels. Mixup creates convex combinations of images by blending them into each other. An illustration of the concept is given in [Figure 5](#). The prediction of the convex combination of the corresponding labels turned out to be beneficial for supervised learning in general [50].

π -model and Temporal Ensembling

Laine & Aila present two similar learning methods with the names π -model and Temporal Ensembling [28]. Both methods use a combination of the supervised CE loss and the unsupervised consistency loss MSE. The first input for the consistency loss in both cases is the output of their network from a randomly augmented input image. The second input is different for each method. In the π -model a augmentation of the same image is used. In Temporal Ensembling an exponential moving average of previous predictions is evaluated. Laine & Aila show that Temporal Ensembling is up to two times faster and and more stable in comparison to the π -model [28]. Illustrations of these methods are given in [Figure 4](#).

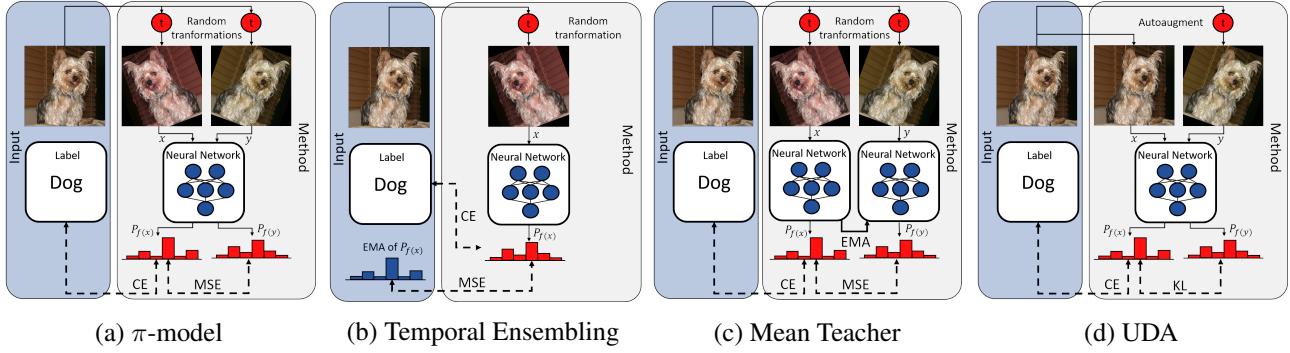


Figure 4: Illustration of four selected semi-supervised methods - The used method is given below each image. The input is given in the blue box on the left side. On the right side an illustration of the method is provided. In general the process is organized from top to bottom. At first the input images are preprocessed by none or two different random transformations. Autoaugment [9] is a special augmentation technique. The following neural network uses these preprocessed images (x, y) as input. The calculation of the loss (dotted line) is different for each method but shares common parts. All methods use the cross-entropy (CE) between label and predicted distribution $P_{f(x)}$ on labeled examples. All methods also use a consistency regularization between different predicted output distributions ($P_{f(x)}, P_{f(y)}$). The creation of these distributions differ for all methods and the details are described in the corresponding entry in [section 3](#). EMA means exponential moving average. The other abbreviations are defined above in [subsection 2.2](#).

$$0.4 \times \begin{array}{c} \text{Cat: 1.0} \\ \text{Dog: 0.0} \end{array} + 0.6 \times \begin{array}{c} \text{Cat: 0.0} \\ \text{Dog: 1.0} \end{array} = \begin{array}{c} \text{Cat: 0.4} \\ \text{Dog: 0.6} \end{array}$$

Figure 5: Illustration of mixup - The images of a cat and a dog are combined with a parametrized blending. The labels are also combined by the same parametrization. The shown images are taken from the dataset STL-10 [7]

Pseudo-Labels

Pseudo-Labels [29] describes a common technique in deep learning and a learning method on its own. For the general technique see above in [subsection 2.2](#). In contrast to many other semi-supervised methods Pseudo-Labels does not use a combination of an unsupervised and a supervised loss. The Pseudo-Labels approach uses the predictions of a neural network as labels for unknown data as described in the general technique. Therefore, the labeled and unlabeled data are used in parallel to minimize the CE loss. The usage of the same loss is a difference to other semi-supervised

methods, but the parallel utilization of labeled and unlabeled data classifies this method as semi-supervised.

Self-Supervised Semi-Supervised Learning (S⁴L)

S⁴L [49] is, as the name suggests, a combination of self-supervised and semi-supervised methods. Zhai et al. split the loss in a supervised and an unsupervised part. The supervised loss is CE while the unsupervised loss is based on the self-supervised techniques using rotation and exemplar prediction [14, 12]. The authors show that their method performs better than other self-supervised and semi-supervised techniques [12, 14, 34, 15, 29]. In their *Mix Of All Models* (MOAM) they combine self-supervised rotation prediction, VAT, entropy minimization, Pseudo-Labels and fine-tuning into a single model with multiple training steps. We count S⁴L as a semi-supervised method due to this combination.

