

Beyond the Final Linear Layer: Enhancing Decision Boundaries

Michael Majurski
Information Technology Lab, NIST
michael.majurski@nist.gov

Sumeet Menon
University of Maryland, Baltimore County

David Chapman
University of Miami

Abstract

[TODO: (majurski) rewrite abstract once paper is done to fill in details] SSL leverages an abundance of unlabeled data to improve deep learning based model performance under limited training data regimes. This paper presents a novel extension to any image classification architecture which improves accuracy in low-label regimes. We extend the FixMatch [24] training scheme with our novel last layers and demonstrate test accuracy improvement. The novelty consists of 2 elements: first we replace the last linear layer with a GMM trained via backprop, and second, we impose class-wise constraints on the embedding space the GMM operates on. These methods match published SOTA 250 label CIFAR-10 [14] results and come close to matching SOTA in the 40 label regime without the significant model complexity of methods like SimMatchV2 [34]. Our method achieves 94.8% and 94.2% accuracy with 250 and 40 CIFAR-10 labels respectively. **[TODO: cleanup the repo]** Our code is available at: <https://github.com/mmajurski/ssl-gmm>

1. Introduction

SSL leverages an abundance of unlabeled data to improve deep learning based model performance under limited training data regimes [12, 17, 35]. Image classification has become a playground for exploring new SSL ideas. The early successes of deep learning based methods relied on large annotated datasets to enable models to learn the relevant features to perform the task, i.e. image classification build on top of ImageNet [9]. With data annotation becoming a significant bottleneck, especially in application domains outside of the standard benchmarks, another learning paradigm was needed.

There are several flavors of SSL. Contrastive learning methods leverage the intuition that similar instances should be close in the representation space, while different in-

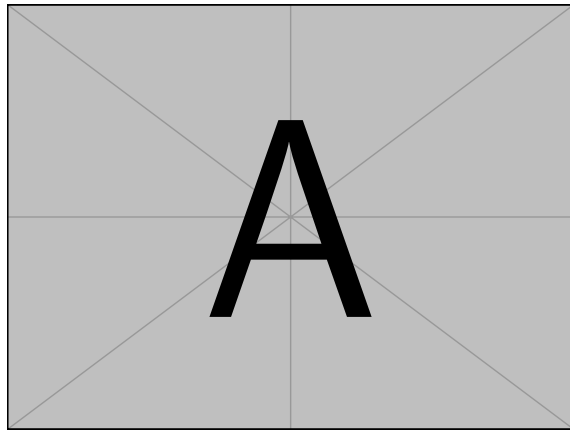


Figure 1. High level overview of our method. **[TODO: (majurski) Create this figure.]**

stances are farther apart [16, 29]. Consistency regularization borrows the intuition that modified views of the same instance should have similar representations and predictions [13, 15, 24, 32]. Pseudo-labeling methods like FixMatch [24] fall within the consistency regularization domain.

This work argues that pseudo-labeling methods can be improved with better calibration of the network logits used to filter the pseudo-labels into reliable and unreliable. Neural networks are known to be overconfident in their predictions [26], and this affects the pseudo-labeling process. Potentially allowing for the inclusion of more incorrect pseudo-labels any specific logit threshold would otherwise have. This work demonstrates the better calibrated replacements for a model's final linear layer can improve the final accuracy of pseudo-labeling based SSL algorithms in very label scarce regimes. This work proposes:

1. Replacing the final linear (fully connected) layer of the neural network with either kmeans [10] or axis-aligned differentiable Gaussian Mixture Model (GMM) trained via back prop, both of which have explicit modeling of class cluster centroids.

2. We explore various constraints on how the embedding space should be structured by adding penalties if the per-class clustering does not conform to between 0 and 4 of the first gaussian moments being identity/zero.
3. We demonstrate that increasing the specificity of how the embedding space should be structure negatively impacts model performance. **[TODO: (chapman) need citation for gaussian moments]**

This paper contributes a simple easy to implement improvement to pseudo-labeling methods where very few annotations are available, replace the last linear layer with a kmeans layer which explicitly models class cluster centers. We demonstrate this methodology using CIFAR-10 and CIFAR-100 [14] with 40 and 400 labels respectively. Additionally, we explore and demonstrate that high level prescriptive constraints on the embedding space produce significantly worse outcomes than allowing the embedding space to take on whatever emergent structure the training process produces. Finally, because the embedding constraint penalties are applied to all unlabeled data and not just the valid pseudo-labels, our method extracts training signal from every unlabeled data point, unlike FixMatch [24] and other methods which only learn from the valid pseudo-labels.

