

A Method of Moments Embedding Constraint for Generative Final Layers

Michael Majurski
Information Technology Lab, NIST
University of Maryland, Baltimore County
`michael.majurski@nist.gov`

Parniyan Farvardin
University of Miami
`pxf291@miami.edu`

Sumeet Menon
University of Maryland, Baltimore County
`sumeet1@umbc.edu`

David Chapman
University of Miami
`dchapman@cs.miami.edu`

Abstract

Discriminative deep learning models with a linear+softmax final layer have a problem: the latent space only predicts the conditional probabilities $p(y|x)$ but not the full joint distribution $p(y,x)$, which necessitates a generative approach. The conditional probability cannot detect outliers, causing sensitivity to them in softmax networks. This exacerbates model over-confidence impacting many problems: from hallucinations, to confounding biases, and dependence on large datasets. We introduce a novel embedding constraint based on the Method of Moments (MoM). We investigate the use of polynomial moments ranging from 1st through 4th order hyper-covariance matrices. Furthermore, we use this embedding constraint to train an Axis-aligned Gaussian Mixture Model (AAGMM) final layer, which learns not only the conditional, but also the joint distribution of the latent space. We demonstrate our approach by extending FixMatch based semi-supervised image classification. We find our MoM constraint with the AAGMM layer is able to match or improve upon the reported FixMatch accuracy, while also modeling the joint distribution, thereby reducing outlier sensitivity. Future work explores potential applications for this layer and embedding constraint, and how/why this MoM technique can overcome theoretical limitations of other existing methods including the approximate KL-divergence constraint of variational autoencoders. Code is available at: https://github.com/*****

1. Introduction

The majority of deep classifiers rely on a softmax final activation layer which predicts the conditional probability $p(y|x)$. When that layer receives input x , the model predicts a soft psuedo-distribution of labels y which argmax

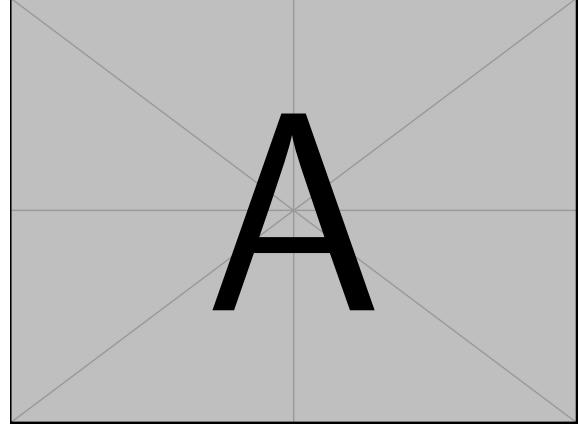


Figure 1. Schematic of the outlier problem, and how generative modeling of the joint probability can improve the situation.

can convert into a hard label. If x is far from the decision boundary, then by definition, softmax assigns a prediction y with high confidence. This works well for inlier samples, well represented by the training distribution. However, when presented with an outlier x , it is likely x will be far from the decision boundary (**[TODO: Figure ?]**). This means that softmax perceptrons, by definition, overconfidently hallucinate when given unexpected inputs [30]. Most deep classifiers use softmax without a safety net and thereby over-confidently predict y . Ideally, when input x is far from the decision boundary and training exemplars, the model should not be confident about the output class label y .

Replacing softmax with a generative method that models the joint probability $p(y,x)$ can improve the capability of deep classifiers. Models using a final layer capable of learning the joint probability $p(y,x)$ can infer the conditional $p(y|x)$. More importantly, such a layer can also infer