Unsupervised Data Augmentation (UDA)

Xie et al. present with UDA a semi-supervised learning algorithm which concentrates on the usage of state-of-the-art augmentation [46]. They use a supervised and an unsupervised loss. The supervised loss is CE while the unsupervised loss is the Kullback Leiber di-

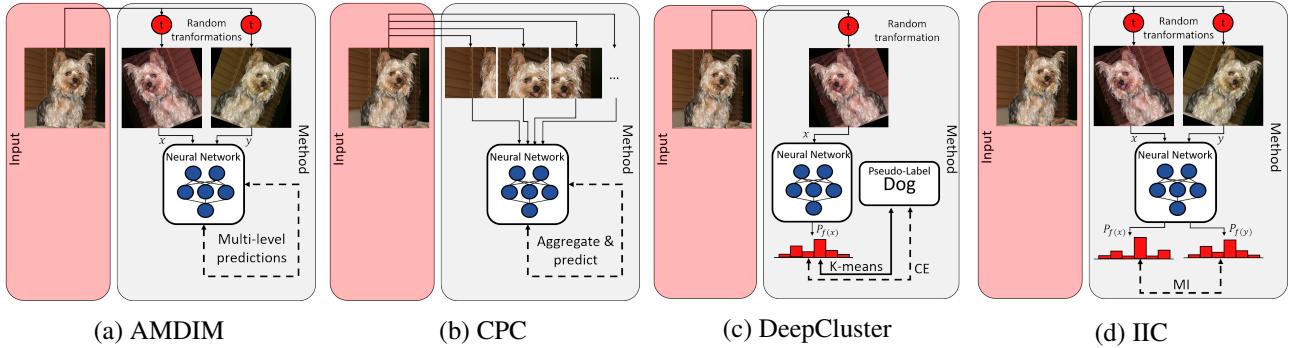


Figure 6: Illustration of four selected self-supervised methods - The used method is given below each image. The input is given in the red box on the left side. On the right side an illustration of the method is provided. The fine-tuning part is excluded. In general the process is organized from top to bottom. At first the input images are either preprocessed by one or two random transformations or are split up. The following neural network uses these preprocessed images (x, y) as input. The calculation of the loss (dotted line) is different for each method. AMDIM and CPC use internal elements of the network to calculate the loss. DeepCluster and IIC use the predicted output distribution ($P_{f(x)}, P_{f(y)}$) to calculate a loss. For further details see the corresponding entry in section 3.

vergence between output predictions. These output predictions are based on an image and an augmented version of this image. For image classification they propose to use the augmentation scheme generated by AutoAugment [9] in combination with Cutout [10]. AutoAugment uses reinforcement learning to create useful augmentations automatically. Cutout is an augmentation scheme where randomly selected regions of the image are masked out. Xie et al. show that this combined augmentation method achieves higher performance in comparison to previous methods on their own like Cutout, Cropping or Flipping. In addition to the different augmentation they propose to use a variety of other regularization methods. They proposed Training Signal Annealing which restricts the influence of labeled examples during the training process in order to prevent overfitting. They use EntMin [15] and a kind of Pseudo-Labeling [29]. We use a kind of Pseudo-Labeling because they do not use the predictions as labels but they use them to filter unsupervised data for outliers. An illustration of this method is given in Figure 4.

Virtual Adversarial Training (VAT)

VAT [34] is not just the name for a regularization technique but it is also a semi-supervised learning method. Miyato et al. used a combination of VAT on unlabeled data and CE on labeled data [34]. They showed that the

adversarial transformation leads to a lower error on image classification than random transformations. Furthermore, they proved that adding EntMin [15] to the loss increased accuracy even more.

3.2. Self-Supervised

Augmented Multiscale Deep InfoMax (AMDIM)

AMDIM [2] maximizes the MI between inputs and outputs of a network. It is an extension of the method DIM [18]. DIM usually maximizes MI between local regions of an image and a representation of the image. AMDIM extends the idea of DIM in several ways. Firstly, the authors sample the local regions and representations from different augmentations of the same source image. Secondly, they maximize MI between multiple scales of the local region and the representation. They use a more powerful encoder and define mixture-based representations to achieve higher accuracies. Bachman et al. fine-tune the representations on labeled data to measure their quality. An illustration of this method is given in Figure 6.

Contrastive Predictive Coding (CPC)

CPC [43, 17] is a self-supervised method which predicts representations of local image regions based on previous image regions. The authors determine the quality of these predictions by identifying the correct prediction out of randomly sampled negative ones.

They use standard CE loss and adopt the loss to the summation over the complete image. This adaption results in their loss InfoNCE [43]. Van den Oord et al. showed that minimizing InfoNCE maximizes the lower bound for MI between the previous image regions and the predicted image region [43]. An illustration of this method is given in [Figure 6](#).

DeepCluster

DeepCluster [4] is a self-supervised technique which generates labels by k-means clustering. Caron et al. iterate between clustering of predicted labels to generate Pseudo-Labels and training with cross-entropy on these labels. They show that it is beneficial to use over-clustering in the pretext task. After the pretext task they fine-tune the network on all labels. An illustration of this method is given in [Figure 6](#).

Deep InfoMax (DIM)

DIM [18] maximizes the MI between local input regions and output representations. Hjelm et al. show that maximizing over local input regions rather than the complete image is beneficial for image classification. In addition, they use a discriminator to match the output representations to a given prior distribution. At the end they fine-tune the network with an additional small fully-connected neural network.

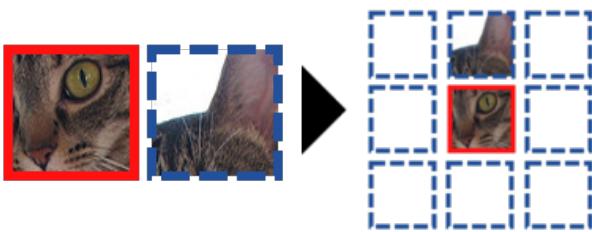


Figure 7: Illustration of the Context pre-text task - A central patch and an adjacent patch from the same image are given. The task is to predict one of the 8 possible relative positions of the second patch to the first one. In the example the correct answer is upper center. The illustration is inspired by [11].

