

# 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 apply this method to the domain of semi-supervised image classification by extending FixMatch with our technique. We find our MoM constraint with the AAGMM layer is able to match 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 rep-

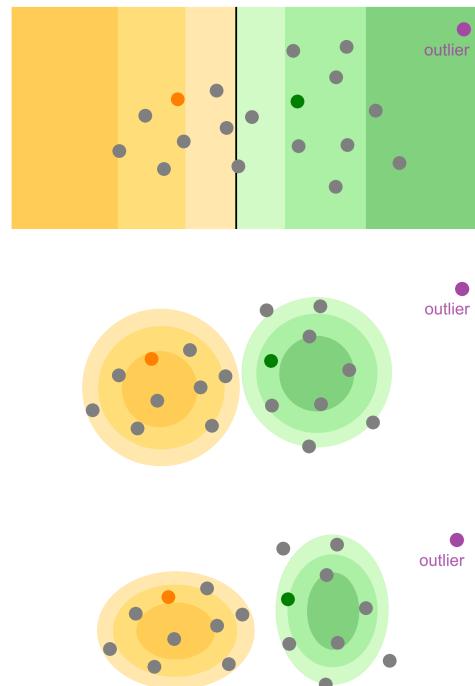


Figure 1. Schematic of the outlier problem, and how generative modeling of the joint probability can improve the situation. Prediction with (top) fully-supervised softmax (middle) semi-supervised KMeans and (bottom) semi-supervised AAGMM.

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

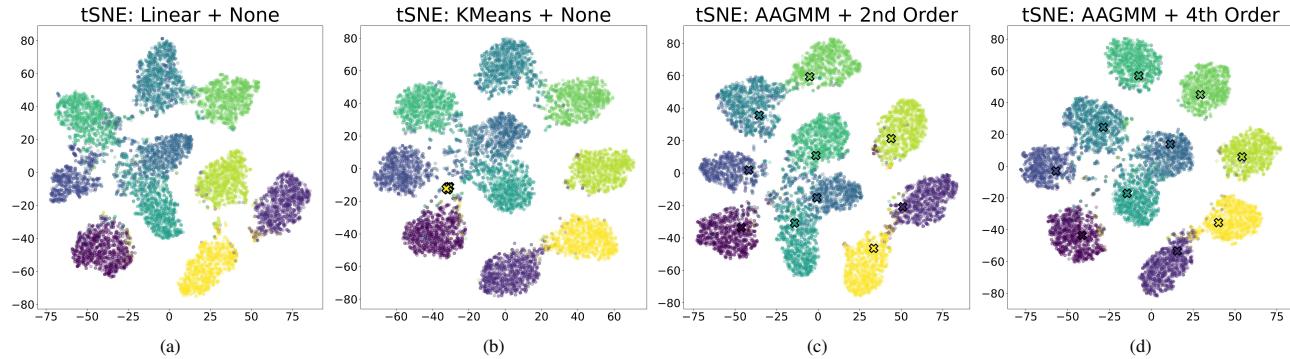


Figure 2. t-SNE[37] plot of the latent embedding space for various final layers with different MoM embedding constraints.

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 joint distribution of the latent space. 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.

We apply this work to the domain of Semi-Supervised Learning (SSL), because over-confident label predictions can cause confounding issues with pseudo-labeling methods [1]. 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. 088
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. 089
3. A visualization of the latent embedding space, showing how it responds to the incorporation of generative final activation layers with and without MoM based embedding constraints. 090

We apply this methodology to the task of semi-supervised CIFAR-10 image classification with 40 and 250 training labels [18]. The embedding constraint penalties are applied to all unlabeled data and not just the valid pseudo-labels. As such our method fits the latent joint distribution across all of the unlabeled data points, an improvement on baseline pseudo-labeling methods (like FixMatch [35]) which only fit the conditional distribution to the high confidence pseudo-labels while removing low confidence pseudo-labels.

