

ConstraintMatch for Semi-constrained Clustering

Jann Goschenhofer^{*†‡}, Bernd Bischl^{*†‡}, Zsolt Kira[§]
LMU Munich^{*}, Fraunhofer Institute for Integrated Circuits[†]
Munich Center for Machine Learning[‡], Georgia Institute of Technology[§]

Abstract—Constrained clustering allows the training of classification models using pairwise constraints only, which are weak and relatively easy to mine, while still yielding full-supervision-level model performance. While they perform well even in the absence of the true underlying class labels, constrained clustering models still require large amounts of binary constraint annotations for training. In this paper, we propose a semi-supervised context whereby a large amount of *unconstrained* data is available alongside a smaller set of constraints, and propose *ConstraintMatch* to leverage such unconstrained data. While a great deal of progress has been made in semi-supervised learning using full labels, there are a number of challenges that prevent a naive application of the resulting methods in the constraint-based label setting. Therefore, we reason about and analyze these challenges, specifically 1) proposing a *pseudo-constraining* mechanism to overcome the confirmation bias, a major weakness of pseudo-labeling, 2) developing new methods for pseudo-labeling towards the selection of *informative* unconstrained samples, 3) showing that this also allows the use of pairwise loss functions for the initial and auxiliary losses which facilitates semi-constrained model training. In extensive experiments, we demonstrate the effectiveness of *ConstraintMatch* over relevant baselines in both the regular clustering and overclustering scenarios on five challenging benchmarks and provide analyses of its several components.

I. INTRODUCTION

Manual annotation of class labels is a tedious and labor-intensive task that can constitute a significant obstacle in applications, particularly in situations where the annotator has to select from a large number of potential class labels or where the annotation is ambiguous due to the task complexity. Additionally, supervised classification models require knowledge of the total number of classes present in the respective application, i.e. the cardinality of the label space. *Constrained clustering* offers a remedy for this as model training in this weakly supervised regime requires only weak, pairwise constraint relations (i.e. similar/dissimilar) which incur less annotation effort compared to instance-specific class labels [47]. These models can also learn meaningful cluster representations even in the overclustering scenario without knowledge of the underlying amount of clusters [14]. The majority of research on constrained clustering focuses on the *constrained* scenario, where each data point is associated with at least one constraint pair. As this setting still requires large sets of given constraints and hence incurs high annotation effort, we focus on the *semi-constrained* setting where a clustering model is trained on both a small dataset of pairwise constraints and a large dataset of unconstrained samples.

While a great deal of progress has been made in semi-supervised learning when class labels are provided, we identify

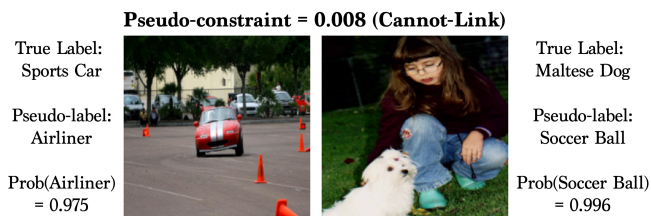


Fig. 1: Illustration of pseudo-constraining. While the model creates overconfident, wrong pseudo-labels for both unlabeled samples, it still yields a semantically correct pseudo-constraint.

through analysis a number of challenges when applying such methods to the constrained clustering setting. One of the most effective methods in semi-supervised learning, pseudo-labeling, utilizes confident predictions on unlabeled data in training and is therefore prone to *confirmation bias*. Specifically, unlabeled samples that were confidently assigned the wrong class label by the model are selected as pseudo-labels, which leads to subsequent model degradation [1]. We analyze this issue in the context of constraint labels and propose a *pseudo-constraining* mechanism that we show can mitigate it, by generating pseudo-constraints from the pseudo-labels (see Fig. 1). Further, we argue that a confidence-based pseudo-label selection criterion is inappropriate in this setting as it leads to the unnecessary de-selection of unconstrained samples that contain valuable information for subsequent pseudo-constraining. We, therefore, propose an entropy-based criterion to select *informative* unconstrained samples and show its superiority. The combination of these two methods, *ConstraintMatch*, facilitates effective pseudo-labeling and unifies the initial and auxiliary learning task. We show that *ConstraintMatch* is able to outperform several state-of-the-art baselines using only a few constraint annotations by substantial margins, even in the more challenging overclustering scenario.

Contributions We 1) propose *ConstraintMatch* as a method for semi-constrained training of clustering models leveraging a large set of unconstrained samples next to a small set of pairwise constraints. Within a series of experiments, we 2) specifically make the case for pseudo-constraints over naive pseudo-labels and provide a detailed analysis of *ConstraintMatch*'s several components. Furthermore, we 3) empirically prove the strong performance of *ConstraintMatch* of up to 16.75% NMI over the constrained baseline on a series of five challenging benchmark datasets in both the regular and the overclustering scenario. Thereby, we evaluate models in

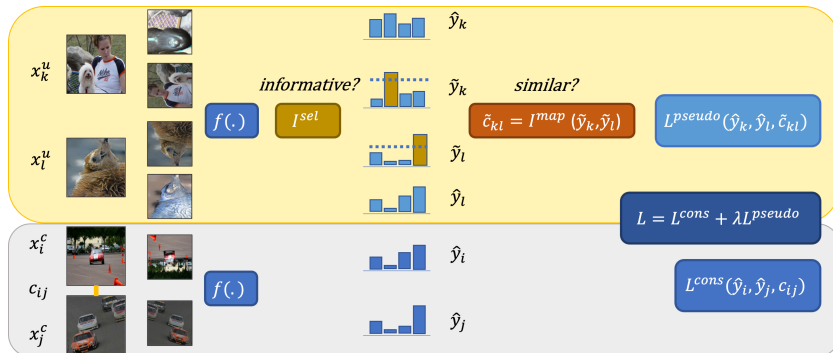


Fig. 2: ConstraintMatch combines pairwise training on constrained (gray) and unconstrained (yellow) samples leveraging weak and strong data augmentations. The criterion \mathcal{L}^{sel} is used to select *informative* pseudo-labels from unconstrained samples which are then mapped to pairwise pseudo-constraints via \mathcal{I}^{map} to overcome the *confirmation bias*. Predictions from model $f(\cdot)$ over strongly augmented versions of these samples serve as inputs to the auxiliary loss \mathcal{L}^{pseudo} to enforce consistency in predicted cluster assignments and the final model is trained on a combination of the pseudo-constrained and the constrained loss \mathcal{L}^{cons} .

different settings to unify the evaluation of modern deep clustering approaches and 4) release our source code¹ for future research on semi-constrained clustering.

II. RELATED WORK

We provide an overview of the context of ConstraintMatch at the intersection of deep clustering, constrained clustering, and semi-supervised learning in the following.