Invariant Information Clustering (IIC)

IIC [21] maximizes the MI between augmented views of an image. The idea is that images should belong to the same class regardless of augmentation. The

augmentation has to be a transformation to which the neural network should be invariant. The authors do not maximize directly over the output distributions but over the class distribution which is approximated for every batch. Ji et al. use auxiliary overclustering on a different output head to increase their performance in the unsupervised case. This idea allows the network to learn subclasses and handle noisy data. Ji et al. use Sobel filtered images as input instead of the original RGB images. Additionally, they show how to extend IIC to image segmentation. Up to this point the method is completely unsupervised. In order to be comparable to other semi-supervised methods they fine-tune their models on a subset of available labels. An illustration of this method is given in [Figure 6](#).

Representation Learning - Context

Doersch et al. propose to use context prediction as a pretext task for visual representation learning [11]. A central patch and an adjacent patch from an image are used as input. The task is to predict one of the 8 possible relative positions of the second patch to the first one. An illustration of the pretext task is given in [Figure 7](#). Doersch et al. argue that this task becomes more easy if you recognize the content of these patches. The authors fine-tune their representations for other tasks and show their superiority in comparison to random initialization. Aside from fine-tuning, Doersch et al. show how their method could be used for Visual Data Mining.

Representation Learning - Exemplar

Dosovitskiy et al. were one of the first to propose a self-supervised pre-text task with additional fine-tuning [12]. They randomly sample patches from different images and augment these patches heavily. Augmentations can be for example rotations, translations, color changes or contrast adjustments. The classification task is to map all augmented versions of a patch to the correct original patch.

Representation Learning - Jigsaw

Noroozi and Favaro propose to solve Jigsaw puzzles as a pre-text task [35]. The idea is that a network has to understand the concept of a presented object in order to solve the puzzle. They prevent simple solutions which

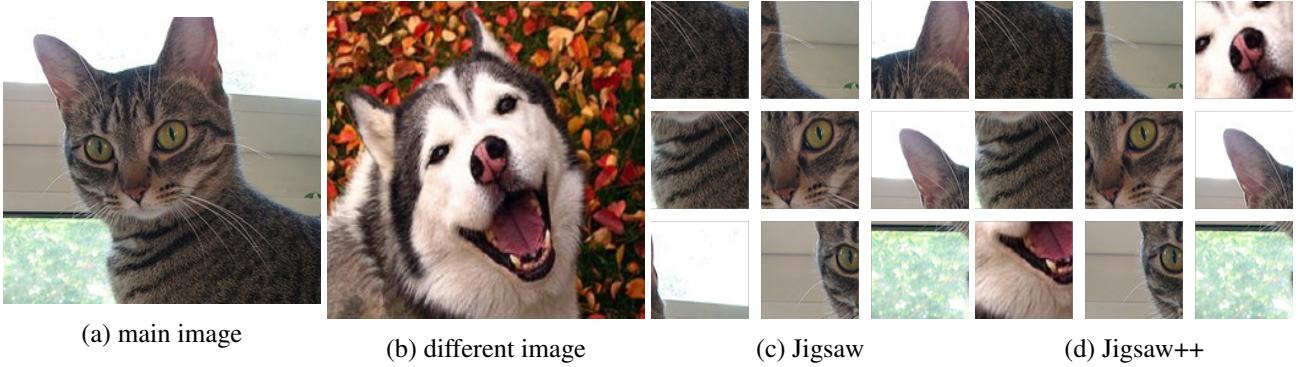


Figure 8: Illustrations of the pretext task Jigsaw and Jigsaw++ - The Jigsaw pretext task consists of solving a simple Jigsaw puzzle generated from the main image. Jigsaw++ augments the Jigsaw puzzle by adding in parts of a different image. The illustrations are inspired by [36].

only look at edges or corners by including small random margins between the puzzle patches. They fine-tune on supervised data for image classification tasks. Noroozi et al. extended the Jigsaw task by adding image parts of a different image [36]. They call the extension Jigsaw++. An example for Jigsaw and Jigsaw++ is given in Figure 8.

Representation Learning - Rotation

Gidaris et al. use a pretext task based on image rotation prediction [14]. They propose to randomly rotate the input image by 0, 90, 180 or 270 degrees and let the network predict the chosen rotation degree. In their work they also evaluate different numbers of rotations but four rotations score the best result. For image classification they fine-tune on labeled data.

3.3. Unsupervised

Deep Adaptive Image Clustering (DAC)

DAC [5] reformulates unsupervised clustering as a pairwise classification. Similar to the idea of Pseudo-Labels Chang et al. predict clusters and use these to retrain the network. **The twist is that they calculate the cosine distance between all cluster predictions.** This distance is used to determine whether the input images are similar or dissimilar with a given certainty. The network is then trained with binary CE on these certain similar and dissimilar input images. During the training process they lower the needed certainty to include more images. As input Chang et al. use a combination of RGB and extracted HOG features. Additionally, they use an auxiliary loss in their source code

which is not reported in the paper.¹

Invariant Information Clustering (IIC)

IIC [21] is described above as a self-supervised learning method. In comparison to other presented self-supervised methods IIC creates usable classifications without fine-tuning the model on labeled data. The reason for this is that the pretext task is constructed in such a way that label predictions can be extracted directly from the model. This leads to the conclusion that IIC can also be interpreted as an unsupervised learning method.

Information Maximizing Self-Augmented Training (IMSAT)

IMSAT [19] maximizes MI between the input and output of the model. As a consistency regularization Hu et al. use CE between an image prediction and an augmented image prediction. They show that the best augmentation of the prediction can be calculated with VAT [34]. The maximization of MI directly on the image input leads to a problem. For datasets like CIFAR-10, CIFAR-100 [25] and STL-10 [7] the color information is too dominant in comparison to the actual content or shape. As a workaround Hu et al. use the features generated by a pretrained CNN on ImageNet [26] as input.

