

A Method of Moments Embedding Constraint and its Application to Semi-Supervised Learning

Anonymous arXiv submission

Paper ID 16555

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 outlier sensitivity 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/*****](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 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 re-

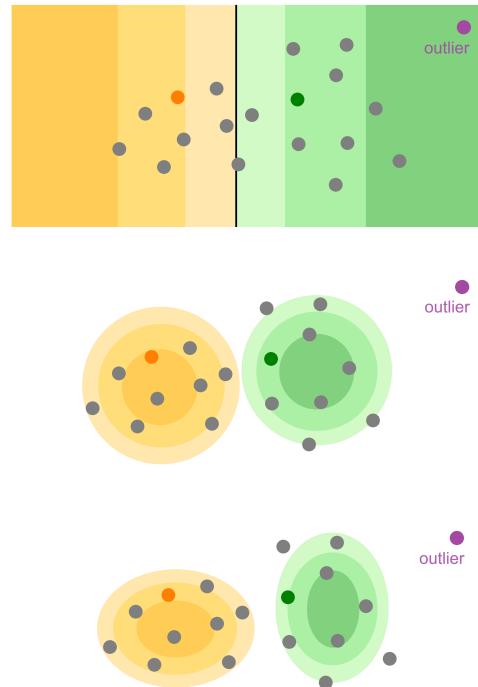


Figure 1. Schematic of the outlier problem, and how generative modeling of the joint probability can improve the situation. Prediction with (top) fully-supervised softmax (middle) semi-supervised KMeans and (bottom) semi-supervised AAGMM.

resented by the training distribution. However, when presented with an outlier x , it is likely x will be far from the decision boundary (Figure 1, top). Therefore softmax perceptrons, by definition, over-confidently hallucinate when given unexpected inputs [42]. 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

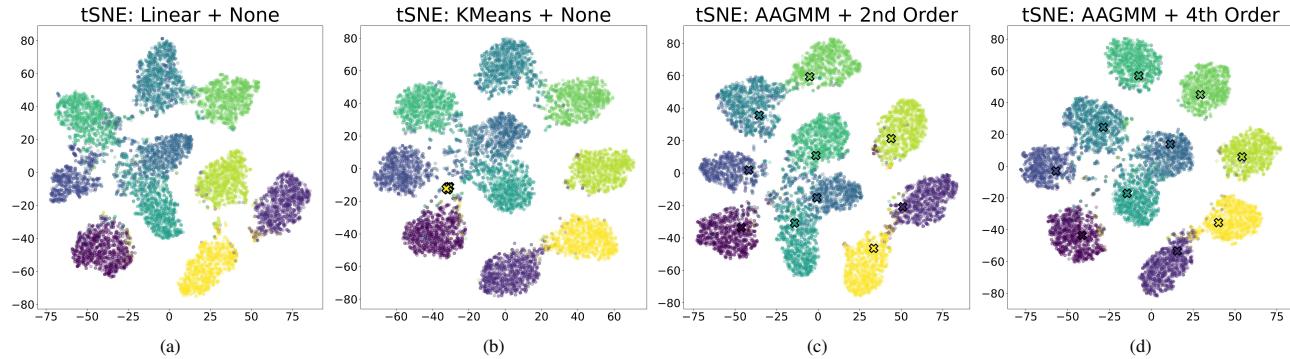


Figure 2. t-SNE[37] plot of the latent embedding space for various final layers with different MoM embedding constraints.

learning the joint probability $p(y, x)$ can infer the conditional $p(y|x)$. More importantly, such a layer can also infer 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.