Deep Clustering Early methods for deep clustering combine a reconstruction target with a clustering loss to learn expressive clustering features via reconstruction [10], [23], [37]. Subsequent approaches shift this focus toward low-level features via alternating cluster assignments with those provided by traditional clustering algorithms [3], [45]. More recent research is directed at mapping the data onto low-dimensional representations which serve as a training target for similarity-based losses and as cluster predictions during model inference [4], [18], [24], [29], [44]. Van Gansbeke et al. [40] found these approaches are prone to learning low-level features which lack meaning for semantic clustering next to heavy dependence on network initialization. Therefore, they propose SCAN as a two-step approach where feature representations learned via contrastive pretext tasks [5], [11] are used to mine nearest neighbors of the unlabeled samples. The model is then trained via a clustering loss which maximizes the alignment of their joint feature representations and enables clustering in the absence of the true underlying amount of clusters. SCAN was decidedly evaluated on test datasets only to prove its efficacy on new, unseen data. While this is appealing from a modeling perspective, it prevents direct comparison with prior work where clustering models are evaluated on the union of training and test datasets. We unify this model comparison in our experiments and find SCAN to perform on par with subsequent approaches TCC [30], CC [24], and MICE [38]. Hence, we use SCAN as a starting point for ConstraintMatch.

Constrained Clustering The introduction of binary instance-level constraints for clustering [41] led to the adaptation of existing clustering methods towards the use of constraints [42], see [9] for an overview. With the proposal of the KCL loss, Hsu et al. [14] introduced constrained clustering and overclustering to deep learning. They further showed its applicability to transfer learning [15] and introduced the Meta-classification-likelihood (MCL) for improved model training with pairwise constraints [16] and both loss formulations can not be used with unconstrained data. Zhang et al. [47] provide a framework to work with various types of constraints.

Semi-supervised Learning In semi-supervised classification, the rationale of consistency regularization lead to substantial improvements over supervised baselines [2], [21], [33], [46]. Among these, FixMatch [33] yields state-of-the-art model performance even in settings with very low supervision. It combines confidence-based pseudo-labeling [22] with consistency regularization using the weak-and strong augmentation scheme over unlabeled samples.

Semi-constrained Clustering With S^3C^2 , a two-stage approach was proposed that leverages pseudo-constraints mined from a siamese network trained on few constraints [32]. While this approach was shown to perform well on simple benchmarks, it lacks end-to-end training and requires the true amount of clusters as input. Similarly, the approach by Fogel et al. [8] requires said amount of true clusters as input next to being a transductive method prohibiting inference on unseen data without access to the training data after training. PCOG [28] is another transductive method that also requires the true amount of cluster while its spectral decomposition component hinders it from scaling to large datasets. The approach of Shukla et al. [31] relies on a few class labels, rendering it non-applicable for scenarios where only constraints are present. In contrast to these approaches, we introduce a semi-constrained clustering method that only relies on constraint annotations, leverages unconstrained data, is inductive, and works well without knowledge w.r.t the true amount of clusters.

¹<https://github.com/slds-lmu/constraintmatch>

III. METHOD

A. Notation

We consider a dataset \mathcal{D} which consists of constrained and unconstrained datasets \mathcal{D}^c and \mathcal{D}^u . \mathcal{D}^c contains n_c constrained pairs of the form $x_{ij}^c = (x_i^c, x_j^c, c_{ij}) \in \mathcal{D}^c$ where x_i^c, x_j^c refer to two input samples and $c_{ij} \in \{0, 1\}$ to the associated binary constraint. These constraints describe that both samples either correspond to the same cluster $c_{ij} = 1$, *Must-Link* constraints (ML), or to different clusters $c_{ij} = 0$, *Cannot-Link* constraints (CL). \mathcal{D}^u consists of n_u unconstrained input samples $x_i^u \in \mathcal{D}^u$. We denote $\mathcal{B} \subset \mathcal{D}, \mathcal{B}^c \subset \mathcal{D}^c, \mathcal{B}^u \subset \mathcal{D}^u$ as batches of input samples x_i of the respective datasets. We refer to true class labels as $y_i \in \mathcal{Y}$ where $K = |\mathcal{Y}|$ describes the amount of true classes, i.e. the amount of underlying clusters K , in the dataset. Note that when K is not known, the model may have a different number of outputs n_{out} than the ground truth number of clusters. We aim at training a clustering model f in the form of a neural network with its final head consisting of n_{out} output neurons followed by a softmax layer, i.e. the model predicts a probability distribution over cluster assignments $\hat{y}_i = f(x_i)$ where \hat{y}_{il} denotes the predicted probability of x_i belonging to cluster $l \in 1, \dots, n_{out}$. Similarly, we refer to pseudo-labels as \tilde{y}_i and to pseudo-constraints as $\tilde{c}_{ij} \in [0, 1]$. We introduce the criterion \mathcal{I}^{sel} which selects a subset of informative pseudo-labels \tilde{y}_i based on their predicted cluster assignments \hat{y}_i from \mathcal{B} . From pairs of selected pseudo-labels \tilde{y}_i, \tilde{y}_j we construct pseudo-constraints \tilde{c}_{ij} using a second criterion \mathcal{I}^{map} .

B. Algorithm

ConstraintMatch is an annotation-efficient method that can leverage large unconstrained (i.e. unlabeled) data \mathcal{D}^u next to few constraint pairs \mathcal{D}^c to train a clustering model f . It uses unsupervised clustering in a pretraining step and combines training strategies from constrained clustering [14], [16] with the semi-supervised method Fixmatch [33], refer to Fig. 2 for illustration. We use SCAN [40] for the pretraining step but other pretraining methods would also be applicable. Specifically, pseudo-labeling (i.e. self-training) has proven itself an effective method for leveraging unlabeled data and is a key component of recent semi-supervised classification models [21], [33]. In naive pseudo-labeling, confident model predictions over unlabeled samples are used as pseudo-targets in an auxiliary classification loss to guide model training next to the initial supervised loss, assuming that model confidence is associated with model correctness [22], [39]. Adapting this concept to constrained clustering, we identified three main weaknesses which we overcome: 1) **Pseudo-constraining**: Prediction errors in the selected pseudo-labels can amplify during training, potentially leading to model degradation, also known as *confirmation bias* [1]. We, therefore, propose the generation of *pseudo-constraints*, as pairwise constraints result in a simpler problem reduction [16]. 2) **Informativeness criterion** to carry information of whether two samples x_i and x_j are predicted to be in the same or a different cluster,

which cannot be done via maximal prediction probability [33] or alternative uncertainty metrics [1], and 3) **Unification of losses** utilizing a constraint-based loss for the unlabeled set.

Our overall algorithm processes unconstrained batches \mathcal{B}^u via an unconstrained branch and constrained batches \mathcal{B}^c within a constrained branch to enable training of clustering model f in this semi-constrained data scenario. The constrained branch is trained via a pairwise objective \mathcal{L}^{cons} which allows the training of the model f on binary pairwise constraints $x_{ij}^c = (x_i^c, x_j^c, c_{ij}) \in \mathcal{D}^c$. Therefore, we combine the predictions from model f over weakly augmented constrained samples $a(x_i^c), a(x_j^c)$ along the associated constraint c_{ij} within the pairwise loss function \mathcal{L}^{cons} . For the unconstrained branch, we build upon the intuition of consistency regularization via weak and strong data augmentation strategies $a()$ and $A()$ [33] as follows. Given a pair of unconstrained samples x_i^u, x_j^u , we use the selection criterion \mathcal{I}^{sel} to select *informative* model predictions over weakly augmented versions of those samples $(f(a(x_i^u)), f(a(x_j^u)))$ as pseudo-labels $(\tilde{y}_i, \tilde{y}_j)$. These are then combined into pseudo-constraints \tilde{c}_{ij} via \mathcal{I}^{map} and used as targets within the auxiliary loss function \mathcal{L}^{pseudo} . Model predictions over strongly augmented versions of the unconstrained pair $(\hat{y}_i, \hat{y}_j) = (f(A(x_i^u)), f(A(x_j^u)))$ serve as inputs for \mathcal{L}^{pseudo} . *ConstraintMatch* is trained using the combined loss function $\mathcal{L} = \mathcal{L}^{cons} + \lambda \mathcal{L}^{pseudo}$. The components of *ConstraintMatch* are explained in the following.