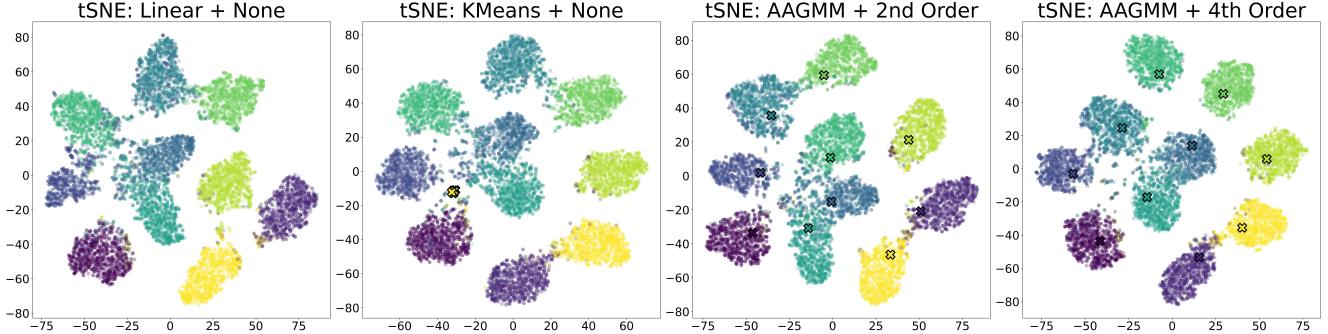


Figure 2. t-SNE plot of the latent embedding space for various final layers with different MoM embedding constraints.

the prior probability $p(x)$. Thus, if x is an unexpected input, then such layer can flag the input as a low-probability outlier, rather than confidently predicting a label.

Although there is related work in creating a generative final activation layer, there is an open question of how best to train such a layer in a deep learning context. **[TODO: (chapman) Citations for the related work...?]** The strawman’s approach would be to minimize the cross entropy between y and y_{pred} as usual. Figure 2 shows actual t-SNE plots for semi-supervised CIFAR-10 [13] classification that beautifully illustrate why such an approach will not work as intended. Figure 2 (top-right) shows the latent space of a classification model that is 93% accurate on the test set, yet the cluster centers are not even remotely located anywhere nearby the data points at all. If one only cares about predictive performance, then one might conclude such a strange model would suffice. But if one needs a robust model, then simply fitting the decision boundary is not adequate, we need to learn the full latent distribution. Figure 2 (bot) shows that with our MoM embedding constraints, the exact same model can achieve comparable if not better accuracy, but with the added benefit of modeling the true location of cluster centers in the latent space.

This work argues that Semi-supervised learning (SSL) pseudo-labeling methods can be improved with better calibration of the network logits used to filter reliable pseudo-labels from the unreliable ones. This is equivalent to improving the accuracy of class inlier determination. SSL is an excellent application to test generative final activation layers which model the joint probability $p(y, x)$, as pseudo-labeling methods are sensitive to outliers.

SSL leverages an abundance of unlabeled data to improve deep learning based model performance under limited training data regimes [11, 16, 39]. Contrastive learning methods leverage the intuition that similar instances should be close in the representation space, while different instances are farther apart [15, 33]. Consistency regularization borrows the intuition that modified views of the same instance should have similar representations and predictions

[12, 14, 27, 36]. Pseudo-labeling methods like FixMatch [27] leverage the ideas of consistency regularization. This work contributes:

1. A replacement of the final linear+softmax final activation layer of the neural network with either an axis-aligned differentiable Gaussian Mixture Model (AAGMM) or an equal variance version named KMeans trained via back prop, both of which have explicit modeling of class cluster centroids.
2. A demonstration on the latent embedding space reacts to various constraints on how it 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 [23].

We demonstrate this methodology using the standard CIFAR-10 benchmark dataset with 40 and 250 labels[13]. 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, an improvement on baseline pseudo-labeling methods (like FixMatch [27]) which only learn from the valid pseudo-labels.

2. Related Work

Semi-Supervised learning has shown great progress in learning high quality models, in some cases matching fully supervised performance for a number of benchmarks [36]. The goal of SSL is to produce a trained model of equivalent accuracy to fully supervised training, with vastly reduced data annotation requirements. Doing so relies on accurately characterizing inlier vs outlier unlabeled samples.