[TODO: (majurski) build t-SNE plots of the embedding spaces for the best models (1 per configuration)]

[TODO: (majurski) purge bibliography of arxiv pre-prints where possible, replacing with their peer reviewed equivalents]

2. Related Work

Semi-Supervised learning has recently show great progress in learning high quality models, in some cases match fully supervised performance, for a number of benchmarks [32]. The goal of SSL is to produce a trained model of equivalent accuracy to fully supervised training, with vastly reduced data annotation requirements.

2.1. Pseudo-Labeling

Self-supervised learning was among the initial approaches employed in the context of semi-supervised learning to annotate unlabeled images. This technique involves the initial training of a classifier with a limited set of labeled samples and incorporates pseudo-labels into the gradient descent process, exceeding a predefined threshold [8, 18, 19, 23, 30, 31]. A closely related method to self-training is co-training, where a given dataset is represented as two distinct feature sets [4]. These independent sample sets are subsequently trained separately using two distinct models, and the sample predictions surpassing predetermined thresholds are utilized in the final model training process [4, 21]. A notably advanced approach to pseudo-labeling is the Mean

Teacher algorithm [25], which leverages exponential moving averages of model parameters to acquire a notably more stable target prediction. This refinement has a substantial impact on enhancing the convergence of the algorithm.

2.2. Consistency Regularization

Consistency regularization operates on the premise that when augmenting an unlabeled sample, its label should remain consistent. This approach implicitly enforces a smoothness assumption, promoting coherence between unlabeled samples and their basic augmentations [27]. In other words, the model should be able to predict the unlabeled sample x exactly the same way it predicts the class for Augmented(x) [2, 3, 20, 24]. In addition to evaluating image-wise augmentations, recent research has demonstrated that incorporating class-wise and instance-based consistencies yields superior performance outcomes [16, 33]. Similarly, using consistencies between augmentations, of the predictions and low-dimensional embeddings of the strong and weak augmentations of the unlabeled images in a graph based setup has shown improvement over class-wise and instance-based consistencies [34]. Finally, pseudo-labeling filtering based on consistence between strongly augmented views shows convergence improvement [13].

2.3. Embedding Clustering/Constraints

Several papers have attempted to improve the quality of pseudo-labels to either improve the final model accuracy, improve the rate of convergence, or avoid confirmation bias [1]. Rizve et al. [22] explores how uncertainty aware pseudo-label selection/filtering can be used to reduce the label noise. Incorrect pseudo-labels can be viewed as a network calibration issue [22] where better network logit calibration might improve results [28]. Other work has attempted to improve the pseudo-labeling process by imposing curriculum [32] or by including a class-aware contrastive term [29]. Previous work has leveraged the concept of explicit class cluster centers for conditioning semantic similarity [33]. Recent work has extended purely clustering based methods like DINO [7] into semi-supervised methods [11].

3. Methodology

In this section, we explore our proposed final layer replacements and embedding space constraints. FixMatch [24] is a simple, well performing, SSL algorithm. As such it serves as a good comparison point for exploring the effect of our contributions. Our methodology is based upon the published FixMatch [24] algorithm, with identical hyperparameters unless otherwise stated. We add to FixMatch an early-stopping condition when the model has not improved for 40 epochs (where epoch size is 1024 batches as in FixMatch), with 2 learning rate reductions by a factor of 0.2

instead of configuring a fixed number of training epochs with a cosine learning rate decay. Additionally, we employ a cyclic learning rate scheduler to vary the learning rate by a factor of ± 2.0 within each epoch to make training less dependent upon exact learning rate value.

Both the linear layer replacements and the embedding constraints explored herein represent increasing levels of prescription about how the final embedding space should be arranged compared to a final linear layer. The idea of leveraging clusters in embedding space is not new [5, 6], but we extend the core idea with a novel differentiable model of learned cluster centroids and GMMs.