1) **Pseudo-Label Selection** Semi-supervised approaches use model confidence as measured via the maximal prediction probability [33] or alternative uncertainty metrics [1] as selection criteria. Model confidence assumes uni-modal model predictions, i.e. the model is confident that sample x_i belongs to class $\hat{y}_i = l$. In contrast to pseudo-labeling, we do not need the information of whether one sample x_i is confidently predicted to be in class $\hat{y}_i = l$ but the information of whether samples x_i and x_j are predicted to be in the same or a different cluster. Filtering for model confidence de-selects multi-modal model predictions that would, for instance, be assigned to two clusters with high probability each - pseudo-constraining allows us to use such multi-modal predictions.

Given a batch of model predictions over weakly-augmented, unconstrained samples, we aim to select those that are important for the subsequent pseudo-constraint generation. We propose measuring such *informativeness* of a probability vector \hat{y}_i using the *normalized entropy*:

$$\mathcal{H}_n(\hat{y}_i) = -\frac{1}{\log(n_{out})} \sum_{l=1}^{n_{out}} p(\hat{y}_{il}) \log(p(\hat{y}_{il})) \quad (1)$$

with $\mathcal{H}_n(\hat{y}_i) \in [0; 1]$ where $\mathcal{H}_n(\hat{y}_i) = 1$ describes the minimum level of information and maximal entropy and $\mathcal{H}_n(\hat{y}_i) = 0$ the maximum level of potential informativeness and minimal entropy where the model places the entire probability mass in one cluster. Hence, we use the normalized entropy in combination with a threshold hyperparameter $\tau \in [0; 1]$ as criterion to select pseudo-labels:

$$\mathcal{I}_\tau^{sel}(\hat{y}_i) = \mathbb{1}(\mathcal{H}_n(\hat{y}_i) < \tau) \quad (2)$$

where $\mathbb{1}$ is an indicator function. In the experiment section, we provide an empirical analysis of the suitability of this criterion next to a sensitivity analysis of τ .

2) Pseudo-Constraining Confirmation bias is a critical problem in pseudo-labeling methods [1], and in constraint-based clustering, there is an opportunity to alleviate this. As an illustration, refer to the two unconstrained samples x_i^u, x_j^u from Fig. 1: the true label y_i of x_i^u would be "sports car" while the model wrongly but confidently assigns it to the "airliner" cluster and similarly x_j^u is assigned the wrong cluster "soccer ball" instead of the true y_j "maltese dog" (see Fig. 6 in the Appendix for more examples). While this prediction error would lead to a wrong prediction target in naive pseudo-labeling and hence confuse model training, the resulting pseudo-constraint $\tilde{c}_{ij} = 0.008$ would still be correctly assigned as Cannot-Link, as $\tilde{y}_j \neq \tilde{y}_i$ in both situations. Therefore, we create pseudo-constraints from the pseudo-labels to drive the loss function \mathcal{L}^{pseudo} . Given a batch of informative pseudo-labels, we next combine pseudo-label pairs \tilde{y}_i, \tilde{y}_j into pseudo-constraints \tilde{c}_{ij} , expressing the (dis-)similarity of those samples. As \tilde{y}_i, \tilde{y}_j are probability vectors, we propose to use a divergence measure to quantify this distance and derive a meaningful pseudo-constraint. The Jensen-Shannon-Distance [25] allows the symmetric mapping of two probability vectors onto a similarity score:

$$JSD(\tilde{y}_i, \tilde{y}_j) = \sqrt{((KL(\tilde{y}_i|m) + KL(\tilde{y}_j|m))/2)} \quad (3)$$

where $m = (\tilde{y}_i + \tilde{y}_j)/2$ and $KL(y_i|m)$ refers to the Kullback-Leibler Distance between \tilde{y}_i and m and $JSD(\tilde{y}_i, \tilde{y}_j) \in [0, 1]$. We exploit this property and use the inverse Jensen-Shannon-Distance to calculate soft pseudo-constraints $\tilde{c}_{ij} = 1 - JSD(\tilde{y}_i, \tilde{y}_j) \in [0; 1]$ where $\tilde{c}_{ij} = 0.0$ resembles a Cannot-Link and $\tilde{c}_{ij} = 1.0$ a Must-Link pseudo-constraint over all pairwise combinations of the informative pseudo-labels. We refer to this inverse Jensen-Shannon-Distance as $\mathcal{I}^{map}(\tilde{y}_i, \tilde{y}_j)$ in Fig. 2. Pseudo-constraints are generated over the combined batch $\mathcal{B} = \mathcal{B}^c \cup \mathcal{B}^u$, treating the samples in \mathcal{B}^c as unconstrained.

3) Pairwise Loss Function There exists a variety of loss functions that can deal with pairwise constraints [47] with the KCL [14] and the MCL [16] being the most prominent ones. Following the findings of Hsu et al. [16] and guided by preliminary experimental results, we propose the use of the MCL as a pairwise loss function, as it was shown to result in higher model performance and smoother model training, and is hyperparameter-free. The MCL is aligned on the binary cross-entropy loss and follows the definition:

$$\mathcal{L}(c_{ij}, \hat{c}_{ij}) = - \sum_{ij} c_{ij} \log(\hat{c}_{ij}) + (1 - c_{ij}) \log(1 - \hat{c}_{ij}) \quad (4)$$

where $\hat{c}_{ij} = \langle \hat{y}_i, \hat{y}_j \rangle$ combines the individual predicted cluster assignment vectors into an alignment score and $c_{ij} \in \{0, 1\}$

refers to the pairwise constraint $c_{ij} = 0$ for Cannot-Link and $c_{ij} = 1$ for Must-Link constraints. The MCL allows training with soft constraints $c_{ij} \in [0, 1]$, similar to the use of soft labels in the cross-entropy loss [27]. We use this property for the processing of soft pseudo-constraints \tilde{c}_{ij} as explained above and analyzed further in the experiments section. ConstraintMatch is trained on the combined loss:

$$\mathcal{L} = \sum_{x_i, x_j \in \mathcal{B}^c} \mathcal{L}^{cons}(c_{ij}, \langle f(a(x_i)), f(a(x_j)) \rangle) + \lambda \sum_{x_k, x_l \in \mathcal{B}} \mathcal{L}^{pseudo}(\tilde{c}_{kl}, \langle f(A(x_k)), f(A(x_l)) \rangle) \quad (5)$$

where hyperparameter $\lambda \geq 0$ controls the impact of the pseudo-constraint loss and is tuned on the validation set.