¹<https://github.com/vector-1127/DAC/blob/master/STL10/stl.py>



(a) CIFAR-10

(b) STL-10

(c) ILSVRC-2012

Figure 9: Examples of four random cats in the different datasets to illustrate the difference in quality

4. Comparison

In this chapter we will analyze which techniques are shared or different between methods. We will compare the performance of all methods with each other on common deep learning datasets.

4.1. Datasets

In this survey we compare the presented methods on a variety of datasets. We selected four datasets that were used in multiple papers to allow a fair comparison. An overview of example images is given in Figure 9.

CIFAR-10 and CIFAR-100 are large datasets of tiny color images with size 32x32 [25]. Both datasets contain 60,000 images belonging to 10 or 100 classes respectively. The 100 classes in CIFAR-100 can be combined into 20 superclasses. Both sets provide 50,000 training labels and 10,000 validation labels. The presented results are only trained with 4,000 labels for CIFAR-10 and 10,000 labels for CIFAR-100 to represent a semi-supervised case. If a method uses all labels this is marked independently.

STL-10 is dataset designed for unsupervised and semi-supervised learning [7]. The dataset is inspired by CIFAR-10 [25] but provides less labels. It only consists of 5,000 training labels and 8,000 validation labels. However, 100,000 unlabeled example image are also provided. These unlabeled examples belong

to the training classes and some different classes. The images are 96x96 color images and were acquired in combination with their labels from ImageNet [26].

ILSVRC-2012 is a subset of ImageNet [26]. The training set consists of 1.2 million images while the validation and the test set include 150,000 images. These images belong to 1000 object categories. Due to this large number of categories, it is common to report Top-5 and Top-1 accuracy. Top-1 accuracy is the classical accuracy where one prediction is compared to one ground truth label. Top-5 accuracy checks if a ground truth label is in a set of at most five predictions. For further details on accuracy see subsection 4.2. The presented results are only trained with 10% of labels to represent a semi-supervised case. If a method uses all labels this is marked independently.

4.2. Evaluation metrics

We compare the performance of all methods based on their classification score. This score is defined differently for unsupervised and all other settings. We follow standard protocol and use the classification accuracy in most cases. For unsupervised strategies we use the cluster accuracy because we need to handle the missing labels during the training. We need to find the best one-to-one permutations from the network cluster predictions to the ground-truth classes.

For vectors $x, y \in \mathbb{Z}^N$ with $N \in \mathbb{N}$ the accuracy is

defined as follows:

$$ACC(x, y) = \frac{\sum_{i=1}^N \mathbb{1}_{y_i=x_i}}{N} \quad (5)$$

For the cluster accuracy we additionally maximize over all possible one-to-one permutations σ .

$$ACC(x, y) = \max_{\sigma} \frac{\sum_{i=1}^N \mathbb{1}_{y_i=\sigma(x_i)}}{N} \quad (6)$$

4.3. Comparison of methods

In this subsection we will compare the methods with regard to their used techniques and performance. We will summarize the presented results and discuss the underlying trends in the next subsection.

Comparison with regard to used techniques

In [Table 1](#) we present all methods and their used techniques. We evaluate only techniques which were used frequently in different papers. Special details such as the mixup algorithm in MixMatch, the different optimizer for fast-SWA or the used distance measure for VAT are excluded. Please see [section 3](#) for further details.

Most methods share similar supervised and unsupervised techniques. All semi-supervised methods use the cross-entropy loss during training. All self-supervised methods use a pretext task and fine-tune in the end. All unsupervised methods do not use any technique which requires ground-truth labels. Due to our definition of the learning strategies this grouping is expected.

The similarities and differences are not that clear for unsupervised techniques. We still step through all techniques based on the significance for differentiating the learning strategies and methods.

MI is not used by any semi-supervised method and most semi-supervised methods do not use a pretext task. Self-supervised methods often rely only on one or both of these techniques. Most unsupervised methods use also MI but do not use a pretext task. S^4L and IIC stand out due to the fact that they use a pretext task. S^4L is the only method which combines a self-supervised pretext task and semi-supervised learning. IIC uses a pretext task in the unsupervised case based on MI. The pretext task creates representations which can be interpreted as classifications without

fine-tuning.

The techniques VAT, EntMin and MSE are often used for semi-supervised methods. VAT and EntMin are more often used together than MSE with either of them. This correlation might be introduced by the selected methods. The methods fast-SWA, Mean Teacher, π -model and Temporal Ensembling are extensions or derivations of each other. Also VAT+EntMin is an extension of VAT.

Pseudo-Labels are used in several different methods. Due to the simple and flexible idea of Pseudo-Labels, it can be used in a variety of different methods. UDA for example uses this technique to filter the unlabeled data for useful images.

Comparison with regard to performance

We compare the performance of the different methods based on their respective reported results or cross-references in other papers. For a better comparability we would have liked to recreate every method in a unified setup but this was not feasible. While using reported values might be the only possible approach, it leads to drawbacks in the analysis.

Kolesnikov et al. showed that changes in the architecture can lead to significant performance boost or drops [\[24\]](#). They state that 'neither [...] the ranking of architectures [is] consistent across different methods, nor is the ranking of methods consistent across architectures' [\[24\]](#). While most methods try to achieve comparability with previous ones by a similar setup, over time small differences still aggregate and lead to a variety of used architectures. Some methods use only early convolutional networks such as AlexNet [\[27\]](#) but others use more modern architectures like Wide ResNet-Architecture [\[48\]](#) or Shake-Shake-Regularization [\[13\]](#). Oliver et al. proposed guidelines to ensure more comparable evaluations in semi-supervised learning [\[37\]](#). They showed that not following these guidelines may lead to changes in the performance [\[37\]](#). While some methods try to follow these guidelines, we cannot guarantee that all methods do so. This impacts the comparability further.