Prior work has explored generative modeling for image classification [15, 17, 22]. Open questions remain of how to best train and utilize generative modeling within a deep learning context. The naive approach of minimizing cross entropy between y and y_pred will not work. Figure 2 (b) shows t-SNE plots for semi-supervised CIFAR-10 [18] image classification which illustrate why the naive approach will not work as intended. The t-SNE[37] plot of "KMeans + None" shows the latent space of a 93% accurate CIFAR-10 classification model. However, the explicitly modeled cluster centers (shown as X's) do not align with the underlying data. While that model has acceptable predictive performance, it does not accurately learn and represent the underlying training data. To construct a robust model, one cannot simply fit a decision boundary. The model needs to learn the full latent distribution. Figure 2 (d) demonstrates that with our proposed AAGMM final layer with 4th order MoM embedding constraints, the exact same model can achieve comparable (if not better accuracy), but with the added benefit of modeling the underlying data clusters 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 [13, 23, 52]. Contrastive learning methods leverage the intuition that similar instances should be close in the representation space, while different

instances are farther apart [20, 45]. Consistency regularization borrows the intuition that modified views of the same instance should have similar representations and predictions [14, 19, 35, 48]. Pseudo-labeling methods like FixMatch [35] leverage the ideas of consistency regularization. This work contributes:

1. A novel Method of Moments (MoM) based embedding constraint that enables the model to not only learn the decision boundary but also the latent joint distribution. Moreover, this constraint ensures that each latent cluster exhibits a well-behaved Gaussian shape.
2. 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 propagation, both of which have explicit modeling of class cluster centroids.
3. A demonstration of the latent embedding space response 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 [30].

We demonstrate this methodology using the standard CIFAR-10 benchmark dataset with 40 and 250 labels[18]. 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 [35]) which only learn from the valid pseudo-labels due to removing low confidence 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 [48]. The goal of SSL is to produce a trained model of equivalent accuracy to fully supervised training, with vastly reduced

124 data annotation requirements. Doing so relies on accurately
125 characterizing inlier vs outlier unlabeled samples.

126 2.1. Pseudo-Labeling

127 Self-supervised learning was among the initial approaches
128 employed in the context of semi-supervised learning to
129 annotate unlabeled images. This technique involves the
130 initial training of a classifier with a limited set of
131 labeled samples and incorporates pseudo-labels into the
132 gradient descent process, exceeding a predefined threshold
133 [9, 24, 25, 27, 34, 46, 47]. A closely related method to self-
134 training is co-training, where a given dataset is represented
135 as two distinct feature sets [4]. These independent sample
136 sets are subsequently trained separately using two distinct
137 models, and the sample predictions surpassing predetermined
138 thresholds are utilized in the final model training
139 process [4, 31]. A notably advanced approach to pseudo-
140 labeling is the Mean Teacher algorithm [36], which leverages
141 exponential moving averages of model parameters to
142 acquire a notably more stable target prediction. This refinement
143 has a substantial impact on enhancing the convergence
144 of the algorithm.

145 Several papers have attempted to enhance the quality
146 of pseudo-labels to either improve the final model accuracy,
147 improve the rate of convergence, or avoid confirmation
148 bias [1]. Rizve et al. [33] explores how uncertainty
149 aware pseudo-label selection/filtering can be used to re-
150 duce the label noise. Incorrect pseudo-labels can be viewed
151 as a network calibration issue [33] where better network
152 logit calibration might improve results [44]. Improvements
153 to the pseudo-labeling process have been demonstrated by
154 imposing curriculum [48] or by including a class-aware
155 contrastive term [45]. Leveraging the concept of explicit
156 class cluster centers for conditioning semantic similarity
157 improves final model accuracy [50]. Additionally, improvements
158 have been found in extended purely clustering based
159 methods like DINO [7] into semi-supervised methods [11].

160 2.2. Consistency Regularization

161 Consistency regularization operates on the premise that
162 when augmenting an unlabeled sample, its label should
163 remain consistent. This approach implicitly enforces a
164 smoothness assumption, promoting coherence between un-
165 labeled samples and their basic augmentations [43]. In
166 other words, the model should be able to predict the un-
167 labeled sample x exactly the same way it predicts the class
168 for $\text{Augmented}(x)$ [2, 3, 28, 35]. In addition to evaluating
169 image-wise augmentations, recent research has demon-
170 strated that incorporating class-wise and instance-based
171 consistencies yields superior performance outcomes [20,
172 50]. Similarly, using consistencies between the predictions
173 and low-dimensional embeddings from the unlabeled im-
174 age strong and weak augmentations in a graph based setup

175 demonstrates improvement over class-wise and instance-
176 based consistencies [51]. Finally, pseudo-labeling filtering
177 based on consistence between strongly augmented views,
178 gaussian filtering and embedding based nearest neighbor fil-
179 tering shows convergence improvement [14, 26].