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 gra-

dient descent process, exceeding a predefined threshold [8, 17, 18, 20, 26, 34, 35]. 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, 24]. A notably advanced approach to pseudo-labeling is the Mean Teacher algorithm [28], 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 [31]. In other words, the model should be able to predict the unlabeled sample x exactly the same way it predicts the class for $\text{Augmented}(x)$ [2, 3, 21, 27]. In addition to evaluating image-wise augmentations, recent research has demonstrated that incorporating class-wise and instance-based consistencies yields superior performance outcomes [15, 37]. Similarly, using consistencies between the predictions and low-dimensional embeddings from the unlabeled image strong and weak augmentations in a graph based setup demonstrates improvement over class-wise and instance-based consistencies [38]. Finally, pseudo-labeling filtering based on consistency between strongly augmented views, gaussian filtering and embedding based nearest neighbor filtering shows convergence improvement [12, 19].

2.3. Latent Embedding Constraints

Several papers have attempted to enhance 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. [25] 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 [25] where better network logit calibration might improve results [32]. Improvements to the pseudo-labeling process have been demonstrated by imposing curriculum [36] or by including a class-aware contrastive term [33]. Leveraging the concept of explicit class cluster centers for conditioning semantic similarity improves final model accuracy [37]. Additionally, improvements have been found in extended purely clustering based methods like DINO [7] into semi-supervised methods [10].

3. Methodology

In this section, we explore our proposed replacement final activation layers and our embedding space constraints. FixMatch [27] 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 [27] algorithm, with identical hyper-parameters unless otherwise stated. We extend FixMatch with a few minor training algorithm modifications explored in the Hyperparameters Section 4.1.

Both the linear layer replacements and the embedding constraints explored herein represent increasing levels of prescription about how the final latent embedding space should be arranged compared to a traditional linear layer. The idea of leveraging clusters in embedding space is not new [5, 6, 9], but we extend the core idea with a novel differentiable model with learned cluster centroids and MoM based constraints.

3.1. Alternate Final Layers

A limitation of traditional final activation layers such as linear+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 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_k and covariance Σ_k . AAGMM is axis aligned because Σ_k is a diagonal matrix, as such the axis-aligned multivariate normal pdf simplifies to the marginal product of Gaussians along each of the D axes as follows,

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

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

As there is one cluster per label category, the joint probability for sample i with label assignment k , $p(Y_{i,k}, X_i)$ 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 Bayesian identity, the posterior probability $\hat{Y}_k = p(Y_k | X_i)$ can therefore be inferred from eq 1 and 3 as follows,

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

The AAGMM layer is implemented as a normal PyTorch [22] module. It has two parameters updated by backprop. (1) the explicit cluster centers, a matrix $\text{num_classes} \times \text{embedding_dim}$ initialized randomly, and (2) the diagonal elements of the Σ matrix, randomly initialized in the range $[0.9, 1.1]$, which contains the diagonal elements of the GMM Sigma matrix for each cluster.

3.1.2 KMeans Layer

We also implement a KMeans final layer which is a more restrictive form of the AAGMM layer. The KMeans layer is additionally constrained such 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.

The KMeans layer is also implemented as a normal PyTorch [22] module. The explicit cluster centers is a learned parameter updated by backprop. See the published codebase for implementation details about the AAGMM and KMeans layers.

3.2. Method of Moments Embedding Constraints

We introduce and evaluate a series of embedding constraints based on the Method of Moments (MoM) [23] in order to fit the semi-parametric latent prior parameters. The latent prior $p(X_i)$ is calculated in equation 1 and then used to infer the posterior $p(Y_{i,k}|X_i)$. As usual, the posterior is trained using cross entropy loss. When embedding constraints are omitted, it is possible for the model to learn an accurate decision boundary for the posterior without modeling the latent prior.