[Table 2](#) shows the collected results for all presented methods. We also provide results for the respective supervised baselines reported by the authors. To keep a fair comparability we did not add state-of-the-art baselines with more complex architectures.

Table 1: Overview of the methods and their used techniques - On the left-hand side the reviewed methods from [section 3](#) are sorted by the learning strategy. The top row lists the possible techniques which have been discussed in [subsection 2.2](#). The techniques are sorted into unsupervised and supervised based on the question if they can be used with or without labeled data. The abbreviations for the techniques are also given in [subsection 2.2](#). Cross-entropy (CE) describes the usage of CE as part of the loss during training. Fine-tuning (FT) describes the usage of cross-entropy for new labels after the initial training e.g. on a pretext task. (X) means that the technique is not used directly but indirectly. The individual explanations are given by the indicated number. 1 - MixMatch does entropy minimization implicitly by sharpening the predictions [3]. 2 - UDA predicts Pseudo-Labels for filtering the unsupervised data. 3 - Minimize mutual information objective as a pretext task e.g. between views [2] or layers [17]. 4 - The loss InfoNCE maximizes the mutual information indirectly [43]. 5 - Deep Cluster uses K-Means to calculate Pseudo-Labels and optimizes the assignment as a pretext task. 6 - DAC uses the cosine distance between elements to estimate similar and dissimilar items. One could say DAC creates Pseudo-Labels for the similarity problem.

	Unsupervised Techniques					Pretext Task	Supervised Techniques	
	VAT	EntMin	MI	MSE	Pseudo		CE	Fine-tuning
Semi-Supervised Methods								
fast-SWA [1]				X			X	
Mean Teacher [42]				X			X	
MixMatch [3]			(X) ¹	X	X		X	
π model [28]				X			X	
Pseudo-Labels [29]					X		X	
S ⁴ L [49]	X	X			X	Rotation	X	X
Temporal Ensembling [28]				X			X	
UDA [46]	X	X			(X) ²		X	
VAT [34]	X						X	
VAT + EntMin [34]	X	X					X	
Self-Supervised Methods								
AMDIM [2]			X		(X) ³		X	
Context [11]					Context		X	
CPC [43, 17]			(X) ⁴		(X) ³		X	
DeepCluster [4]				(X) ⁵	(X) ⁵		X	
DIM [18]			X		(X) ³		X	
Exemplar [12]					Augmentation		X	
IIC [21]			X		(X) ³		X	
Jigsaw [35]					Jigsaw		X	
Rotation [14]					Rotation		X	
Unsupervised Methods								
DAC [5]				(X) ⁶				
IIC [21]			X		(X) ³			
IMSAT [19]	X		X					

Considering the above mentioned limitations, we do not focus on small differences but look for general trends and specialities instead.

In general, the used architectures become more complex and the accuracies rise over time. This behavior is expected as new results are often improvements of earlier works. The changes in architecture may have led to the improvements. However, many papers include

ablation studies and comparisons to only supervised methods to show the impact of their method. We believe that a combination of more modern architecture and more advanced methods leads to the improvement. For the CIFAR-10 dataset almost all semi- and self-supervised methods reach about or over 90% accuracy. The best method MixMatch reaches an accuracy of about 95% and is roughly two percent worse than the

Table 2: Overview of the reported accuracies - The first column states the used method. For the supervised baseline we used the best reported results which were considered as baseline from the other methods. The original paper is given in brackets after the score. The architecture and their reference is given in the second column. The third column shows the year of publication or the release year of the preprint. The last four columns report the Top-1 accuracy score in % for the respective dataset (See [subsection 4.2](#) for further details). If the results are not reported in the original paper the reference is given after the result. A blank entry represents the fact that no result was reported. Result for Top-5 accuracies are marked with \star . Normally a fully connected layer is used for fine-tuning. If multiple layers are used as a multilayer perceptron the results are marked with \ddagger . The amount of annotated label for each dataset is defined in [subsection 4.1](#). If all labels were used to calculate the result they are marked with \dagger . 1 - Architecture includes Shake-Shake regularization. 2 - Input is 4x wider. 3 - It uses ten random classes out of the default 1000 classes. 4 - It only predicts 20 superclasses instead of the default 100 classes. 5 - Inputs are pretrained ImageNet features.