IV. EXPERIMENTS

In this section, we compare the performance of ConstraintMatch with prior work on five challenging benchmark datasets and provide empirical evidence for the effectiveness of pseudo-constraining. This includes i) the relative improvement of ConstraintMatch across various baselines, ii) an empirical analysis of the benefit of pseudo-constraining and iii) its robustness w.r.t annotation noise, iv) analyses of the algorithmic choices made for its several components, and v) an evaluation of ConstraintMatch in the overclustering scenario.

A. Experimental Setup

Datasets and Constraint Mining We use the Cifar10 [19], Cifar100 [19], STL10 [6], ImageNet-10 [4] and ImageNet-Dogs [4] datasets to demonstrate the effectiveness of the proposed method. We use the 20 superclasses in Cifar100 as ground truth labels for constraint mining following prior work [24], [40], declaring it as Cifar100-20 in the following. We adhere to the provided train/test splits to enable comparison with prior work evaluated on separate test datasets [16], [40]. Recent deep clustering approaches instead are evaluated on the training set or the union of training and test datasets [24], [30], [38]. To be comparable with both bodies of literature, we provide benchmark results in both settings marked as "Test" and "Train(+Test)" in Table I following [24]. An overview of the used datasets, the amount of sampled constraints as well as the training and validation splits for hyperparameter tuning is provided in Table V in the Appendix. For constraint-sampling, we randomly sample n_c constraints from each dataset via the following procedure: n_c samples are randomly sampled without replacement as constraint members (x_i^c, y_i) and for each of those samples, a second pair member (x_j^c, y_j) is randomly chosen with replacement from the remaining training samples to create the constraint pair (x_i^c, x_j^c, c_{ij}) , where $c_{ij} = 1$, if $y_i = y_j$ and $c_{ij} = 0$, if $y_i \neq y_j$. This results in a dataset \mathcal{D}^c of n_c constrained samples (x_i^c, x_j^c, c_{ij}) . To account for randomness in the constraint sampling process, we report performance averaged over five random sampling repetitions.

TABLE I: Comparison of ConstraintMatch with relevant baselines (B), competitors (C), and upper bound models (U) across datasets and varying amounts of constraints n_c . Performance metrics were averaged over five folds and calculated on separate test splits in the upper part and in the Train(+Test) setting in the lower part. Best results comparing ConstraintMatch with the baseline and competitor models are shown in bold and \dagger denotes values reported in the literature. Statistical significance for differences in model performance between ConstraintMatch and the constrained competitor for $n_c \in \{5k, 10k\}$ respectively established using the Wilcoxon signed-rank test [7], [43] (significance code $*$: $p < 0.05$).

Split	Model	n_c	Cifar 10			Cifar 100-20			STL 10			ImageNet 10			ImageNet Dogs		
			ACC	NMI	ARI	ACC	NMI	ARI	ACC	NMI	ARI	ACC	NMI	ARI	ACC	NMI	ARI
Test	Supervised \dagger	U	93.80	86.20	87.00	80.00	68.00	63.20	80.60	65.90	63.10	-	-	-	-	-	-
	Fully Constrained	U	94.86	88.39	89.11	77.99	68.37	61.94	90.49	80.94	80.47	96.64	93.45	92.60	67.10	73.52	58.13
	SCAN \dagger [40]	B 0	87.60	78.70	75.80	45.90	46.80	30.10	76.70	68.00	61.60	86.20	81.57	75.71	47.20	55.42	35.87
	Constrained	C 5k	90.12	80.52	80.02	50.99	46.23	32.63	85.90	74.62	72.51	93.12	87.43	85.49	44.08	43.27	28.92
	ConstraintMatch	5k	92.23*	84.64*	84.27*	54.19*	52.74*	37.84*	88.20*	78.20*	76.68*	94.68*	90.43*	88.61*	49.43*	55.23*	38.12*
	Constrained	C 10k	90.89	81.73	81.47	52.45	46.57	33.79	88.21	77.63	76.38	94.90	90.42	89.04	45.52	44.44	30.10
	ConstraintMatch	10k	93.17*	85.88*	85.92*	57.15*	54.37*	40.59*	90.08*	80.57*	79.81*	95.68	92.09	90.70	50.73*	54.92*	38.34*
Train (+Test)	PICA \dagger [17]	B 0	69.60	59.10	51.20	33.70	31.00	17.10	71.30	61.10	53.10	87.00	80.20	76.10	35.20	35.20	20.10
	MICE \dagger [38]	B 0	83.50	73.70	69.80	44.00	43.60	28.00	75.20	63.50	57.50	-	-	-	43.90	42.30	28.60
	CC \dagger [24]	B 0	79.00	70.50	63.70	42.90	43.10	26.60	85.00	76.40	72.60	89.30	85.90	82.20	42.90	44.50	27.40
	TCC \dagger [30]	B 0	90.60	79.00	73.30	49.10	47.90	31.20	81.40	73.20	68.90	89.70	84.80	82.50	59.50	55.40	41.70
	SCAN [40]	B 0	88.53	80.09	77.72	50.67	47.72	33.07	81.28	70.15	65.22	91.63	84.00	82.93	44.06	45.09	30.75
	Constrained	C 5k	91.14	82.30	82.05	51.63	46.58	33.37	80.51	68.38	63.68	95.09	88.42	89.49	43.25	38.82	28.93
	ConstraintMatch	5k	92.67*	85.12*	84.98*	54.16*	52.68*	37.79*	82.97*	71.13*	67.80*	95.61*	89.64*	90.59*	47.63*	47.82*	35.95*
	Constrained	C 10k	92.21	83.07	84.08	53.13	47.20	34.86	89.90	80.12	79.49	96.47	91.18	92.36	44.17	40.07	30.13
	ConstraintMatch	10k	93.61*	86.55*	86.80*	57.18*	53.37*	40.30*	91.30*	82.42*	82.13*	96.68*	91.59*	92.80*	49.34*	49.16*	37.37*

Implementation Details In accordance with prior work [40], we used a ResNet-18 backbone architecture [12] for the experiments with the Cifar10, Cifar100-20, and STL10 datasets and a ResNet-34 backbone [12] following [24] for the ImageNet datasets. We used model weights that were pre-trained via SCAN [40] for the initialization of the model backbone as ConstraintMatch benefits from expressive feature representations as a warm start. Next to the model weights released by [40] for Cifar10, Cifar100-20, and STL10 we used the authors’ codebase² to pretrain the ResNet-34 backbone via SCAN and then used these resulting model weights for model initialization. For model training, we used a standard SGD optimizer with momentum set to 0.9 and weight decay regularization [36] and all models were trained for a total of 20000 optimization steps unless noted otherwise. We used a cosine learning rate scheduler [26] which updates the learning rate at each update step to $\eta \cos\left(\frac{7\pi t}{16T}\right)$ with η being the initial learning rate, t the current training step and $T = 20000$ the total amount of training steps following [33]. Hyperparameters were tuned via a grid search on constraints mined from the validation datasets with hyperparameter ranges shown in Table IV and more details on the validation splits are given in Table V in the Appendix. The size of constrained/unconstrained batches was set to 200/600 respectively for Cifar10 and Cifar100-20 and to 100/300 for the other datasets.