Our novel last layer is semi-parametric, because the prior is a parametric model of the latent distribution Z which is the result of a neural network feature extraction $F(X; \theta)$. Therefore, attempting to fit the GMM directly to Z using Maximum Likelihood (ML) or simple Expectation Maximization (EM), is not appropriate, because doing so would fail to learn an appropriate feature space for discrimination. MoM solves these problems and is an appropriate strategy for semi-parametric models including ours.

The MoM relies on the use of *consistent estimators*, which asymptotically share sample and population statistics. Assume that z is a finite sample of n elements drawn from infinite population Z , then a series of P well-behaved sample statistics g_p should very closely approximate their k population statistic as follows,

$$\forall p = 1 \dots P \quad \frac{1}{n} \sum_{i=1}^n g_p(z_i) \approx E(g_p(Z)) \quad (5)$$

We can therefore constrain the latent representation of our model to approximate an independent joint Gaussian distribution. In the univariate gaussian case, the p^{th} order centralized moment constraint is the following.

$$E[(Z - \mu)^p] = \begin{cases} 0 & \text{if } p \text{ is odd} \\ \sigma^p(p-1)!! & \text{if } p \text{ is even} \end{cases} \quad (6)$$

By this formula, the univariate unit gaussian has mean 0, standard deviation 1, skew 0, and kurtosis 3.

In the joint multivariate case, each dimension is independent by definition. As such, if we redefine Z , μ , and p to be all D dimensional, then the centralized joint gaussian moment can be defined as follows,

$$E[g_p(Z - \mu)] = E\left[\prod_{d=1}^D (Z_d - \mu_d)_d^p\right] \quad (7)$$

Due to independence of the axes, this moment can be represented as a product of univariate moments of the individual gaussians as follows,

$$E\left[\prod_{d=1}^D (Z_d - \mu_d)_d^p\right] = \prod_{d=1}^D E[(Z_d - \mu_d)_d^p] \quad (8)$$

The error (loss) term associated with the embedding constraint for any moment p is equal to the L2 difference between the sample and population statistics as follows,

$$\varepsilon_p = \left(\frac{1}{n} \sum_{i=1}^n g_p(z_i) - E(g_p(Z)) \right)^2 \quad (9)$$

Some moments are more important than others, and must be weighted more heavily. First order moments are simply the sample mean, and should be given the greatest weight as an embedding constraint. The second order moments form a sample covariance matrix, which ideally should be equal to the identity matrix, but the diagonal terms should be given greater weight than the off-diagonal terms. This is because, in a $D \times D$ covariance matrix, there are $D(D-1)$ off diagonal terms, but only D , diagonal terms. The p^{th} order sample moments form a $p-1$ dimensional hyper-covariance matrix, with terms residing on the intersection of anywhere between 0 and $p-1$ hyper-diagonals. To prevent over-representation of off-diagonal terms and encourage representation of on-diagonal terms, the loss function we use for any given moment term is inversely proportional to the number moment terms that share the same number of hyper-diagonals. This heuristic weighting scheme ensures that the overall contribution of each moment order is not overly influenced by the off-diagonal terms, and that the error weighting is therefore diagonally dominant.

4. Experiments

We evaluate both AAGMM and KMeans linear layer replacements and the embedding space constraints using our modified FixMatch[27] on the common SSL benchmarks CIFAR-10 [13] at 40 and 250 labels (4 and 25 labels per class). We randomly selected 5 seed a priori for evaluation. For each algorithm configuration tested one model was trained per seed. During each run, the required number of labeled samples are drawn without replacement from the training population of the dataset in a deterministic manner (reproducible with the same seed). All data not in this labeled subset is used as unlabeled data (i.e. the labels are discarded).

AAGMM+None on CIFAR-10 at 40 Labels

Run Number	1	2	3	4	5
Test Accuracy %	94.6	92.8	92.8	89.9	86.4

Table 1. Test accuracy for showing the run-to-run variance depending on the quality of the 40 labels selected from the full population.