## 2. Related Work

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

## 121 2.1. Pseudo-Labeling

122 Self-training was among the initial approaches employed  
123 in the context of SSL to annotate unlabeled images. This  
124 technique involves the initial training of a classifier with  
125 a limited set of labeled samples and incorporates pseudo-  
126 labels into the gradient descent process, exceeding a pre-  
127 defined threshold [9, 24, 25, 27, 34, 46, 47]. A closely related  
128 method to self-training is co-training, where a given dataset  
129 is represented as two distinct feature sets [4]. These inde-  
130 pendent sample sets are subsequently trained separately us-  
131 ing two distinct models, and the sample predictions surpass-  
132 ing predetermined thresholds are utilized in the final model  
133 training process [4, 31]. A notably advanced approach to  
134 pseudo-labeling is the Mean Teacher algorithm [36], which  
135 leverages exponential moving averages of model parame-  
136 ters to acquire a notably more stable target prediction. This  
137 refinement has a substantial impact on enhancing the con-  
138 vergence of the algorithm.

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

## 154 2.2. Consistency Regularization

155 Consistency regularization operates on the premise that  
156 when augmenting an unlabeled sample, its label should  
157 remain consistent. This approach implicitly enforces a  
158 smoothness assumption, promoting coherence between un-  
159 labeled samples and their basic augmentations [43]. In  
160 other words, the model should be able to predict the un-  
161 labeled sample  $x$  exactly the same way it predicts the class  
162 for  $\text{Augmented}(x)$  [2, 3, 28, 35]. In addition to evaluat-  
163 ing image-wise augmentations, recent research has demon-  
164 strated that incorporating class-wise and instance-based  
165 consistencies yields superior performance outcomes [20,  
166 50]. Similarly, using consistencies between the predictions  
167 and low-dimensional embeddings from the unlabeled im-  
168 age strong and weak augmentations in a graph based setup  
169 demonstrates improvement over class-wise and instance-  
170 based consistencies [51]. Finally, pseudo-labeling filtering  
171 based on consistence between strongly augmented views,  
172 gaussian filtering and embedding based nearest neighbor fil-

173 tering shows convergence improvement [14, 26].

## 174 2.3. Latent Embedding Constraints

175 A notable latent embedding constraint that is related yet  
176 substantially different from our approach is the Evidence  
177 Lower Bound (ELBO). ELBO approximates a latent sam-  
178 ple with a variational distribution and constrains the KL-  
179 divergence between the variational distribution and a target  
180 shape which is typically a multivariate standard normal dis-  
181 tribution [15]. The main drawback of this approach is that  
182 the true KL-divergence is intractable to calculate. As such  
183 the posterior must take on a simplified form. Most prac-  
184 tical implementations use a diagonal posterior which can  
185 only penalize simple differences in shape such as mean and  
186 standard deviation. Arbitrarily complex posteriors are nev-  
187 ertheless possible using the method of Normalizing flows,  
188 which provides an iterative framework based on change of  
189 variables although this method is quite involved [8, 16, 32].  
190 Our MoM constraint is relatively simple but can also penali-  
191 ze complex differences in shape by constraining 2nd, 3rd,  
192 and 4th order hyper-covariance matrices, although we do so  
193 by comparing the moments directly. This greatly simplifies  
194 implementation as we do not need to explicitly construct a  
195 posterior distribution.

196 Another notable embedding constraint is the Maximum  
197 Mean Discrepancy (MMD), also known as the two-sample  
198 test [12]. MMD was used in the Generative Moment Match-  
199 ing Network [21] and has since been used extensively for  
200 the problem of domain adaptation [38, 39], in order to con-  
201 strain the latent projections of the source and target distri-  
202 butions to follow the same distribution. MMD is a moment-  
203 matching constraint based on the kernel trick, and can there-  
204 fore constrain any difference in shape between two sam-  
205 ples including very high order moments. Due to the ker-  
206 nel trick requiring proper inner products, MMD can only  
207 be used to constrain one sample to another sample. It can-  
208 not directly constrain sample statistics to population statis-  
209 tics, although it is possible to approximate populations nu-  
210 merically via monte-carlo sampling [49]. Like MMD, our  
211 method is based on MoM, but it does not involve the kernel  
212 trick, and instead penalizes polynomial moments explicitly  
213 thereby enabling the sample embedding to be constrained  
214 to an exact target distribution.