180 2.3. Latent Embedding Constraints

181 A notable latent embedding constraint that related to our ap-
182 proach is Evidence Lower Bound (ELBO). ELBO approxi-
183 mates a latent sample with a variational distribution and
184 constrains the KL-divergence between the variational dis-
185 tribution and a target shape which is typically a multivariate
186 Gaussian [15]. The main drawback of this approach is
187 that the true KL-divergence is intractable to calculate. As
188 such the posterior must take on a simplified form. Most
189 implementations use a diagonal posterior which only penalizes
190 simple differences in shape such as mean and standard
191 deviation. The method of Normalizing Flows provides an
192 iterative framework based on change of variables to con-
193 struct arbitrarily complex posteriors [8, 16, 32]. Our MoM
194 constraint can also penalize complex differences in shape
195 by constraining 2nd, 3rd, and 4th order hyper-covariance
196 matrices, although we do so by comparing the moments di-
197 rectly. This greatly simplifies implementation as we do not
198 need to explicitly construct a posterior distribution.

199 Another notable embedding constraint is the Maximum
200 Mean Discrepancy (MMD), also known as the two-sample
201 test [12]. MMD was used in the Generative Moment Match-
202 ing Network [21] and has since been used extensively for
203 the problem of domain adaptation [38, 39], in order to con-
204 strain the latent projections of the source and target distri-
205 butions to follow the same distribution. MMD is a moment-
206 matching constraint based on the kernel trick, and can there-
207 fore constrain any difference in shape between two distribu-
208 tions including very high order moments. Typically, MMD
209 is only be used to constrain one sample to another sample,
210 although numerical methods can extend this technique to
211 include sample-to-distribution comparisons [49].

212 3. Methodology

213 In this section, we explore our proposed replacement final
214 activation layers and our embedding space constraints. Fix-
215 Match [35] is a simple, well performing SSL algorithm. As
216 such, it serves as a good comparison point for exploring
217 the effect of our contributions. Our methodology is based
218 upon the published FixMatch [35] algorithm, with identical
219 hyper-parameters unless otherwise stated. We extend Fix-
220 Match with a few minor training algorithm modifications
221 explored in the Hyperparameters Section 4.1. While we
222 keep the FixMatch valid pseudo-label selection, our embed-
223 ding constraints operate on all data, labeled and unlabeled.

224 Both the linear layer replacements and the embedding
225 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, 10], but we extend the core idea with a novel differentiable model with learned cluster centroids and MoM based constraints. The MoM constraints do not impose any assumptions outside of applying l2 penalties as described in Section 3.2.

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,

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

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,

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

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

As there is one cluster per label category, the joint probability for sample i with label assignment k , $p(Y_{i,k}, X_i)$ is 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 $\hat{Y}_{i,k} = p(Y_{i,k}|X_i)$ 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 [29] module. It has two parameters updated by backprop. (1) the explicit cluster centers, a matrix $num_classes \times embedding_dim$ initialized randomly, and (2) the diagonal elements of the *Sigma* 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 [29] 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) [30] 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. This weighting scheme supports using 0 to 4th order MoM constraints seamlessly and is not a hyper-parameter we expect to require tuning.

4. Experiments

We evaluate both AAGMM and KMeans linear layer replacements and the embedding space constraints using our modified FixMatch[35] on the common SSL benchmarks CIFAR-10 [18] 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 [35] 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