As prior work [27] has noted, the resulting model quality is highly variable when only 4 samples are selected per class, as the quality and usefulness of the specific 4 samples

can vary drastically. Table 1 shows final test accuracy for the 5 AAGMM model runs with no embedding constraints, with the accuracy varying from 86% to 94%. Due to the potential for significant variance in the final model test accuracy, it can be informative to compare mean performance with max performance over the $N = 5$ runs. This explores both how well a method can be expected to do on average with random label sampling, vs how well it can potentially do with a more representative subset of labeled data.

4.1. Hyper-Parameters

The CIFAR-10 models were trained with the standard benchmark WideResNet28-2 architecture. Matching the FixMatch [27] hyper-parameters, we use SGD with Nesterov momentum and $\lambda_u = 1$, $\beta = 0.9$, $\tau = 0.95$, $\mu = 7$, $B = 64$, and epoch size = 1024 regardless of the number of images in the labeled dataset. We modify the stock FixMatch [27] training algorithm with an early-stopping condition when the model has not improved for 50 epochs (where epoch size is defined as 1024 batches). We also use a $learning_rate(\eta) = 0.01$ and an exponential moving average of the model weights with a decay 0.999. Our replacement of a fixed number of training steps with an early stopping criteria prevents the use of a cosine decay schedule. Therefore, we replace that with a plateau learning rate scheduler which multiplies the learning rate by 0.2 every time the early stopping criteria is met (before being reset) for a max of 2 reductions. 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. Additionally, due to the higher training instability of the AAGMM layers compared to a linear layer, if the training loss is greater than 1.0 we clip the gradient norm to 1.0. Despite computing the AAGMM and embedding constraints as numerically stable as we could, they are still less stable during backprop than a simple linear layer.

As part of this work we explore how various latent embedding dimensionality affects the generative model linear layer replacement. As such we modify the model architectures with a single additional linear layer before the output to project the base model embedding dimension (128 for WideResNes28-2) down to a reduced 8 dimensional space. We evaluate our method both with and without this reduced embedding space, to explore whether the generative capabilities of the AAGMM improve when its not fighting the curse of dimensionality. Results listed with an embedding dimensionality of 128 do not include the additional linear layer which reduces the latent dimensionality. Therefore, results with 128D embedding represents an unmodified network architecture.

Due to exponential GPU memory requirements with each successive MoM moment, only the 8D embedding can

CIFAR-10 Mean Test Accuracy

Last Layer	Emb Dim	40 Labels (5 trials)			
Embedding Constraint		None	1st Order	2nd Order	3rd Order
Linear (i.e. FullyConnected)	128	82.19 \pm 6.54			
	8	88.40 \pm 3.54			
AAGMM	128	91.23 \pm 2.89	89.23 \pm 2.50	90.22 \pm 3.42	
	8	88.34 \pm 2.34	82.39 \pm 9.96	87.97 \pm 3.34	82.51 \pm 9.42
KMeans	128	81.13 \pm 2.17	89.74 \pm 2.62	89.89 \pm 2.79	
	8	78.28 \pm 9.87	73.27 \pm 12.1	75.80 \pm 10.69	80.96 \pm 5.94
Last Layer	Emb Dim	250 Labels (5 trials)			
Embedding Constraint		None	1st Order	2nd Order	3rd Order
Linear (i.e. FullyConnected)	128	94.61 \pm 0.11			
	8	93.72 \pm 0.57			
AAGMM	128	94.09 \pm 0.34	94.20 \pm 0.62	94.42 \pm 0.14	
	8	94.17 \pm 0.57	94.01 \pm 0.73	94.33 \pm 0.37	93.77 \pm 0.90
KMeans	128	92.83 \pm 1.16	93.29 \pm 0.92	93.16 \pm 1.25	
	8	93.71 \pm 0.89	94.09 \pm 0.50	93.64 \pm 0.99	94.31 \pm 0.39

Table 2. Mean test accuracy % for CIFAR-10 SSL benchmark comparing various configurations of our method. The FixMatch results in the table is our reproduction of the published results, using our training pipeline modifications. For CIFAR-10 the WideResNet model used by FixMatch has an embedding size of 128 dimension. Due to exponential GPU memory requirements only the 8D embedding can operate with higher order MoM embedding constraints. Results for a given order of embedding constraint include all lower constraints. [TODO: (majurski) update table with final few missing runs to bring N=5]

operate with higher order MoM embedding constraints. Results for any given order of embedding constraint include all lower constraints.