## 215 3. Methodology

216 In this section, we explore our proposed replacement final  
217 activation layers and our embedding space constraints. Our  
218 methodology is based upon the published FixMatch [35] al-  
219 gorithm, with identical hyper-parameters unless otherwise  
220 stated. We extend FixMatch with a few minor training algo-  
221 rithm modifications explored in the Hyperparameters, Sec-  
222 tion 4.1. FixMatch [35] is a simple, well performing SSL

223 algorithm. As such, it serves as a good comparison point  
 224 for exploring the effect of our contributions.

225 Both the linear layer replacements and the embedding  
 226 constraints explored herein represent increasing levels of  
 227 prescription about how the final latent embedding space  
 228 should be arranged compared to a traditional linear layer.  
 229 The idea of leveraging clusters in embedding space is not  
 230 new [5, 6, 10], but we extend the core idea with a novel  
 231 differentiable model with learned cluster centroids and MoM  
 232 based constraints. The MoM constraints do not impose any  
 233 assumptions outside of applying l2 penalties as described in  
 234 Section 3.2.

### 235 3.1. Alternate Final Layers

236 A limitation of traditional final activation layers such as linear+softmax  
 237 is that they are fully discriminative; i.e. they  
 238 estimate the conditional probability  $p(Y|X)$ , but do not attempt  
 239 to model the prior distribution  $p(X)$  or the joint probabilities  
 240  $p(Y, X)$ . To overcome this limitation, we present  
 241 two generative final activation layers (a) the Axis Aligned  
 242 GMM (AAGMM) layer, and (b) an equal variance version  
 243 of AAGMM that we henceforth call the KMeans activation  
 244 layer due to the similarity of the objective function with a  
 245 gradient based KMeans.

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

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

#### 255 3.1.1 Axis Aligned Gaussian Mixture Model Layer

256 The AAGMM layer defines a set of  $K$  trainable clusters,  
 257 one cluster per label category. Each cluster  $k = 1 \dots K$   
 258 has a cluster center  $\mu_k$  and cluster covariance  $\Sigma_k$ . The prior  
 259 probability of any given sample  $X_i$  is defined by the mixture  
 260 of cluster probability densities over the  $D$ -dimensional latent  
 261 representation  $Z_i$  as follows,

$$265 p(X_i) = \sum_{k=1}^K \mathcal{N}(Z_i, \mu_k, \Sigma_k) \quad (1)$$

266 where  $Z_i = F(X_i, \theta_F)$

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

268 is axis aligned because  $\Sigma_k$  is a diagonal matrix, as such  
 269 the axis-aligned multivariate normal pdf simplifies to the  
 270 marginal product of Gaussians along each of the  $D$  axes as  
 271 follows,

$$272 \mathcal{N}(X_i, \mu_k, \Sigma_k) = \prod_{d=1}^D \frac{1}{\sigma_{k,d}\sqrt{2\pi}} \exp\left(\frac{(Z_{i,d} - \mu_{k,d})^2}{\sigma_{k,d}^2}\right) \quad (2)$$

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

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

$$276 p(Y_{i,k}, X_i) = \mathcal{N}(Z_i, \mu_k, \Sigma_k) \quad (3)$$

277 By Bayesian identity, the conditional probability  $\hat{Y}_{i,k} =$   
 278  $p(Y_{i,k}|X_i)$  can therefore be inferred from eq 1 and 3 as follows,

$$279 \hat{Y}_{i,k} = p(Y_{i,k}|X_i) = \frac{p(Y_{i,k}, X_i)}{p(X_i)} \quad (4)$$