Model Comparison We compare ConstraintMatch with different baselines (B), competitors (C), and upper bound models (U). This includes deep clustering models SCAN [40], TCC [30], CC [24], MICE [38] and PICA [17] as baselines

and a constrained clustering model that was trained on D^c using the MCL [16] as competitor. As upper bounds, we compare with a fully constrained clustering model trained on a fully constraint version of training dataset D and a supervised baseline where the backbone was trained on the fully labeled training set D as reported by [40]. The authors of MICE [17], PICA [38], CC [24] and TCC [30] used a ResNet-34 backbone for the Cifar10, Cifar100-20 and STL10 datasets. We used a ResNet-18 backbone for these three datasets in adherence with SCAN [40]. A comparison with the semi-constrained approaches was not possible due to a lack of open-source code and performance metrics on established benchmarks.

B. Results

We summarize our main results in Table I measuring model performance in Accuracy (ACC), Normalized Mutual Information (NMI) [35] and the Adjusted Rand Index (ARI) [34], as standard in (constrained) clustering [47]. As established in the literature [14], [30], [40], we use the Hungarian Assignment method to optimally map the resulting cluster predictions to the true cluster labels [20]. As expected and shown in previous work [14], [16], we find training with pairwise constraints to be a valid option to train strong-performing clustering models. This benchmark is the first attempt to compare SCAN with subsequent deep clustering methods on the union of train and test datasets showing that SCAN is competitive with those methods (lower part of Table I). Further, fine-tuning SCAN via a subset of constraints improves model performance across all datasets but the fine-grained ImageNet-Dogs dataset. Overall, ConstraintMatch outperforms the (un-)constrained baselines and competitors across all datasets in both evaluation settings

²<https://github.com/wvangansbeke/Unsupervised-Classification>

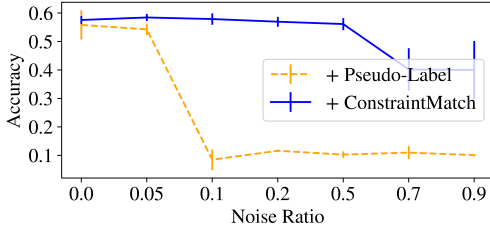


Fig. 3: Robustness of the pseudo-labeling baseline and ConstraintMatch towards pseudo-label noise.

except ImageNet-Dogs, a task with semantically very similar classes, where it falls behind TCC in the Train(+Test) evaluation. This fine-grainedness makes training of constrained clustering models challenging, as the distinction between Must- and Cannot-link losses expressiveness which also explains the worse performance of the constrained competitor compared to SCAN. We interpret the finding that ConstraintMatch, in turn, outperforms both models by substantial margins in ACC and ARI as further proof of the effectiveness of pseudo-constraining. Relative (absolute) performance gains are the largest for Cifar100-20 with ConstraintMatch increasing model performance over the constrained baseline with 10k constraints by 8.96% (4.70 percentage points) Accuracy, 16.75% (7.80pp) NMI and 20.12% (6.80pp) ARI on the test dataset. We attribute these large performance gains to the complexity of the task and the efficient use of pseudo-constraints in this complex 20-cluster setting. Further, both the constrained competitor and ConstraintMatch benefit from more constraints n_c , with a larger relative performance increase for ConstraintMatch. We yield a statistically significant difference in model performance for ConstraintMatch and the constrained competitor across all datasets for all performance metrics using a Wilcoxon signed-rank test [7], [43] ($p < 0.05$) for $n_c \in \{5k, 10k\}$. Those empirical results confirm ConstraintMatch as a suitable method for semi-constrained clustering.

C. The Empirical Case for Pseudo-Constraints

We propose pseudo-constraining to overcome *confirmation bias*. To support this claim, we conducted a simulation experiment to evaluate the robustness of naive pseudo-labeling against confirmation bias in comparison to subsequent pseudo-constraining within ConstraintMatch. This naive pseudo-labeling baseline differs from ConstraintMatch in the handling of unconstrained samples, similar to the processing of unlabeled data in FixMatch [33]: weakly augmented, unconstrained samples are selected via a confidence threshold over their predicted cluster assignment and the major predicted cluster is chosen as pseudo-label. Predictions over strong augmented versions of these samples then serve as input for an auxiliary cross-entropy loss function (see Fig. 7 in the Appendix). We introduce a mode-flip function $m(\hat{y}_i)$ that swaps the position of the two largest predicted probabilities within the model prediction \hat{y}_i . This simulates a prediction

TABLE II: Ablation study on ConstraintMatch, results averaged over 5 folds with $n_c = 10000$.

Model	Test Performance		
	ACC	NMI	ARI
SCAN	45.90	46.80	30.10
+ Constrained	52.45	46.57	33.79
+ Pseudo-Labeling	55.38	53.49	39.98
+ Pseudo-Constraining	57.15	54.37	40.59

error where the model "confuses" two cluster assignments within the pseudo-labeling of x_i . We randomly apply $m()$ to a noise fraction ρ of the unconstrained samples $x_i \in \mathcal{B}^u$ and train both models in this setting. The results in Fig. 3 confirm our intuition as naive pseudo-labeling already degrades at $\rho \geq 0.1$ while ConstraintMatch can cope with $\rho \leq 0.5$.

Pseudo-constraining further allows to use the same pairwise loss function as both the auxiliary and the initial objective for model training. We provide an ablation study on Cifar100-20 to quantify this benefit where we subsequently add constrained training, naive pseudo-labeling, and finally pseudo-constraining to the SCAN model, each with tuned hyperparameters. The results in Table II show that while the use of naive pseudo-labeling leads to a substantial performance gain over the constrained baseline, the subsequent application of pseudo-constraining within ConstraintMatch enables further model improvements. We conclude its effectiveness is grounded in both the robustness w.r.t. the confirmation bias and the similarity in training objectives.

D. Additional Analyses

In this section, we analyze the several components of ConstraintMatch. Unless noted otherwise, these analyses and experiments were run on the Cifar100-20 dataset with $n_c = 10000$ using the optimal hyperparameters obtained for the main experiments and we report results on the test splits.

Pseudo-Constraining We use soft pseudo-constraints $\tilde{c}_{ij} \in [0, 1]$ in ConstraintMatch. One alternative would be the separation into hard pseudo-constraints $\hat{c}_{ij}^h \in \{0, 1\}$ using a threshold μ such that $\hat{c}_{ij}^h = 1$ if $c_{ij} \geq \mu$ and $\hat{c}_{ij}^h = 0$ if $c_{ij} < \mu$. We find that while ConstraintMatch is still outperforming the constrained baseline using hard pseudo-constraints, it benefits further from the use of soft pseudo constraints as shown in Fig. 4a over different values of μ . Using soft constraints also eliminates the need to tune μ . We hypothesize that the model can effectively use the continuous information provided via the soft pseudo-constraints within the MCL loss, similar to the use of soft labels in supervised classification [27].