4.2. CIFAR-10

The CIFAR-10 SSL benchmark was used to explore the full configuration space of our method. While both 40 and 250 label counts were used, the 250 label case SOTA is close to fully supervised accuracy. We include 250 performance to document our result is approximately equivalent to SOTA. The 40 label case provides a far more challenging task, though recent results have demonstrated accuracies that nearly match fully-supervised performance on CIFAR-10 (similar to 250 label CIFAR-10).

Table 2 summarizes the relative performance of our various configurations for both 40 and 250 labels. We reproduced FixMatch [27] using our hyper-parameters and didn't quite matching the published performance at 40 labels (250 labels matched). Hyper-parameter selection for our training algorithm is likely sub-optimal for baseline FixMatch. The "Linear (i.e. FullyConnected)" rows in table 2 represent the baseline fully connected linear last layer without additional embedding dimensionality projection.

For CIFAR-10 with 250 labels, all last layers perform reasonably close to semi-supervised SOTA, which itself is almost identical to the fully supervised CIFAR-10 test accuracy of 95.38% [29].

In addition to average performance, it is informative to

examine the max test accuracy over the $N = 5$ random trials to understand how well the algorithm can do, with samples that are representative of the larger dataset. Table 3 demonstrates that in the best case, the AAGMM can get within 1% of SOTA [38] performance.

The modeled cluster centers vary in quality between individual model runs of the AAGMM layer due to the stochasticity of the training process. Figure 3 top row showcases degenerate cluster centers. The bottom left AAGMM model learned cluster centers that are an ok approximation of the underlying data. However, the embedding constraints encourage cluster centers which are better aligned with the underlying data (Figure 3 bottom right). It is worth noting that we did not observed the KMeans layer learning non-degenerate cluster centers without an embedding constraint. In contrast, the AAGMM layer can, under some circumstances, learn viable cluster centers.

5. Discussion

[TODO: (chapman) need citation for VAE methods] The most similar embedding constraint to our MoM as commonly seen in deep learning is the use of variational inference, as applied in Variational Auto-Encoder models (VAE). VAE constrains the latent distribution using a KL-divergence penalty $D_{KL}(Q||P)$, where $Q(Z)$ is the variational distribution designed to approximate $p(Z|X)$, and P is the target distribution, typically the multivariate standard normal distribution. Although the use of KL-divergence is

CIFAR-10 Max Test Accuracy

Last Layer	Emb Dim	40 Labels (5 trials)				
Embedding Constraint		None	1st Order	2nd Order	3rd Order	4th Order
Linear (i.e. FullyConnected)	128	91.01				
	8	92.10				
AAGMM	128	94.64	91.33	92.74		
	8	90.40	93.25	92.11	91.95	92.01
KMeans	128	84.08	91.62	92.22		
	8	93.21	91.22	89.35	85.96	92.59

Table 3. Max test accuracy (%) for CIFAR-10 SSL benchmark with 40 labels comparing various configurations. This table shows the best-case performance of our various methods; without the effect of poorly representative labels selected for each class.