CIFAR-10 Mean 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	77.14 \pm 9.09				
	8	88.40 \pm 3.54				
AAGMM	128	91.30 \pm 2.89	89.23 \pm 2.50	90.22 \pm 3.42		
	8	88.40 \pm 2.34	82.39 \pm 9.96	87.97 \pm 3.35	82.51 \pm 9.42	82.57 \pm 6.94
KMeans	128	81.13 \pm 2.19	89.74 \pm 2.62	89.89 \pm 2.79		
	8	78.28 \pm 9.87	73.27 \pm 12.01	75.80 \pm 10.69	80.96 \pm 5.94	80.01 \pm 7.53
Last Layer	Emb Dim	250 Labels (5 trials)				
Embedding Constraint		None	1st Order	2nd Order	3rd Order	4th Order
Linear (i.e. FullyConnected)	128	94.06 \pm 0.87				
	8	93.69 \pm 0.50				
AAGMM	128	94.30 \pm 0.51	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	94.10 \pm 0.51
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 0.59	93.64 \pm 0.99	94.10 \pm 0.55	94.09 \pm 0.59

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.

391 with random label sampling, vs how well it can potentially
 392 do with a more representative subset of labeled data.

4.1. Hyper-Parameters

394 All CIFAR-10 models were trained with the standard
 395 benchmark WideResNet28-2 architecture. This work lever-
 396 aged the published FixMatch [35] hyper-parameters; us-
 397 ing SGD with Nesterov momentum, $\lambda_u = 1$, $\beta = 0.9$,
 398 $\tau = 0.95$, $\mu = 7$, $B = 64$, and epoch size = 1024
 399 batches regardless of the number of images in the labeled
 400 dataset. Model weights were updated as the moving
 401 average of the training weights with an exponential moving
 402 average (EMA) decay of 0.999. The training algorithm was
 403 modified from stock FixMatch to include an early-stopping
 404 condition when the model has not improved for 50 epochs.
 405 The $learning_rate(\eta) = 0.01$. Replacing the fixed number
 406 of training steps with an early stopping criteria prevents the
 407 use of a cosine decay schedule. Therefore, it was replaced
 408 with a plateau learning rate scheduler which multiplies the
 409 learning rate by 0.2 every time the early stopping criteria is
 410 met (before being reset) for a max of 2 reductions. To re-
 411 duce the training algorithm dependence on specific learning
 412 rate values, a cyclic learning rate scheduler was employed
 413 to vary the learning rate by a factor of ± 2.0 within each
 414 epoch. Additionally, due to the higher training instability of
 415 the AAGMM layers compared to a linear layer, if the train-
 416 ing loss is greater than 1.0, the gradient norm was clipped to
 417 1.0. Despite specific attention to computing the AAGMM
 418 and embedding constraints in a numerically stable manner,

419 they are still less stable during backprop than a simple linear
 420 layer.

421 This work includes an exploration of how various la-
 422 tent embedding dimensionalities affects the generative lin-
 423 ear layer replacement. As such, the model architecture
 424 was modified with a single additional linear layer before
 425 the output to project the baseline model embedding dimen-
 426 sion (128 for WideResNes28-2) down to a reduced 8 dimen-
 427 sional space. The AAGMM and KMeans replacement lay-
 428 ers with and without this reduced embedding space, were
 429 evaluated to determine whether the generative capabilities
 430 improve when not fighting the curse of dimensionality. Re-
 431 sults listed with an embedding dimensionality of 128 do not
 432 include the additional linear layer which reduces the latent
 433 dimensionality. Therefore, results with 128D embedding
 434 represents an unmodified network architecture.

435 Due to exponential GPU memory requirements with
 436 each successive MoM moment, only the 8D embedding can
 437 operate with higher order MoM embedding constraints. Re-
 438 sults for any given order of embedding constraint include all
 439 lower constraints.

4.2. CIFAR-10

440 The CIFAR-10 SSL benchmark was used to explore the
 441 full configuration space of our method. While both 40 and
 442 250 label counts were used, the 250 label case SOTA is
 443 close to fully supervised accuracy. We include 250 per-
 444 formance to document our result is approximately equiv-
 445 alent to SOTA. The 40 label case provides a far more chal-
 446 lenge.

CIFAR-10 Max Test Accuracy

Last Layer	Emb Dim	40 Labels (5 trials)			
Embedding Constraint		None	1st Order	2nd Order	3rd Order
Linear (i.e. FullyConnected)	128	91.01			
	8	92.08			
AAGMM	128	94.64	91.33	92.74	
	8	90.40	93.25	92.11	91.95
KMeans	128	84.08	91.62	92.22	
	8	93.21	91.22	89.35	85.96

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.

lenging 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 [35] 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% [41].