Pseudo-Label Selection We argue for the selection of *informative* samples as pseudo-labels over that of *confident* samples. Fig. 4b contrasts the use of both with a fixed value $\tau = 0.2$ for the informativeness criterion and with varying thresholds for confidence-based selection showing that informative samples enable ConstraintMatch to leverage unconstrained samples more effectively. Further, we provide a sensitivity analysis of the threshold τ within \mathcal{I}_τ^{sel} in

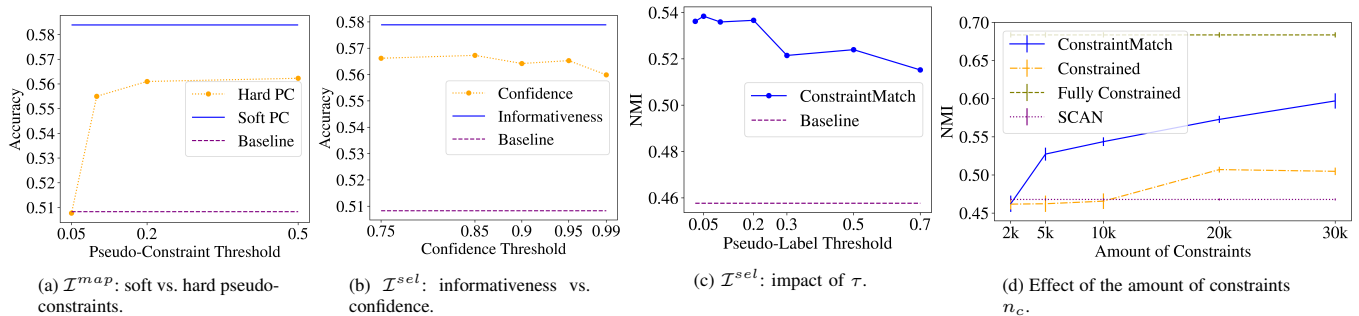


Fig. 4: Further analyses of ConstraintMatch.

Fig. 4c. This reveals that the sensitivity of ConstraintMatch towards τ lies within a reasonable margin and we recommend $\tau \in [0.025, 0.2]$ as a range for tuning.

Amount of Constraints Fig. 4d shows the effect of the amount of constraints n_c on the model performance as measured in NMI over five folds. The constrained clustering competitor model performs worse than the SCAN baseline for $n_c \leq 10000$ which might be due to the fact that $n_c = 10000$ results in 500 Must-Link constraints only in the Cifar100-20 scenario, allowing the MCL loss to overfit those few pairs quickly. In contrast to that, ConstraintMatch successfully overcomes this issue via its pseudo-constraining mechanism for $n_c \geq 5000$ with relative gains increasing for increasing n_c . The comparably low performance of ConstraintMatch for $n_c = 2000$ indicates that it still requires a certain degree of supervision to produce reliable pseudo-constraints.

Robustness w.r.t Noisy Constraints Next to the robustness of ConstraintMatch over noisy pseudo-labels, we further investigate its robustness towards noise in the annotation of the known ground truth constraints. This simulates the situation where the annotators might erroneously flip the constraint annotation, similar to the concept of label noise in supervised classification [13]. Therefore, we randomly flipped a varying percentage of the known constraint annotations and compared the effect of this annotation noise on model training. The results in Fig. 5 show that ConstraintMatch is more robust towards higher levels of annotation noise than the constrained competitor. We attribute this increased robustness to the stabilizing effect of the pseudo-constraining mechanism.

E. Overclustering

We further evaluate ConstraintMatch for overclustering, where the true amount of clusters K is unknown and the model can assign more clusters than inherently present in the data, $n_{out} \gg K$ [14]. Therefore, we compare ConstraintMatch with the constrained competitor and the SCAN baseline with $n_{out} = 5K$ resulting in 100 potential clusters for Cifar100-20 and 50 for Cifar10. Models were evaluated using the Hungarian Assignment [20] with cluster predictions that do not match a corresponding ground truth cluster counting as an error. As shown in Table III, we find that the constrained competitor achieves strong performance gains over the unsupervised baseline despite the challenging task. ConstraintMatch’s performance gains translate well to this overclustering

TABLE III: Overclustering results averaged over five folds.

Dataset n_{out}	Cifar-10 50			Cifar 100-20 100		
	ACC	NMI	ARI	ACC	NMI	ARI
SCAN	34.68	61.56	34.52	29.88	47.35	23.23
Constrained	82.24	75.40	73.80	39.41	44.34	27.86
ConstraintMatch	88.89	83.04	82.06	43.65	52.37	34.88

scenario yielding a relative (absolute) performance gain over the constrained competitor of 18.11% (8.03pp) NMI and 10.75% (4.24pp) Accuracy on Cifar100-20.

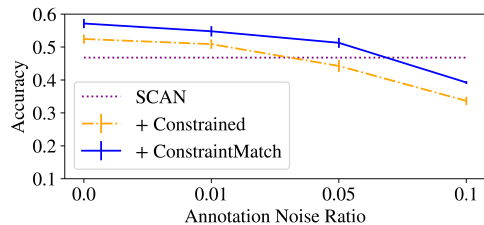


Fig. 5: Impact of ground truth constraint annotation noise.

V. CONCLUSION

ConstraintMatch is a novel method for training clustering models in a semi-constrained setting, using a combination of large amounts of unconstrained data and a limited number of constraint pairs. Therefore, it selects informative pseudo-labels processed within a pseudo-constraining mechanism that allows training the model on a unified loss function to overcome the limitations of naive pseudo-labeling in this setting. With empirical results across five benchmarks, we demonstrate ConstraintMatch’s strong performance, outperforming baselines and competitors by substantial margins, even in challenging overclustering scenarios. We furthermore support our algorithmic choices with empirical evidence and empirically showed that pseudo-constraining leads to increased model robustness towards different sources of annotation noise.

While we have shown the application in scenarios with a cluster cardinality of up to $K \leq 20$, we did not yet investigate the applicability to high-cardinality settings. Also, the use of alternative constraints such as triplets or continuous constraints as well as experiments with non-uniform annotation noise would be interesting for future research.

ACKNOWLEDGMENTS

JG and BB were partially supported by the Bavarian Ministry of Economic Affairs, Regional Development, and Energy through the Center for Analytics – Data – Applications (ADA-Center) within the framework of BAYERN DIGITAL II (20-3410-2-9-8) and the German Federal Ministry of Education and Research (BMBF) under Grant No. 01IS18036A, Munich Center for Machine Learning (MCML). ZK was partly supported by DARPA’s Learning with Less Labels (LwLL) program under agreement HR0011-18-S-0044