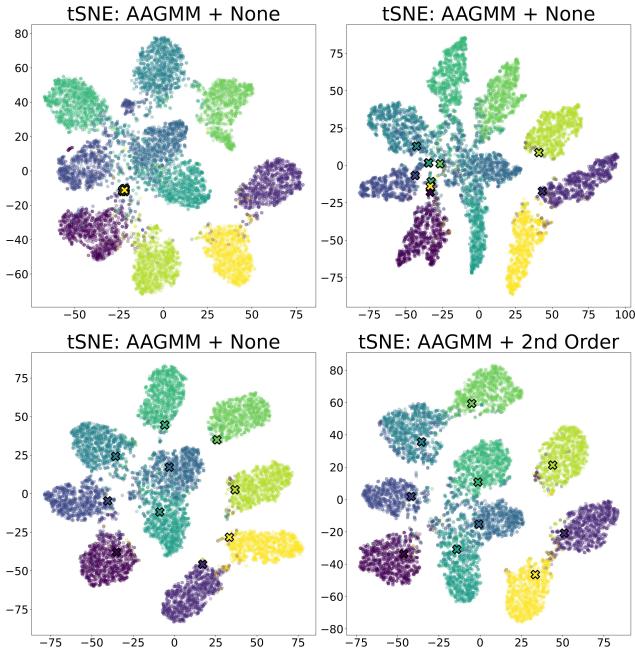


Figure 3. t-SNE plot of fully trained and reasonably accuracy AAGMM model’s latent embedding space (with the learned cluster centers marked with X’s). Depending on the run, the AAGMM cluster centers will be degenerate (top left), non-degenerate but still mis-aligned with the clusters (top right), acceptably aligned (bottom left), or well aligned with the underlying clusters when a 2nd order constraint is employed (bottom right).

strongly rooted in information theory, true calculation of KL-divergence is intractable, and thus approximations are often employed. As such, the choice of the variational distribution $Q(Z)$ hides a great deal of error in approximation, limiting the VAE embedding constraint to penalize simple differences in the shape of the feature space.

In most VAE models the variational distribution $Q(Z)$ constrains the first order (mean) and second order diagonal moments (standard deviation). Second order off-diagonal terms are usually discarded by VAE models. This means

that most practical VAE models have no ability to constrain even the linear orthogonality of the feature space, because the cross-covariance terms are omitted to ‘simplify’ KL-divergence. By comparison, our MoM constraints not only constrain linearly orthogonality, but also penalize more complicated non-linear shape descriptions including skew, kurtosis, as well as complex multivariate hyper-covariance descriptions of shape. By penalizing higher order off-diagonal terms, we can produce feature spaces that are not only linearly orthogonal, but closely resemble true statistical independence.

One significant downside to our MoM embedding constraints is the exponential GPU memory requirements. The unreasonable amount of GPU memory required limits their practical application. Additional optimization and/or avoiding the explicit creation of both the n th order moment and its target value on device would likely improve the usability.

Semi-supervised learning is highly sensitive to both which samples are selected for the labeled population [27] and the stochasticity of the training process. The stochasticity of training runs with the same starting seed (and identical labeled samples) will quickly diverge due to the randomness inherent in the training process. Anecdotally its worse in semi-supervised than fully supervised. To characterize this variance, and hence how much we should trust the error bars for our 5 runs, we took a few of our configurations and ran them $N = 5$ times with the same seed with CIFAR-10 at 40 labeled samples. This produced the potentially widely varying final model accuracies seen in Table 4.

Future work in this area should explore both accuracy improvements as well as implementation optimization to ensure the novel final layers are not prohibitively memory expensive to use. Additionally, one should explore how to best take advantage of the better behaved latent embedding space to improve data efficiency for model training.

We demonstrate a fully differentiable Axis-Aligned Gaussian Mixture Model with Method of Moments based latent embedding space constraints to improve the generative inlier/outlier performance of image classification deep

Identical Seed Runs on CIFAR-10 at 40 Labels

	AAGMM +2nd Order	KMeans +2nd Order	AAGMM +None	KMeans +None
Run 1	87.5	86.9	71.2	$xx.y$
Run 2	85.6	91.2	67.9	$xx.y$
Run 3	84.8	77.7	$xx.y$	$xx.y$
Run 4	85.3	87.6	$xx.y$	$xx.y$
Run 5	85.2	$xx.y$	$xx.y$	$xx.y$

Table 4. Test accuracy for independent training runs with the same random seed for various final layer configurations. All runs use the 128D (default) model latent embedding. All runs use the same labeled samples. [TODO: (majurski) update table with final few currently running experiments]

learning models. This preliminary work constructs those novel layers with the associated constraints, and demonstrates reasonable performance on challenging benchmark semi-supervised learning tasks.