280 The AAGMM layer is implemented as a normal PyTorch  
 281 [29] module. It has two parameters updated by backprop.  
 282 (1) the explicit cluster centers, a matrix  $num\_classes \times$   
 283  $embedding\_dim$  initialized randomly, and (2) the diagonal  
 284 elements of the  $\Sigma_k$  matrix, randomly initialized in the  
 285 range  $[0.9, 1.1]$ , which contains the diagonal elements of the  
 286 GMM Sigma matrix for each cluster.

#### 287 3.1.2 KMeans Layer

288 We also implement a KMeans final layer which is a more  
 289 restrictive form of the AAGMM layer. The KMeans layer  
 290 is additionally constrained such that the gaussian covariance  
 291 matrix  $\Sigma_k$  for each cluster center  $k$  is the  $[D \times D]$  identity  
 292 matrix. This constraint yields spherical cluster centers; similar  
 293 to how the traditional KMeans algorithm also assumes  
 294 spherical clusters.

295 The KMeans layer is also implemented as a normal PyTorch  
 296 [29] module. The explicit cluster centers is a learned  
 297 parameter updated by backprop. See the published codebase  
 298 for implementation details about the AAGMM and  
 299 KMeans layers.

### 300 3.2. Method of Moments Embedding Constraints

301 We introduce and evaluate a series of embedding constraints  
 302 based on the Method of Moments (MoM) [30]. For each  
 303 sample  $i$ , the joint  $p(Y_{i,k}, X_i)$  is calculated as in equation  
 304 3 and then used to infer the prior  $p(X_i)$  and the conditional  
 305  $p(Y_{i,k}|X_i)$ . As usual, the conditional probability is trained  
 306 using cross entropy loss. When embedding constraints are  
 307 omitted, it is possible for the model to learn an accurate

309 decision boundary for the conditional probability without  
 310 modeling the latent joint distribution. MoM solves these  
 311 problems and is an appropriate strategy for semi-parametric  
 312 models.

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

$$319 \forall p = 1 \dots P \quad \frac{1}{n} \sum_{i=1}^n g_p(z_i) \approx E(g_p(Z)) \quad (5)$$

320 We can therefore constrain the latent representation of  
 321 our model to approximate a multivariate standard normal  
 322 distribution. In the univariate standard normal case, the  $p^{th}$   
 323 order centralized moment constraint is the following.

$$324 E[(Z - \mu)^p] = \begin{cases} 0 & \text{if } p \text{ is odd} \\ \sigma^p(p-1)!! & \text{if } p \text{ is even} \end{cases} \quad (6)$$

325 Where '!!' represents the double factorial operator. By  
 326 this formula, the univariate unit gaussian has mean 0, stan-  
 327 dard deviation 1, skew 0, and kurtosis 3.

328 The multivariate standard normal distribution is the  
 329 marginal product of the univariate standard normal distri-  
 330 butions. As such, if we redefine  $Z$ ,  $\mu$ , and  $p$  to be all  $D$   
 331 dimensional, then the centralized marginal product moment  
 332 can be defined as follows,

$$333 E[g_p(Z - \mu)] = E\left[\prod_{d=1}^D (Z_d - \mu_d)^{p_d}\right] \quad (7)$$

334 Due to independence of the axes, this multivariate popu-  
 335 lation moment can be represented as a product of univariate  
 336 moments of the individual standard normal distributions as  
 337 follows,

$$338 E\left[\prod_{d=1}^D (Z_d - \mu_d)^{p_d}\right] = \prod_{d=1}^D E[(Z_d - \mu_d)^{p_d}] \quad (8)$$

339 The error (loss) term associated with the embedding con-  
 340 straint for any moment  $p$  is equal to the L2 difference be-  
 341 tween the sample and population statistics as follows,

$$342 \varepsilon_p = \left( \frac{1}{n} \sum_{i=1}^n g_p(z_i) - E(g_p(Z)) \right)^2 \quad (9)$$