REFERENCES

- [1] E. Arazo, D. Ortego, P. Albert, N. O’Connor, and K. McGuinness. Pseudo-labeling and confirmation bias in deep semi-supervised learning. *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2020.
- [2] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. Raffel. Mixmatch: A holistic approach to semi-supervised learning. *NeurIPS*, 32, 2019.
- [3] M. Caron, P. Bojanowski, A. Joulin, and Ma. Douze. Deep clustering for unsupervised learning of visual features. *Proceedings of the European conference on computer vision (ECCV)*, pages 132–149, 2018.
- [4] J. Chang, L. Wang, G. Meng, S. Xiang, and C. Pan. Deep adaptive image clustering. *IEEE/CVF CVPR Proceedings*, pages 5879–5887, 2017.
- [5] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. *International Conference on Machine Learning*, pages 1597–1607. PMLR, 2020.
- [6] A. Coates, A. Ng, and H. Lee. An analysis of single-layer networks in unsupervised feature learning. *Proceedings of the fourteenth International Conference on Artificial Intelligence and Statistics*, pages 215–223. JMLR, 2011.
- [7] D. Janeczko. Statistical comparisons of classifiers over multiple data sets *Journal of Machine Learning Research*, pages 1–30. JMLR, 2006.
- [8] S. Fogel, H. Averbuch-Elor, D. Cohen-Or, and J. Goldberger. Clustering-driven deep embedding with pairwise constraints. *IEEE Computer Graphics and Applications*, 39(4):16–27, 2019.
- [9] P. Gançarski, T. Dao, B. Crémilleux, G. Forestier, and T. Lampert. Constrained clustering: Current and new trends. *Guided Tour of Artificial Intelligence Research*, pages 447–484. Springer, 2020.
- [10] X. Guo, L. Gao, X. Liu, and J. Yin. Improved deep embedded clustering with local structure preservation. *IJCAI*, pages 1753–1759, 2017.
- [11] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual representation learning. *IEEE/CVF CVPR Proceedings*, pages 9729–9738, 2020.
- [12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *IEEE/CVF CVPR Proceedings*, pages 770–778, 2016.
- [13] M. Hedderich, D. Zhu, and D. Klakow. Analysing the noise model error for realistic noisy label data. *AAAI Proceedings*, pages 7675–7684, 2021.
- [14] Y. Hsu and Z. Kira. Neural network-based clustering using pairwise constraints. *International Conference on Learning Representations Workshop*, 2016.
- [15] Y. Hsu, Z. Lv, and Z. Kira. Learning to cluster in order to transfer across domains and tasks. *International Conference on Learning Representations*, 2018.
- [16] Y. Hsu, Z. Lv, J. Schlosser, P. Odom, and Z. Kira. Multi-class classification without multi-class labels. *International Conference on Learning Representations*, 2019.
- [17] J. Huang, S. Gong, and X. Zhu. Deep semantic clustering by partition confidence maximisation. *IEEE/CVF CVPR Proceedings*, pages 8849–8858, 2020.
- [18] X. Ji, J. Henriques, and A. Vedaldi. Invariant information clustering for unsupervised image classification and segmentation. *IEEE/CVF CVPR Proceedings*, pages 9865–9874, 2019.
- [19] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. *Advances in Neural Information Processing Systems*, 2009.
- [20] H. Kuhn. The Hungarian Method for the Assignment Problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.
- [21] C. Kuo, C. Ma, J. Huang, and Z. Kira. Featmatch: Feature-based Augmentation for Semi-supervised Learning. *Proceedings of the European conference on computer vision (ECCV)*. Springer, 2020.
- [22] D. Lee. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. *Workshop on challenges in representation learning, International Conference on Machine Learning*, page 896, 2013.
- [23] F. Li, H. Qiao, and B. Zhang. Discriminatively boosted image clustering with fully convolutional auto-encoders. *Pattern Recognition*, 83:161–173, 2018.
- [24] Y. Li, P. Hu, Z. Liu, D. Peng, J. Zhou, and X. Peng. Contrastive clustering. *AAAI Proceedings*, 2021.
- [25] J. Lin. Divergence measures based on the shannon entropy. *IEEE Transactions on Information theory*, 37(1):145–151, 1991.
- [26] I. Loshchilov and F. Hutter. SGDR: Stochastic gradient descent with warm restarts. *International Conference on Learning Representations*, 2017.
- [27] C. Meister, E. Salesky, and R. Cotterell. Generalized entropy regularization or: There’s nothing special about label smoothing. *ACL*, 2020.
- [28] F. Nie, H. Zhang, R. Wang, and X. Li. Semi-supervised clustering via pairwise constrained optimal graph. *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, pages 3160–3166, 2021.
- [29] C. Niu and G. Wang. Spice: Semantic pseudo-labeling for image clustering. *arXiv preprint arXiv:2103.09382*, 2021.
- [30] Y. Shen, Z. Shen, M. Wang, J. Qin, P. Torr, and L. Shao. You never cluster alone. *Advances in Neural Information Processing Systems*, 34, 2021.
- [31] A. Shukla, G. Cheema, and S. Anand. Semi-supervised clustering with neural networks. *2020 IEEE Sixth International Conference on Multimedia Big Data (BigMM)*, pages 152–161. IEEE, 2020.
- [32] M. Śmieja, L. Struski, and M. Figueiredo. A classification-based approach to semi-supervised clustering with pairwise constraints. *Neural Networks*, 127:193–203, 2020.
- [33] K. Sohn, D. Berthelot, C. Li, Z. Zhang, N. Carlini, E. Cubuk, A. Kurakin, H. Zhang, and C. Raffel. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *CoRR*, abs/2001.07685, 2020.
- [34] D. Steinley. Properties of the hubert-arable adjusted rand index. *Psychological methods*, 9(3):386, 2004.
- [35] A. Strehl and J. Ghosh. Cluster ensembles - a knowledge reuse framework for combining multiple partitions. *Journal of machine learning research*, 3(Dec):583–617, 2002.
- [36] I. Sutskever, J. Martens, G. Dahl, and G. Hinton. On the importance of initialization and momentum in deep learning. *International Conference on Machine Learning*, pages 1139–1147. PMLR, 2013.
- [37] K. Tian, S. Zhou, and J. Guan. Deepcluster: A general clustering framework based on deep learning. *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 809–825. Springer, 2017.
- [38] T. Tsai, C. Li, and J. Zhu. Mice: Mixture of contrastive experts for unsupervised image clustering. *International Conference on Learning Representations*, 2021.
- [39] J. Van Engelen and H. Hoos. A survey on semi-supervised learning. *Machine Learning*, 109(2):373–440, 2020.
- [40] W. Van Gansbeke, S. Vandenhende, S. Georgoulis, M. Proesmans, and L. Van Gool. Scan: Learning to classify images without labels. *European Conference on Computer Vision*, pages 268–285. Springer, 2020.
- [41] K. Wagstaff and C. Cardie. Clustering with instance-level constraints. *AAAI Proceedings*, 1097:577–584, 2000.
- [42] K. Wagstaff, C. Cardie, S. Rogers and S. Schrödl. Constrained k-means clustering with background knowledge. *International Conference on Machine Learning*, volume 1, pages 577–584, 2001.
- [43] F. Wilcoxon. Individual comparisons by ranking methods. *Biometrics*, volume 1, pages 80–83, 1945.
- [44] J. Wu, K. Long, F. Wang, C. Qian, C. Li, Z. Lin, and H. Zha. Deep comprehensive correlation mining for image clustering. *IEEE/CVF CVPR Proceedings*, pages 8150–8159, 2019.
- [45] J. Yang, D. Parikh, and D. Batra. Joint unsupervised learning of deep representations and image clusters. *IEEE/CVF CVPR Proceedings*, pages 5147–5156, 2016.
- [46] B. Zhang, Y. Wang, W. Hou, H. Wu, J. Wang, M. Okumura, and T. Shinozaki. Flexmatch: Boosting semi-supervised learning with curriculum pseudo labeling. *Advances in Neural Information Processing Systems*, volume 34, pages 18408–18419.
- [47] H. Zhang, T. Zhan, S. Basu, and I. Davidson. A framework for deep constrained clustering. *Data Mining and Knowledge Discovery*, 35(2):593–620, 2021.