	Architecture	Publication	CIFAR-10	CIFAR-100	STL-10	ILSVRC-2012
Supervised (100% labels)	Best reported	-	97.14[42]	79.82[2]	68.7 [18]	78.57 [49] / 94.10 \star [49]
Semi-Supervised Methods						
fast-SWA [1]	CONV-13 [28]	2019	90.95	66.38		
fast-SWA [1]	ResNet-26 ¹ [13]	2019	93.72			
Mean Teacher [42]	CONV-13 [28]	2017	87.69			
Mean Teacher [42]	ResNet 26 and 152	2017	93.72			90.89 \star
MixMatch [3]	Wide ResNet-28 [48]	2019	95.05	74.12	94.41	
π model [28]	CONV-13 [28]	2017	87.64			
Pseudo-Label [29]	ResNet50v2 ² [16]	2013				82.41 \star [49]
S ⁴ L [49]	ResNet50v2 ² [24]	2019				73.21 / 91.23 \star
Temporal Ensembling [28]	CONV-13 [28]	2017	87.84			
UDA [46]	Wide ResNet-28 [48]	2019	94.7			68.66 / 88.52 \star
VAT [34]	CONV-13 [28]	2018	88.64			
VAT [34]	ResNet50v2 [16]	2018				82.78 \star [49]
VAT [34] + EntMin [15]	CONV-13 [28]	2018	89.45 [34]			
VAT [34] + EntMin [15]	ResNet50v2 [16]	2018	86.41 [49]			83.39 \star [49]
Self-Supervised Methods						
AMDIM [2]	ResNet18 [16]	2019	91.3 \dagger / 93.6 \ddagger	70.2 \dagger / 73.8 \ddagger	93.6 / 93.8 \ddagger	60.2 \dagger / 60.9 \ddagger
Context [11]	ResNet50 [16]	2015				51.4 \dagger [24]
CPC [43, 17]	ResNet-170 [17]	2019	77.45 \dagger [18]		77.81 \dagger [18]	61.0 / 84.88 \star
DeepCluster [4]	AlexNet [27]	2018			73.4 [21]	41 \dagger
DIM [18]	AlexNet [27]	2019			72.57 \ddagger	
DIM [18]	GAN Discriminator [39]	2019	75.21 \ddagger	49.74 \ddagger		
Exemplar [12]	ResNet50 [16]	2016				46.0 \dagger [24] / 81.01 \star [49]
IIC [21]	ResNet34 [16]	2019			88.8	
Jigsaw [35]	AlexNet [27]	2016				44.6 \dagger [24]
Rotation [14]	AlexNet [27]	2018				55.4 \dagger [24]
Rotation [14]	ResNet50v2 [16]	2018				78.53 \star [49]
Unsupervised Methods						
DAC [5]	All-ConvNet [41]	2017	52.18	23.75	46.99	52.72 ³
IIC [21]	ResNet34 [16]	2019	61.7	25.7 ⁴	61.0	
IMSAT [19]	Autoencoder ⁵	2017	45.6	27.5	94.1	

fully supervised baseline. For the CIFAR-100 dataset fewer results are reported. MixMatch is with about 74% on this dataset the best method in comparison to the fully supervised baseline of about 80%.

For the STL-10 dataset most methods report a better

result than the supervised baseline. These results are possible due to the unlabeled part of the dataset. The unlabeled data can only be utilized by semi-, self- or unsupervised methods. MixMatch achieves the best results with about 94%.

The ILSVRC-2012 dataset is the most difficult dataset based on the reported Top-1 accuracies. Most methods achieve only a Top-1 accuracy which is roughly 20% worse than the reported supervised baseline with around 79%. Only the methods S⁴L and UDA achieve an accuracy which is less than 10% worse than the baseline. S⁴L achieves the best accuracy with a Top-1 accuracy of about 73% and a Top-5 accuracy of around 92%.

The unsupervised methods are separated from the supervised baseline by a clear margin of up to 50%. IIC achieves the best results of about 61% on CIFAR-10 and STL-10. IMSAT reports an accuracy of about 94% on STL-10. Due to the fact that IMSAT uses pretrained ImageNet features, a superset of STL-10, the results are not directly comparable.

4.4. Discussion

In this subsection we discuss the presented results of the previous subsection. We divide our discussion into three major trends which we identified. All these trends lead to possible future research opportunities.

1. Trend: Real World Applications

Previous methods were not scalable to real world images and applications and used workarounds e.g. extracted features [19] to process real world images. Many methods can report a result of over 90% on CIFAR-10, a simple low-resolution dataset. Only two methods are able to achieve a Top-5 accuracy of over 90% on ILSVRC-2012, a high-resolution dataset. We conclude that most methods are not scalable to real world image classification problems. However, the best reported methods like MixMatch and S⁴L surpassed the point of only scientific usage and can be applied to real world applications.

This conclusion applies to real world image classification tasks with balanced and clearly separated classes. This conclusion also implicates which real world issues need to solved in future research. Class imbalance or noisy labels are not treated by the presented methods. Datasets with also few unlabeled data points are not considered.

2. Trend: Needed supervision is decreasing

We see that the gap between reduced supervised and supervised methods is shrinking. For CIFAR-10,

CIFAR-100 and ILSVRC-2012 we have a gap of less than 5% left between total supervised and reduced supervised learning. For STL-10 the reduced supervised methods even surpass the total supervised case by about 30% due to the additional set of unlabeled data. We conclude that reduced supervised learning reaches comparable results while using only roughly 10% of the labels.

A lot of newly proposed methods are semi- or self-supervised in comparison to unsupervised ones. Unsupervised methods like IIC still reach results of over 60% and show that this kind of training can be beneficial for semi-supervised learning [21]. However, the results are still surpassed by semi- or self-supervised methods by a large margin e.g. over 30% on CIFAR-10. The integration of the knowledge of some labels into the training process seems to be crucial.

In general we considered a reduction from 100% to 10% of all labels. Some methods [49, 3] try even to train their models with only 1% of all labels. For ILSVRC-2012 this is equivalent to about 13 images per class. We expect that future research will concentrate on achieving comparable results for only 1% or even fewer of all labels. In the end research fields like few-shot, single-shot and semi-supervised learning might even merge.

We assume that in parallel to the reduction of needed labels, the usage of total unsupervised methods will decrease further. The benefit of even some labels as guiding reference for learning methods is too valuable to ignore. We believe that for many real world applications a very low supervision in form of archetypes might be sufficient. This will lead to a shift in the corresponding research efforts.

3. Trend: Combination of techniques

In the comparison we identified that semi-supervised and self-supervised methods share few unsupervised techniques.