343 Some moments are more important than others, and must  
 344 be weighted more heavily. First order moments are simply  
 345 the sample mean, and should be given the greatest weight as

346 an embedding constraint. The second order moments form a  
 347 sample covariance matrix, which ideally should be equal to  
 348 the identity matrix, but the diagonal terms should be given  
 349 greater weight than the off-diagonal terms. This is because,  
 350 in a  $D \times D$  covariance matrix, there are  $D(D-1)$  off diag-  
 351 onal terms, but only  $D$ , diagonal terms. The  $p^{th}$  order sample  
 352 moments form a  $p-1$  dimensional hyper-covariance matrix,  
 353 with terms residing on the intersection of anywhere between  
 354 0 and  $p-1$  hyper-diagonals. To prevent over-representation  
 355 of off-diagonal terms and encourage representation of on-  
 356 diagonal terms, the loss function we use for any given mo-  
 357 ment term is inversely proportional to the number moment  
 358 terms that share the same number of hyper-diagonals. This  
 359 heuristic weighting scheme ensures that the overall contribu-  
 360 tion of each moment order is not overly influenced by the  
 361 off-diagonal terms, and that the error weighting is therefore  
 362 diagonally dominant. This weighting scheme supports us-  
 363 ing 0 to 4th order MoM constraints seamlessly and is not a  
 364 hyper-parameter we expect to require tuning.

## 4. Experiments

365 We evaluate both AAGMM and KMeans linear layer re-  
 366 placements and the embedding space constraints using our  
 367 modified FixMatch[35] on the common SSL benchmarks  
 368 CIFAR-10 [18] at 40 and 250 labels (4 and 25 labels per  
 369 class). We randomly selected 5 seed a priori for evalua-  
 370 tion. For each algorithm configuration tested one model  
 371 was trained per seed. During each run, the required num-  
 372 ber of labeled samples are drawn without replacement from  
 373 the training population of the dataset in a deterministic man-  
 374 ner (reproducible with the same seed). All data not in this  
 375 labeled subset is used as unlabeled data (i.e. the labels are  
 376 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

377 Table 1. Test accuracy for showing the run-to-run variance de-  
 378 pending on the quality of the 40 labels selected from the full pop-  
 379 ulation.

380 As prior work [35] has noted, the resulting model qual-  
 381 ity is highly variable when only 4 samples are selected per  
 382 class, as the quality and usefulness of the specific 4 samples  
 383 can vary drastically. Table 1 shows final test accuracy for  
 384 the 5 AAGMM model runs with no embedding constraints,  
 385 with the accuracy varying from 86% to 94%. Due to the  
 386 potential for significant variance in the final model test ac-  
 387 curacy, it can be informative to compare mean performance  
 388 with max performance over the  $N = 5$  runs. This explores  
 389 both how well a method can be expected to do on average  
 390 with random label sampling, vs how well it can potentially

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	<b>91.30</b> $\pm$ 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	<b>89.89</b> $\pm$ 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.

389 do with a more representative subset of labeled data.

## 4.1. Hyper-Parameters

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

layer.

This work includes an exploration of how various latent embedding dimensionalities affects the generative linear layer replacement. As such, the model architecture was modified with a single additional linear layer before the output to project the baseline model embedding dimension (128 for WideResNes28-2) down to a reduced 8 dimensional space. The AAGMM and KMeans replacement layers with and without this reduced embedding space, were evaluated to determine whether the generative capabilities improve when 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 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 ac-

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	<b>94.64</b>	91.33	92.74	
	8	90.40	93.25	92.11	91.95
KMeans	128	84.08	91.62	92.22	
	8	<b>93.21</b>	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.

445 curacies that nearly match fully-supervised performance on  
 446 CIFAR-10 (similar to 250 label CIFAR-10).