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 [51] 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 (a & b) showcases degenerate cluster centers. The Figure 3 (c) 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 (d). 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.

Table 4 puts these results in context with the current SSL SOTA for CIFAR-10 at 40 labels and demonstrates that this methodology still requires improvement before it is competitive with the latest methods.

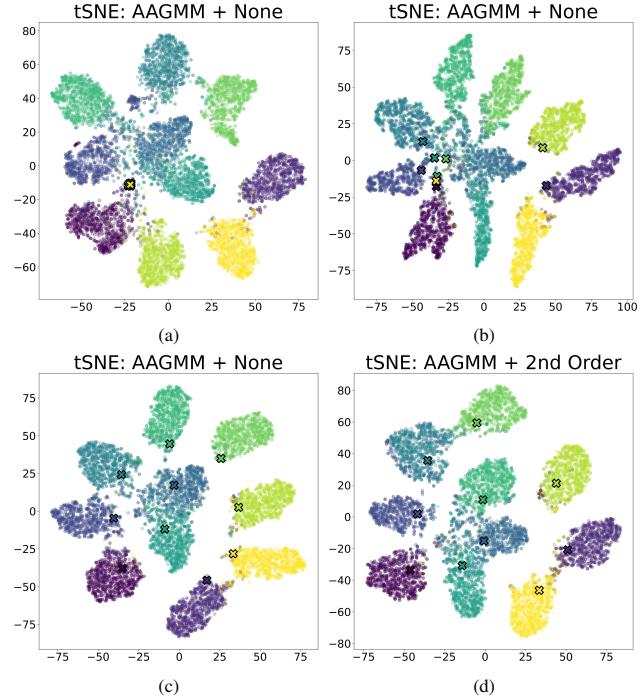


Figure 3. t-SNE plot of fully trained and reasonably accuracy AAGMM model's latent embedding space (with the learned cluster centers marked with X's). Depending on the run, the AAGMM cluster centers will be degenerate (top left), non-degenerate but still mis-aligned with the clusters (top right), acceptably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

5. Discussion

The proposed MoM embedding constraint has at least one significant downside, it requires exponentially increasing amounts of GPU memory for each successive moment penalty included. This limits the current practicality of these MoM constraints. 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 us-

Method	CIFAR-10	
Label Count	40	250
FixMatch[35]	13.81 \pm 3.37	5.07 \pm 0.65
FlexMatch[48]	4.97 \pm 0.06	4.98 \pm 0.09
FreeMatch[41]	4.90 \pm 0.29	4.98 \pm 0.09
SimMatchV2[51]	4.90 \pm 0.04	5.04 \pm 0.09
Ours (AAGMM+None)	8.77 \pm 2.89	5.91 \pm 0.34
Ours (KMeans+2ndOrder)	10.11 \pm 2.79	6.84 \pm 1.25

Table 4. Error rate % for CIFAR-10 SSL benchmark comparing to state of the art results. Results for previously published methods are drawn from USB [40] except for FreeMatch[41] and SimMatchV2[51] publications.

ability.

Semi-supervised learning is highly sensitive to both which samples are selected for the labeled population [35] and the stochasticity of the training process itself. Given identical starting random seeds, and identical labeled samples, training stochasticity will quickly cause models to diverge, resulting in vastly different final results. Anecdotally its appears worse in semi-supervised methods than fully supervised models. To characterize this variance, and hence how much one can trust the error bars for Table 2 and 3, we took a few final layer configurations and ran them $N = 5$ times with the same seed. Table 5 showcases the run-to-run variances for models that started out identical. Interestingly enough, the AAGMM models converge with much lower variance than the KMeans models. This contrasts with the KMeans layer being significantly simpler than AAGMM, both mathematically and implementation-wise.