We believe there is only little overlap between semi- and self-supervised learning due to the different aims of the respective authors. Many self-supervised papers focus on creating good representations. They fine-tune their results only to be comparable. Semi-supervised papers aim for the best accuracy scores with as few labels as possible.

The comparison showed that MixMatch and S⁴L are

the best methods or, if we consider the above mentioned limitations due to architecture differences, one of the best. Both methods combine several techniques in their method. S^4L stands out as it combines many different techniques and the combined approach is even called "Mix of all models" [49]. We assume that this combination is the reason for their improved performance. This assumption is supported by the included comparisons in the original papers. S^4L shows the impact of each method separately as well as the combination of all [49].

IIC is the only method which can be used as an unsupervised or self-supervised method with fine-tuning. This flexibility allows approaches with a smooth transition between no and small supervision.

The comparison showed that the technique Pseudo-Labels can be applied to a variety of methods. However, only few methods use this technique.

We believe that the combination of different basic techniques is a promising future research field due to the fact that many combinations are not yet explored.

5. Conclusion

In this paper, we provided an overview over semi-, self- and unsupervised techniques. We analyzed their difference, similarities and combinations based on 21 different methods. This analysis led to the identification of several trends and possible research fields.

We based our analysis on the definition of the different learning strategies (semi-, self- or unsupervised) and common techniques in these strategies. We showed how the methods work in general, which techniques they use and as what kind of strategy they could be classified. Despite the difficult comparison of the methods' performances due to different architectures and implementations we identified three major trends. Results of over 90% Top-5 accuracy on ILSVRC-2012 with only 10% of the labels show that semi-supervised methods are applicable for real world problems. However, issues like class imbalance are not considered. Future research has to address these issues.

The performance gap between supervised and semi- or self-supervised methods is closing. For one dataset it is even surpassed by about 30%. The number of labels to get comparable results to fully supervised learning is decreasing. Future research can lower the number of needed labels even further. We noticed that, as time

progresses, unsupervised methods are used less often. These two conclusions lead us to the assumption that unsupervised methods will lose significance for real world image classification in the future.

We concluded that semi- and self-supervised learning strategies mainly use a different set of techniques. In general, both strategies use a combination of different techniques but there are few overlaps in these techniques. S^4L is the only presented method which gaps this separation. We identified the trend that a combination of different techniques is beneficial to the overall performance. In combination with the small overlap between the techniques we identified possible future research opportunities.

References

- [1] B. Athiwaratkun, M. Finzi, P. Izmailov, and A. G. Wilson. There are many consistent explanations of unlabeled data: Why you should average. In *International Conference on Learning Representations*, 2019. [6](#), [13](#), [14](#)
- [2] P. Bachman, R. D. Hjelm, and W. Buchwalter. Learning representations by maximizing mutual information across views. In *Advances in Neural Information Processing Systems*, pages 15509–15519, 2019. [4](#), [5](#), [8](#), [13](#), [14](#)
- [3] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. A. Raffel. Mixmatch: A holistic approach to semi-supervised learning. In *Advances in Neural Information Processing Systems*, pages 5050–5060, 2019. [6](#), [13](#), [14](#), [15](#)
- [4] M. Caron, P. Bojanowski, A. Joulin, and M. Douze. Deep clustering for unsupervised learning of visual features. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 132–149, 2018. [5](#), [9](#), [13](#), [14](#)
- [5] J. Chang, L. Wang, G. Meng, S. Xiang, and C. Pan. Deep adaptive image clustering. *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 5880–5888, 2017. [4](#), [10](#), [13](#), [14](#)
- [6] V. Cheplygina, M. de Bruijne, and J. P. Pluim. Not-so-supervised: A survey of semi-supervised, multi-instance, and transfer learning in medical image analysis. *Medical Image Analysis*, 54:280296, May 2019. [3](#)

- [7] A. Coates, A. Ng, and H. Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 215–223, 2011. 7, 10, 11
- [8] T. M. Cover and J. A. Thomas. *Elements of information theory*. John Wiley & Sons, 2012. 5
- [9] E. D. Cubuk, B. Zoph, D. Mane, V. Vasudevan, and Q. V. Le. Autoaugment: Learning augmentation strategies from data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 113–123, 2019. 7, 8
- [10] T. Devries and G. W. Taylor. Improved regularization of convolutional neural networks with cutout. *ArXiv*, abs/1708.04552, 2017. 8
- [11] C. Doersch, A. Gupta, and A. A. Efros. Unsupervised visual representation learning by context prediction. *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 1422–1430, 2015. 9, 13, 14
- [12] A. Dosovitskiy, P. Fischer, J. T. Springenberg, M. Riedmiller, and T. Brox. Discriminative unsupervised feature learning with exemplar convolutional neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(9):17341747, Sep 2016. 7, 9, 13, 14
- [13] X. Gastaldi. Shake-shake regularization. *arXiv preprint arXiv:1705.07485*, 2017. 12, 14
- [14] S. Gidaris, P. Singh, and N. Komodakis. Unsupervised representation learning by predicting image rotations. In *International Conference on Learning Representations*, 2018. 7, 10, 13, 14
- [15] Y. Grandvalet and Y. Bengio. Semi-supervised learning by entropy minimization. In *Advances in neural information processing systems*, pages 529–536, 2005. 5, 7, 8, 14
- [16] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 14
- [17] O. J. Hénaff, A. Razavi, C. Doersch, S. Eslami, and A. v. d. Oord. Data-efficient image recognition with contrastive predictive coding. *arXiv preprint arXiv:1905.09272*, 2019. 8, 13, 14
- [18] R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman, A. Trischler, and Y. Bengio. Learning deep representations by mutual information estimation and maximization. In *International Conference on Learning Representations*, 2019. 8, 9, 13, 14
- [19] W. Hu, T. Miyato, S. Tokui, E. Matsumoto, and M. Sugiyama. Learning discrete representations via information maximizing self-augmented training. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1558–1567. JMLR.org, 2017. 5, 10, 13, 14, 15
- [20] P. Izmailov, D. Podoprikin, T. Garipov, D. Vetrov, and A. G. Wilson. Averaging weights leads to wider optima and better generalization. In *Conference on Uncertainty in Artificial Intelligence*, 2018. 6
- [21] X. Ji, J. F. Henriques, and A. Vedaldi. Invariant information clustering for unsupervised image classification and segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9865–9874, 2019. 2, 4, 5, 9, 10, 13, 14, 15
- [22] L. Jing and Y. Tian. Self-supervised visual feature learning with deep neural networks: A survey, 2019. 2
- [23] M. Kaya and H. Ş. Bilge. Deep metric learning: a survey. *Symmetry*, 11(9):1066, 2019. 2
- [24] A. Kolesnikov, X. Zhai, and L. Beyer. Revisiting self-supervised visual representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1920–1929, 2019. 12, 14
- [25] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009. 10, 11
- [26] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. 1, 3, 10, 11
- [27] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional