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. [3](#)
- [2] 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. [3](#)
- [3] David Berthelot, Nicholas Carlini, Ekin D. Cubuk, Alex Kurakin, Kihyuk Sohn, Han Zhang, and Colin Raffel. Rematch: Semi-supervised learning with distribution matching and augmentation anchoring. In *International Conference on Learning Representations*, 2020. [3](#)
- [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. [3](#)
- [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. [3](#)
- [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. [3](#)
- [9] Joseph Enguehard, Peter O’Halloran, and Ali Gholipour. Semi-supervised learning with deep embedded clustering for image classification and segmentation. *Ieee Access*, 7: 11093–11104, 2019. [3](#)
- [10] 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. [3](#)
- [11] Mohamed Farouk Abdel Hady and Friedhelm Schwenker. Semi-supervised learning. *Handbook on Neural Information Processing*, pages 215–239, 2013. [2](#)
- [12] 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. [2, 3](#)
- [13] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, 2009. [2, 5](#)
- [14] 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. [2](#)
- [15] 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. [2, 3](#)
- [16] Yu-Feng Li and De-Ming Liang. Safe semi-supervised learning: a brief introduction. *Frontiers of Computer Science*, 13: 669–676, 2019. [2](#)
- [17] 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. [3](#)
- [18] 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. [3](#)
- [19] Sumeet Menon. *Semi-Supervised Expectation Maximization with Contrastive Outlier Removal*. PhD thesis, 2022. AAI29167191. [3](#)
- [20] Sumeet Menon, David Chapman, Phuong Nguyen, Yelena Yesha, Michael Morris, and Babak Saboury. Deep expectation-maximization for semi-supervised lung cancer screening. 2019. [3](#)
- [21] 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. [3](#)

- [22] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. 4
- [23] Karl Pearson. Method of moments and method of maximum likelihood. *Biometrika*, 28(1/2):34–59, 1936. 2, 4
- [24] V Jothi Prakash and Dr LM Nithya. A survey on semi-supervised learning techniques. *arXiv preprint arXiv:1402.4645*, 2014. 3
- [25] 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. In *International Conference on Learning Representations*, 2021. 3
- [26] 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. 3
- [27] 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. 2, 3, 5, 6, 7
- [28] 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. 3
- [29] Yidong Wang, Hao Chen, Qiang Heng, Wenxin Hou, Yue Fan, Zhen Wu, Jindong Wang, Marios Savvides, Takahiro Shinozaki, Bhiksha Raj, Bernt Schiele, and Xing Xie. Freematch: Self-adaptive thresholding for semi-supervised learning. In *The Eleventh International Conference on Learning Representations*, 2023. 6
- [30] 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
- [31] 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. 3
- [32] Chen Xing, Sercan Arik, Zizhao Zhang, and Tomas Pfister. Distance-based learning from errors for confidence calibration. In *International Conference on Learning Representations*, 2020. 3
- [33] 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. 2, 3
- [34] David Yarowsky. Unsupervised word sense disambiguation rivaling supervised methods. In *33rd annual meeting of the association for computational linguistics*, pages 189–196, 1995. 3
- [35] 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. 3
- [36] 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. 2, 3
- [37] Mingkai Zheng, Shan You, Lang Huang, Fei Wang, Chen Qian, and Chang Xu. Simmatch: Semi-supervised learning with similarity matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14471–14481, 2022. 3
- [38] Mingkai Zheng, Shan You, Lang Huang, Chen Luo, Fei Wang, Chen Qian, and Chang Xu. Simmatchv2: Semi-supervised learning with graph consistency. pages 16432–16442, 2023. 3, 6
- [39] Xiaojin Zhu and Andrew B Goldberg. *Introduction to semi-supervised learning*. Springer Nature, 2022. 2