3.1. Alternative Final Layers

A limitation of traditional final activation layers such as softmax is that they are fully discriminative; i.e. they estimate the posterior $p(Y|X)$, but do not attempt to model the sample distribution $p(X)$ or the joint probabilities $p(Y, X)$. To overcome this limitation, we present two semi-parametric final activation layers (a) the Axis Aligned GMM (AAGMM) layer, and (b) an equal variance version of AAGMM that we henceforth call the KMeans activation layer due to the similarity of the objective function with a gradient based KMeans.

These activation layers are fully differentiable and integrated into the neural network architecture as a module in the same way as a traditional final linear layer. As such, they do not require or depend on external training and do not use expectation maximization. They are drop in replacements for the final linear layer.

Importantly, these activation layers exhibit both discriminative and generative properties. The neural network model $F(X; \theta_F)$ transforms the data X into a latent space $Z = F(X; \theta_F)$, and the final activation layer estimates the probability densities $p(X)$, $p(Y; X)$ and $p(Y|X)$ by fitting a parametric model to the latent representation Z .

3.1.1 Axis Aligned Gaussian Mixture Model Layer

The AAGMM layer defines a set of K trainable clusters, one cluster per label category. Each cluster $k = 1 \dots K$ has a cluster center μ_k and cluster covariance Σ_k . The prior probability of any given sample X_i is defined by the mixture of cluster probability densities over the latent representation Z_i as follows,

$$p(X_i) = \sum_{k=1}^K \mathcal{N}(Z_i, \mu_k, \Sigma_k) \quad (1)$$

$$\text{where } Z_i = F(X_i, \theta_F)$$

Where $\mathcal{N}(Z_i, \mu_k, \Sigma_k)$ represents the multivariate gaussian pdf with centroid μ and covariance Σ_k . AAGMM is

axis aligned because Σ_k is a diagonal matrix, as such the Normal pdf simplifies to the joint density of the pdfs along each of the D axes as follows,

$$\mathcal{N}(X_i, \mu_k, \Sigma_k) = \prod_{d=1}^D \frac{1}{\sigma_{k,d} \sqrt{2\pi}} \exp\left(-\frac{(Z_{i,d} - \mu_{k,d})^2}{\sigma_{k,d}^2}\right) \quad (2)$$

$$\text{where } \sigma_{k,d}^2 = \Sigma_{k,d,d}$$

As there is one cluster per label category, the joint probability for sample i with label assignment k , $p(Y_k, X_i)$ is the given by the normal pdf of the k^{th} cluster,

$$p(Y_{i,k}, X_i) = \mathcal{N}(Z_i, \mu_k, \Sigma_k) \quad (3)$$

By simple Bayesian identity, the posterior probability $\hat{Y}_k = p(Y_k|X_i)$ can therefore be inferred from eq 1 and 3 as follows,

$$\hat{Y}_{i,k} = p(Y_{i,k}|X_i) = \frac{p(Y_{i,k}, X_i)}{p(X_i)} \quad (4)$$

3.1.2 KMeans Layer

The KMeans final layer is a more restrictive form of the AAGMM layer, in the sense that we impose an additional constraint that the gaussian covariance matrix Σ_k for each cluster center k is the $[D \times D]$ identity matrix. This constraint yields spherical cluster centers similar to how the traditional KMeans algorithm also assumes spherical clusters.

3.2. Semi-parametric Embedding Constraints

We introduce and evaluate a series of embedding constraints based on the Method of Moments (MoM). As our goal is semi-supervised classification, the primary goal is to minimize the expected difference between Y and \hat{Y} through cross entropy loss. However, if one were to train the model without any embedding constraints, then it is possible that the model could learn a good decision boundary for the posteriors $p(Y|X)$ but without actually modeling the probability of samples $p(X)$. The first four gaussian moments being identity/zero increasingly constrain the model latent embedding space.