- neural networks. *Commun. ACM*, 60:84–90, 2012. [12](#), [14](#)
- [28] S. Laine and T. Aila. Temporal ensembling for semi-supervised learning. In *International Conference on Learning Representations*, 2017. [4](#), [6](#), [13](#), [14](#)
- [29] D.-H. Lee. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on challenges in representation learning, ICML*, volume 3, page 2, 2013. [6](#), [7](#), [8](#), [13](#), [14](#)
- [30] J. MacQueen et al. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 281–297. Oakland, CA, USA, 1967. [2](#)
- [31] D. Mahajan, R. Girshick, V. Ramanathan, K. He, M. Paluri, Y. Li, A. Bharambe, and L. van der Maaten. Exploring the limits of weakly supervised pretraining. *Lecture Notes in Computer Science*, page 185201, 2018. [1](#)
- [32] A. Mey and M. Loog. Improvability through semi-supervised learning: A survey of theoretical results. *arXiv preprint arXiv:1908.09574*, 2019. [3](#)
- [33] E. Min, X. Guo, Q. Liu, G. Zhang, J. Cui, and J. Long. A survey of clustering with deep learning: From the perspective of network architecture. *IEEE Access*, 6:39501–39514, 2018. [2](#)
- [34] T. Miyato, S.-i. Maeda, S. Ishii, and M. Koyama. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 2018. [4](#), [5](#), [7](#), [8](#), [10](#), [13](#), [14](#)
- [35] M. Noroozi and P. Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European Conference on Computer Vision*, pages 69–84. Springer, 2016. [9](#), [13](#), [14](#)
- [36] M. Noroozi, A. Vinjimoor, P. Favaro, and H. Pirsiavash. Boosting self-supervised learning via knowledge transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9359–9367, 2018. [10](#)
- [37] A. Oliver, A. Odena, C. A. Raffel, E. D. Cubuk, and I. Goodfellow. Realistic evaluation of deep semi-supervised learning algorithms. In *Advances in Neural Information Processing Systems*, pages 3235–3246, 2018. [12](#)
- [38] G.-J. Qi and J. Luo. Small data challenges in big data era: A survey of recent progress on unsupervised and semi-supervised methods, 2019. [2](#)
- [39] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *International Conference on Learning Representations*, 2016. [14](#)
- [40] L. Schmarje, C. Zelenka, U. Geisen, C.-C. Glüer, and R. Koch. 2d and 3d segmentation of uncertain local collagen fiber orientations in shg microscopy. In *DAGM GCPR*, pages 374–386, 2019. [1](#), [3](#)
- [41] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. A. Riedmiller. Striving for simplicity: The all convolutional net. *CoRR*, abs/1412.6806, 2014. [14](#)
- [42] A. Tarvainen and H. Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *ICLR*, 2017. [4](#), [6](#), [13](#), [14](#)
- [43] A. van den Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding, 2018. [8](#), [9](#), [13](#), [14](#)
- [44] J. E. van Engelen and H. H. Hoos. A survey on semi-supervised learning. *Machine Learning*, pages 1–68, 2019. [2](#)
- [45] J. Xie, R. B. Girshick, and A. Farhadi. Unsupervised deep embedding for clustering analysis. In *ICML*, 2015. [2](#), [4](#)
- [46] Q. Xie, Z. Dai, E. Hovy, M.-T. Luong, and Q. V. Le. Unsupervised data augmentation for consistency training, 2019. [2](#), [5](#), [7](#), [13](#), [14](#)
- [47] R. Xu and D. C. Wunsch. Survey of clustering algorithms. *IEEE Transactions on Neural Networks*, 16:645–678, 2005. [2](#)
- [48] S. Zagoruyko and N. Komodakis. Wide residual networks. *Proceedings of the British Machine Vision Conference 2016*, 2016. [12](#), [14](#)

- [49] X. Zhai, A. Oliver, A. Kolesnikov, and L. Beyer. S4l: Self-supervised semi-supervised learning. In *Proceedings of the IEEE international conference on computer vision*, pages 1476–1485, 2019. [2](#), [4](#), [5](#), [7](#), [13](#), [14](#), [15](#), [16](#)
- [50] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018. [6](#)
- [51] X. Zhu. Semi-supervised learning literature survey. *Comput Sci, University of Wisconsin-Madison*, 2, 07 2008. [2](#)