APPENDIX

A. Hyperparameter Tuning

We conducted a grid search over the validation splits detailed in Table V for one fold of sampled training constraints for hyperparameter tuning. We used the validation loss on the constraints from the validation splits to select the optimal hyperparameters for each dataset and model combination with the lowest final validation loss as performance criterion. Final models were then trained on these optimal hyperparameters on five repeated folds of the respective constrained and unconstrained training samples and final performance metrics were reported for both the "Test" and the "Train(+Test)" settings. The *shared* parameters were used in and tuned for all trained models and the specific hyperparameters for ConstraintMatch and the naive pseudo-labeling baseline were tuned over a grid of different values, see Table IV.

TABLE IV: Hyperparameters and their respective values considered in the grid search for the different models.

Parameter	Search Values
Shared	
Weight decay	0.001, 0.0001, 0.00001
Learning rate	0.03, 0.01, 0.003, 0.001, 0.0001
naive Pseudo-labeling	
λ	1.0, 0.5, 0.1, 0.05
τ	0.7, 0.8, 0.9, 0.95, 0.99
ConstraintMatch	
λ	1.0, 0.5, 0.1, 0.05
τ	0.05, 0.1, 0.2, 0.3

B. Data Augmentation

ConstraintMatch follows the rationale of consistency regularization via weak and strong augmentations $a()$ and $A()$. As weak augmentations $a()$, we used random cropping and horizontal flipping. For strong augmentations, $A()$, we used the RandAugment strategy with the data augmentation procedures used in FixMatch and described in Appendix D of [33].

C. Datasets

Table V provides an overview of the datasets used in the experimental section IV-A alongside their splits and sizes. The final column describes the exact dataset splits that were used in the Train(+Test) evaluation setting following [24].

D. Visualization of the Confirmation Bias

In Fig. 6, we visualized four samples from the unconstrained part of the ImageNet-10 dataset which suffer from the confirmation bias similar to Fig. 1, i.e. samples for which the model confidently predicted the wrong cluster assignment. These unconstrained samples were selected as high-confidence (max. predicted probability > 0.98) but wrongly predicted examples. We can observe that pseudo-labeling would lead to wrong prediction targets (e.g. cluster "Airship" instead of the true cluster "Soccer Ball" in the bottom left example) and

TABLE V: Datasets used in the experiments including the respective training, validation, and test splits. We also mention the evaluation dataset for the Train(+Test) setting in the last column following [24].

Dataset	K	Samples			Constraints		Train(+Test)
		Train	Val	Test	Train	Val	
Cifar10	10	45k	5k	10k	5/10k	10k	Train + Test
Cifar100-20	20	45k	5k	10k	5/10k	10k	Train + Test
STL10	10	4k	1k	8k	5/10k	5k	Train + Test
ImageNet-10	10	12k	1k	500	5/10k	1k	Train
ImageNet-Dogs	15	18.5k	1k	750	5/10k	1k	Train

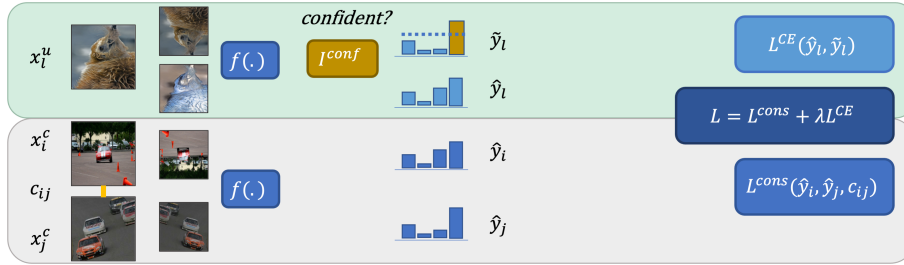


Fig. 6: Illustration of pseudo-labeling failure cases due to confirmation bias. Pseudo-constraints generated on top of these wrong pseudo-labels are still semantically correct.

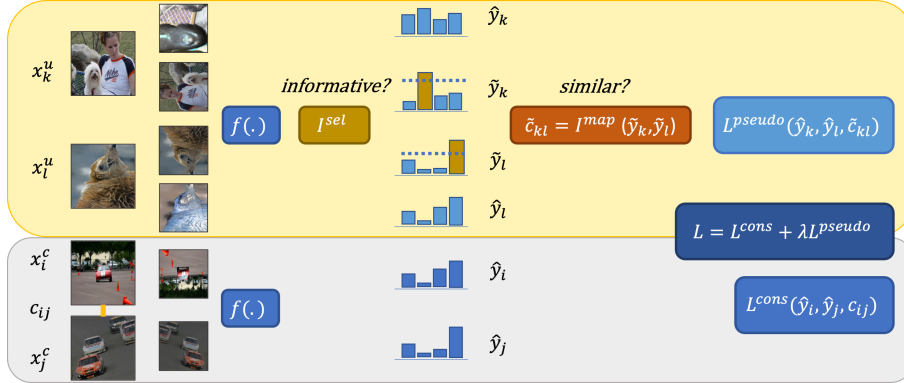
hence confuse model training. On the other hand, pseudo-constraints generated on top of pairs of these wrongly assigned pseudo-labels still are semantically correct and can support model training on these unconstrained samples. This does not only hold for Cannot-Link (bottom) but also for Must-Link (top) pseudo-constraints where both samples with the same true cluster affiliation are assigned the same wrong cluster by the model. These samples were selected from a random batch of unconstrained samples B^u from the ImageNet-10 dataset from ConstraintMatch trained for 500 training steps with a ResNet-34 backbone.

E. Naive Pseudo-Labeling Baseline

In Fig. 7, we visualize the naive pseudo-labeling baseline, a simplified version of ConstraintMatch, with which we compared ConstraintMatch in the results Section IV-B. This baseline follows the use of unlabeled samples in FixMatch [33] and similarly leverages the weak-strong augmentation scheme for consistency regularization. Concretely, weakly augmented, unconstrained samples are selected via a confidence threshold over their predicted cluster assignments, and



(a) The naive pseudo-labeling baseline combines training on pairwise constrained (gray) and individual unconstrained (green) samples leveraging weak and strong data augmentations following the FixMatch approach [33]. A cross-entropy loss function is used as an auxiliary loss function for the pseudo-labeled unlabeled samples.



(b) ConstraintMatch combines pairwise training on constrained (gray) and unconstrained (yellow) samples leveraging weak and strong data augmentations. It extends the naive pseudo-labeling baseline by the generation of pseudo-constraints from *informative* pseudo-labels to overcome the confirmation bias as detailed in the methods section of the paper.

Fig. 7: Illustration of a) the naive pseudo-labeling baseline and b) ConstraintMatch.

the predicted clusters with the highest assigned probability are subsequently chosen as pseudo-labels. This confidence-based selection criterion is depicted as \mathcal{I}^{conf} in Fig. 7a and the associated threshold $\tau \in [0, 1]$ is a hyperparameter that we tuned on the validation set as described above and listed in Table IV. Predictions over strong augmented versions of these samples then serve as input for an auxiliary cross-entropy loss function, referred to as \mathcal{L}^{CE} in Fig. 7a. Similar to ConstraintMatch, the constrained loss \mathcal{L}^{cons} is calculated over model predictions on pairwise samples and their corresponding constraint annotations. The naive pseudo-labeling baseline is then trained via the final loss $\mathcal{L} = \mathcal{L}^{cons} + \lambda \mathcal{L}^{CE}$ as a weighted linear combination of both losses where hyperparameter λ controls the impact of the unconstrained samples.