3.2.1 Moment1

This constraint encourages the learned cluster means to have zero distance to the observed centroids. Define $pred_k$ as the set of indices $i = 1 \dots K$ such that are predicted to be of category k as follows,

$$pred_k = \{i \mid \operatorname{argmax}_j \hat{Y}_{i,j} \text{ is equal to } k\} \quad (5)$$

This embedding constraint ensures that cluster centers μ_k approximate the sample distribution of latent points that are part of cluster k , $E(Z_i | \hat{Y}_{i,k} = 1)$. The $L2$ norm is used as a loss function to constrain cluster centers as follows,

$$\mathcal{L}_{L2} = \sum_{k=1}^K \left\| \mu_k - \bar{Z} \right\|_2 \quad (6)$$

where $\bar{Z} = E_{i \in \text{pred}_k} Z_i$

3.2.2 Moment2

Additionally, the $L2$ norm was also used to constraint the cluster covariance terms to the identity matrix. Although the AAGMM assumes a diagonal covariance matrix Σ_k , the mini-batch sample estimate of the covariance may exhibit off-diagonal terms, which represent

3.2.3 Moments 3 and 4

Moment3 encourages clusters to have zero skew.

Moment4 encourages cluster kurtosis to be the identity matrix.

[TODO: figure out how to talk about the moments 3 and 4, since all of our testing indicates that the GPU memory requirement is horrifying, and they are unstable as hell. As such they are not useful but as a baseline to compare the feasibility.]

4. Experiments

We evaluate our drop in final linear layer replacements on common SSL

[TODO: (majurski) continue writing here]

4.1. Ablation Study

5. Conclusions

References

- [1] Eric Arazo, Diego Ortego, Paul Albert, Noel E O'Connor, and Kevin McGuinness. Pseudo-labeling and confirmation bias in deep semi-supervised learning. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2020. 2
- [2] David Berthelot, Nicholas Carlini, Ekin D Cubuk, Alex Kurakin, Kihyuk Sohn, Han Zhang, and Colin Raffel. Remixmatch: Semi-supervised learning with distribution alignment and augmentation anchoring. *arXiv preprint arXiv:1911.09785*, 2019. 2
- [3] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A Raffel. Mixmatch: A holistic approach to semi-supervised learning. *Advances in Neural Information Processing Systems*, 32, 2019. 2
- [4] Avrim Blum and Tom Mitchell. Combining labeled and unlabeled data with co-training. In *Proceedings of the eleventh annual conference on Computational learning theory*, pages 92–100, 1998. 2
- [5] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *Proceedings of the European conference on computer vision (ECCV)*, pages 132–149, 2018. 3
- [6] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural information processing systems*, 33:9912–9924, 2020. 3
- [7] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 2
- [8] Olivier Chapelle, Bernhard Scholkopf, and Alexander Zien. Semi-supervised learning (chapelle, o. et al., eds.; 2006)[book reviews]. *IEEE Transactions on Neural Networks*, 20(3):542–542, 2009. 2
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 1
- [10] Debidatta Dwibedi, Yusuf Aytar, Jonathan Tompson, Pierre Sermanet, and Andrew Zisserman. With a little help from my friends: Nearest-neighbor contrastive learning of visual representations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9588–9597, 2021. 1
- [11] Enrico Fini, Pietro Astolfi, Karteek Alahari, Xavier Alameda-Pineda, Julien Mairal, Moin Nabi, and Elisa Ricci. Semi-supervised learning made simple with self-supervised clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3187–3197, 2023. 2
- [12] Mohamed Farouk Abdel Hady and Friedhelm Schwenker. Semi-supervised learning. *Handbook on Neural Information Processing*, pages 215–239, 2013. 1
- [13] Jiwon Kim, Youngjo Min, Daehwan Kim, Gyuseong Lee, Junyoung Seo, Kwangrok Ryoo, and Seungryong Kim. Comatch: Semi-supervised learning with confidence-guided consistency regularization. In *European Conference on Computer Vision*, pages 674–690. Springer, 2022. 1, 2
- [14] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009. 1, 2
- [15] Doyup Lee, Sungwoong Kim, Ildoo Kim, Yeongjae Cheon, Minsu Cho, and Wook-Shin Han. Contrastive regularization for semi-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3911–3920, 2022. 1
- [16] Junnan Li, Caiming Xiong, and Steven CH Hoi. Comatch: Semi-supervised learning with contrastive graph regularization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9475–9484, 2021. 1, 2