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	86.5
Run 2	85.6	91.2	67.9	85.4
Run 3	84.8	77.7	75.8	86.0
Run 4	85.3	87.6	69.5	87.6
Run 5	85.2	83.7	78.9	85.7

Table 5. Test accuracy for independent training runs with the same random seed for various final layer configurations. All runs use the 128D (baseline) model latent embedding dimensionality. All runs use the same labeled samples.

Future work in this area should explore both accuracy improvements as well as implementation optimization to ensure the proposed novel final layers are not prohibitively memory expensive. 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 novel 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 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. Remixmatch: 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. 4
- [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. 4
- [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] Anthony Caterini, Rob Cornish, Dino Sejdinovic, and Arnaud Doucet. Variational inference with continuously-indexed normalizing flows. In *Uncertainty in Artificial Intelligence*, pages 44–53. PMLR, 2021. 3
- [9] 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
- [10] 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. 4
- [11] Enrico Fini, Pietro Astolfi, Kartek 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*

- 573 on *Computer Vision and Pattern Recognition*, pages 3187–
574 3197, 2023. 3
- 575 [12] Arthur Gretton, Kenji Fukumizu, Choon Teo, Le Song, Bern-
576 hard Schölkopf, and Alex Smola. A kernel statistical test of
577 independence. *Advances in neural information processing systems*, 20, 2007. 3
- 578 [13] Mohamed Farouk Abdel Hady and Friedhelm Schwenker.
579 Semi-supervised learning. *Handbook on Neural Information
580 Processing*, pages 215–239, 2013. 2
- 581 [14] Jiwon Kim, Youngjo Min, Daehwan Kim, Gyuseong Lee,
582 Junyoung Seo, Kwangrok Ryoo, and Seungryong Kim. Comatch:
583 Semi-supervised learning with confidence-guided
584 consistency regularization. In *European Conference on
585 Computer Vision*, pages 674–690. Springer, 2022. 2, 3
- 586 [15] Diederik P Kingma and Max Welling. Auto-encoding variational
587 bayes. *arXiv preprint arXiv:1312.6114*, 2013. 2, 3
- 588 [16] Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen,
589 Ilya Sutskever, and Max Welling. Improved variational in-
590 ference with inverse autoregressive flow. *Advances in neural
591 information processing systems*, 29, 2016. 3
- 592 [17] Diederik P Kingma, Max Welling, et al. An introduction to
593 variational autoencoders. *Foundations and Trends® in Ma-
594 chine Learning*, 12(4):307–392, 2019. 2
- 595 [18] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple
596 layers of features from tiny images. Technical report, 2009.
597 2, 5
- 598 [19] Doyup Lee, Sungwoong Kim, Ildoo Kim, Yeongjae Cheon,
599 Minsu Cho, and Wook-Shin Han. Contrastive regulariza-
600 tion for semi-supervised learning. In *Proceedings of the
601 IEEE/CVF Conference on Computer Vision and Pattern
602 Recognition*, pages 3911–3920, 2022. 2
- 603 [20] Junnan Li, Caiming Xiong, and Steven CH Hoi. Comatch:
604 Semi-supervised learning with contrastive graph regulariza-
605 tion. In *Proceedings of the IEEE/CVF International Confer-
606 ence on Computer Vision*, pages 9475–9484, 2021. 2, 3
- 607 [21] Yujia Li, Kevin Swersky, and Rich Zemel. Generative mo-
608 ment matching networks. In *International conference on ma-
609 chine learning*, pages 1718–1727. PMLR, 2015. 3
- 610 [22] Yang Li, Quan Pan, Suhang Wang, Haiyun Peng, Tao Yang,
611 and Erik Cambria. Disentangled variational auto-encoder for
612 semi-supervised learning. *Information Sciences*, 482:73–85,
613 2019. 2
- 614 [23] Yu-Feng Li and De-Ming Liang. Safe semi-supervised learn-
615 ing: a brief introduction. *Frontiers of Computer Science*, 13:
616 669–676, 2019. 2
- 617 [24] Ioannis E Livieris, Konstantina Drakopoulou, Vassilis T