- [17] Yu-Feng Li and De-Ming Liang. Safe semi-supervised learning: a brief introduction. *Frontiers of Computer Science*, 13: 669–676, 2019. 1
- [18] Ioannis E Livieris, Konstantina Drakopoulou, Vassilis T Tampakas, Tassos A Mikropoulos, and Panagiotis Pintelas. Predicting secondary school students’ performance utilizing a semi-supervised learning approach. *Journal of educational computing research*, 57(2):448–470, 2019. 2
- [19] David McClosky, Eugene Charniak, and Mark Johnson. Reranking and self-training for parser adaptation. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pages 337–344, 2006. 2
- [20] Aamir Mustafa and Rafał K Mantiuk. Transformation consistency regularization—a semi-supervised paradigm for image-to-image translation. In *European Conference on Computer Vision*, pages 599–615. Springer, 2020. 2
- [21] V Jothi Prakash and Dr LM Nithya. A survey on semi-supervised learning techniques. *arXiv preprint arXiv:1402.4645*, 2014. 2
- [22] Mamshad Nayeem Rizve, Kevin Duarte, Yogesh S Rawat, and Mubarak Shah. In defense of pseudo-labeling: An uncertainty-aware pseudo-label selection framework for semi-supervised learning. *arXiv preprint arXiv:2101.06329*, 2021. 2
- [23] Chuck Rosenberg, Martial Hebert, and Henry Schneiderman. Semi-supervised self-training of object detection models. In *2005 Seventh IEEE Workshops on Applications of Computer Vision (WACV/MOTION’05) - Volume 1*, pages 29–36, 2005. 2
- [24] Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in neural information processing systems*, 33:596–608, 2020. 1, 2
- [25] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in neural information processing systems*, 30, 2017. 2
- [26] Hongxin Wei, Renchunzi Xie, Hao Cheng, Lei Feng, Bo An, and Yixuan Li. Mitigating neural network overconfidence with logit normalization. In *International Conference on Machine Learning*, pages 23631–23644. PMLR, 2022. 1
- [27] Qizhe Xie, Zihang Dai, Eduard Hovy, Thang Luong, and Quoc Le. Unsupervised data augmentation for consistency training. *Advances in neural information processing systems*, 33:6256–6268, 2020. 2
- [28] Chen Xing, Sercan Arik, Zizhao Zhang, and Tomas Pfister. Distance-based learning from errors for confidence calibration. In *International Conference on Learning Representations*, 2020. 2
- [29] Fan Yang, Kai Wu, Shuyi Zhang, Guannan Jiang, Yong Liu, Feng Zheng, Wei Zhang, Chengjie Wang, and Long Zeng. Class-aware contrastive semi-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14421–14430, 2022. 1, 2
- [30] David Yarowsky. Unsupervised word sense disambiguation rivaling supervised methods. In *33rd annual meeting of the association for computational linguistics*, pages 189–196, 1995. 2
- [31] Xiaohua Zhai, Avital Oliver, Alexander Kolesnikov, and Lucas Beyer. S4l: Self-supervised semi-supervised learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1476–1485, 2019. 2
- [32] Bowen Zhang, Yidong Wang, Wenxin Hou, Hao Wu, Jindong Wang, Manabu Okumura, and Takahiro Shinozaki. Flexmatch: Boosting semi-supervised learning with curriculum pseudo labeling. *Advances in Neural Information Processing Systems*, 34:18408–18419, 2021. 1, 2
- [33] Mingkai Zheng, Shan You, Lang Huang, Fei Wang, Chen Qian, and Chang Xu. Simmatch: Semi-supervised learning with similarity matching, 2022. 2
- [34] Mingkai Zheng, Shan You, Lang Huang, Chen Luo, Fei Wang, Chen Qian, and Chang Xu. Simmatchv2: Semi-supervised learning with graph consistency. *arXiv preprint arXiv:2308.06692*, 2023. 1, 2
- [35] Xiaojin Zhu and Andrew B Goldberg. *Introduction to semi-supervised learning*. Springer Nature, 2022. 1