618 Tampakas, Tassos A Mikropoulos, and Panagiotis Pintelas.
619 Predicting secondary school students’ performance utilizing
620 a semi-supervised learning approach. *Journal of educational
621 computing research*, 57(2):448–470, 2019. 3
- 622 [25] David McClosky, Eugene Charniak, and Mark Johnson.
623 Reranking and self-training for parser adaptation. In *Pro-
624 ceedings of the 21st International Conference on Compu-
625 tational Linguistics and 44th Annual Meeting of the Associa-
626 tion for Computational Linguistics*, pages 337–344, 2006. 3
- 627 [26] Sumeet Menon. *Semi-Supervised Expectation Maximiza-
628 tion with Contrastive Outlier Removal*. PhD thesis, 2022.
629 AAI29167191. 3
- 630 [27] Sumeet Menon, David Chapman, Phuong Nguyen, Ye-
631 lena Yesha, Michael Morris, and Babak Saboury. Deep
632 expectation-maximization for semi-supervised lung cancer
633 screening. 2019. 3
- 634 [28] Aamir Mustafa and Rafał K Mantiuk. Transformation
635 consistency regularization—a semi-supervised paradigm for
636 image-to-image translation. In *European Conference on
637 Computer Vision*, pages 599–615. Springer, 2020. 3
- 638 [29] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer,
639 James Bradbury, Gregory Chanan, Trevor Killeen, Zeming
640 Lin, Natalia Gimelshein, Luca Antiga, Alban Desma-
641 son, Andreas Kopf, Edward Yang, Zachary DeVito, Mar-
642 tin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit
643 Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch:
644 An imperative style, high-performance deep learning library.
645 In *Advances in Neural Information Processing Systems* 32,
646 pages 8024–8035. Curran Associates, Inc., 2019. 4
- 647 [30] Karl Pearson. Method of moments and method of maximum
648 likelihood. *Biometrika*, 28(1/2):34–59, 1936. 2, 4
- 649 [31] V Jothi Prakash and Dr LM Nithya. A survey on
650 semi-supervised learning techniques. *arXiv preprint
651 arXiv:1402.4645*, 2014. 3
- 652 [32] Danilo Rezende and Shakir Mohamed. Variational inference
653 with normalizing flows. In *International conference on ma-
654 chine learning*, pages 1530–1538. PMLR, 2015. 3
- 655 [33] Mamshad Nayeem Rizve, Kevin Duarte, Yogesh S Rawat,
656 and Mubarak Shah. In defense of pseudo-labeling: An
657 uncertainty-aware pseudo-label selection framework for
658 semi-supervised learning. In *International Conference on
659 Learning Representations*, 2021. 3
- 660 [34] Chuck Rosenberg, Martial Hebert, and Henry Schneiderman.
661 Semi-supervised self-training of object detection models. In
662 *2005 Seventh IEEE Workshops on Applications of Computer
663 Vision (WACV/MOTION'05) - Volume 1*, pages 29–36, 2005.
664 3
- 665 [35] Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao
666 Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk,
667 Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying
668 semi-supervised learning with consistency and confidence.
669 *Advances in neural information processing systems*, 33:596–
670 608, 2020. 2, 3, 5, 6, 7, 8
- 671 [36] Antti Tarvainen and Harri Valpola. Mean teachers are better
672 role models: Weight-averaged consistency targets improve
673 semi-supervised deep learning results. *Advances in neural
674 information processing systems*, 30, 2017. 3
- 675 [37] Laurens Van der Maaten and Geoffrey Hinton. Visualizing
676 data using t-sne. *Journal of machine learning research*, 9
677 (11), 2008. 2
- 678 [38] Mei Wang and Weihong Deng. Deep visual domain adap-
679 tation: A survey. *Neurocomputing*, 312:135–153, 2018. 3
- 680 [39] Wei Wang, Haojie Li, Zhengming Ding, Feiping Nie, Jun-
681 yang Chen, Xiao Dong, and Zhihui Wang. Rethinking maxi-
682 mum mean discrepancy for visual domain adaptation. *IEEE
683 Transactions on Neural Networks and Learning Systems*, 34
684 (1):264–277, 2023. 3
- 685 [40] Yidong Wang, Hao Chen, Yue Fan, Wang Sun, Ran Tao,
686 Wenxin Hou, Renjie Wang, Linyi Yang, Zhi Zhou, Lan-Zhe
687 9

- 688 Guo, et al. Usb: A unified semi-supervised learning bench-
689 mark for classification. *Advances in Neural Information Pro-*
690 *cessing Systems*, 35:3938–3961, 2022. 8
- 691 [41] Yidong Wang, Hao Chen, Qiang Heng, Wenxin Hou, Yue
692 Fan, Zhen Wu, Jindong Wang, Marios Savvides, Takahiro
693 Shinozaki, Bhiksha Raj, Bernt Schiele, and Xing Xie.
694 Freematch: Self-adaptive thresholding for semi-supervised
695 learning. In *The Eleventh International Conference on*
696 *Learning Representations*, 2023. 7, 8
- 697 [42] Hongxin Wei, RENCHUNZI Xie, Hao Cheng, Lei Feng, Bo An,
698 and Yixuan Li. Mitigating neural network overconfidence
699 with logit normalization. In *International Conference on*
700 *Machine Learning*, pages 23631–23644. PMLR, 2022. 1
- 701 [43] Qizhe Xie, Zihang Dai, Eduard Hovy, Thang Luong, and
702 Quoc Le. Unsupervised data augmentation for consistency
703 training. *Advances in neural information processing systems*,
704 33:6256–6268, 2020. 3
- 705 [44] Chen Xing, Sercan Arik, Zizhao Zhang, and Tomas Pfister.
706 Distance-based learning from errors for confidence calibra-
707 tion. In *International Conference on Learning Representa-*
708 *tions*, 2020. 3
- 709 [45] Fan Yang, Kai Wu, Shuyi Zhang, Guannan Jiang, Yong Liu,
710 Feng Zheng, Wei Zhang, Chengjie Wang, and Long Zeng.
711 Class-aware contrastive semi-supervised learning. In *Pro-*
712
713 *and Pattern Recognition*, pages 14421–14430, 2022. 2, 3
- 714 [46] David Yarowsky. Unsupervised word sense disambiguation
715 rivaling supervised methods. In *33rd annual meeting of the*
716 *association for computational linguistics*, pages 189–196,
717 1995. 3
- 718 [47] Xiaohua Zhai, Avital Oliver, Alexander Kolesnikov, and Lu-
719 cas Beyer. S4l: Self-supervised semi-supervised learning. In
720 *Proceedings of the IEEE/CVF International Conference on*
721 *Computer Vision*, pages 1476–1485, 2019. 3
- 722 [48] Bowen Zhang, Yidong Wang, Wenxin Hou, Hao Wu, Jin-
723 dong Wang, Manabu Okumura, and Takahiro Shinozaki.
724 Flexmatch: Boosting semi-supervised learning with curricu-
725 lum pseudo labeling. *Advances in Neural Information Pro-*
726 *cessing Systems*, 34:18408–18419, 2021. 2, 3, 8
- 727 [49] Shengjia Zhao, Jiaming Song, and Stefano Ermon. Info-
728 vae: Balancing learning and inference in variational autoen-
729 coders. In *Proceedings of the aaai conference on artificial*
730 *intelligence*, pages 5885–5892, 2019. 3
- 731 [50] Mingkai Zheng, Shan You, Lang Huang, Fei Wang, Chen
732 Qian, and Chang Xu. Simmatch: Semi-supervised learning
733 with similarity matching. In *Proceedings of the IEEE/CVF*
734 *Conference on Computer Vision and Pattern Recognition*,
735 pages 14471–14481, 2022. 3
- 736 [51] Mingkai Zheng, Shan You, Lang Huang, Chen Luo, Fei
737 Wang, Chen Qian, and Chang Xu. Simmatchv2: Semi-
738 supervised learning with graph consistency. pages 16432–
739 16442, 2023. 3, 7, 8
- 740 [52] Xiaojin Zhu and Andrew B Goldberg. *Introduction to semi-*
741 *supervised learning*. Springer Nature, 2022. 2