447 Table 2 summarizes the relative performance of our  
 448 various configurations for both 40 and 250 labels. We repro-  
 449 duced FixMatch [35] using our hyper-parameters and didn't  
 450 quite matching the published performance at 40 labels (250  
 451 labels matched). Hyper-parameter selection for our training  
 452 algorithm is likely sub-optimal for baseline FixMatch. The  
 453 "Linear (i.e. FullyConnected)" rows in table 2 represent the  
 454 baseline fully connected linear last layer without additional  
 455 embedding dimensionality projection.

456 For CIFAR-10 with 250 labels, all last layers perform  
 457 reasonably close to semi-supervised SOTA, which itself is  
 458 almost identical to the fully supervised CIFAR-10 test ac-  
 459 curacy of 95.38% [41].

460 In addition to average performance, it is informative to  
 461 examine the max test accuracy over the  $N = 5$  random  
 462 trials to understand how well the algorithm can do, with  
 463 samples that are representative of the larger dataset. Table  
 464 3 demonstrates that in the best case, the AAGMM can get  
 465 within 1% of SOTA [51] performance.

466 The modeled cluster centers vary in quality between  
 467 individual model runs of the AAGMM layer due to the  
 468 stochasticity of the training process. Figure 3 (a & b) show-  
 469 cases degenerate cluster centers. The Figure 3 (c) AAGMM  
 470 model learned cluster centers that are an ok approximation  
 471 of the underlying data. However, the embedding constraints  
 472 encourage cluster centers which are better aligned with the  
 473 underlying data, Figure 3 (d). It is worth noting that we  
 474 did not observed the KMeans layer learning non-degenerate  
 475 cluster centers without an embedding constraint. In con-  
 476 trast, the AAGMM layer can, under some circumstances,  
 477 learn viable cluster centers.

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

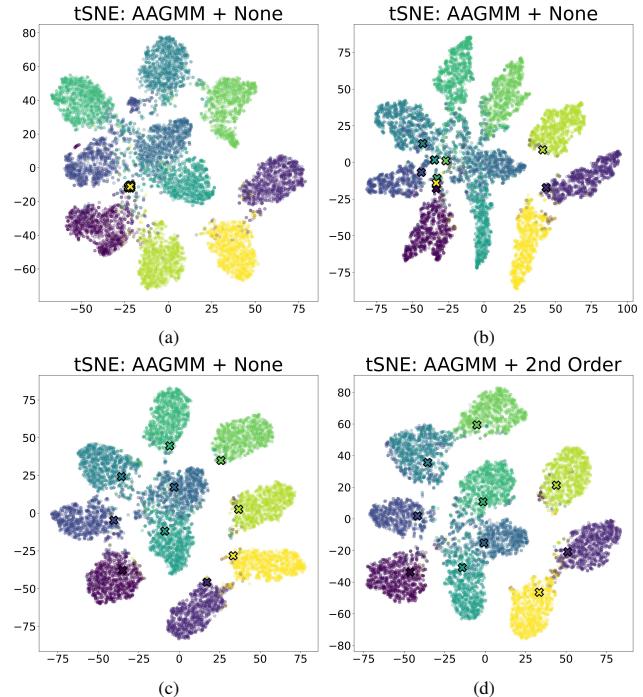


Figure 3. t-SNE plot of fully trained and reasonably accuracy AAGMM model's latent embedding space (with the learned cluster centers marked with X's). Depending on the run, the AAGMM cluster centers will be degenerate (top left), non-degenerate but still mis-aligned with the clusters (top right), acceptably 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. 522
- [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. 527
- [3] David Berthelot, Nicholas Carlini, Ekin D. Cubuk, Alex Kurakin, Kihyuk Sohn, Han Zhang, and Colin Raffel. Remmixmatch: Semi-supervised learning with distribution matching and augmentation anchoring. In *International Conference on Learning Representations*, 2020. 531
- [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. 536
- [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. 540
- [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. 544
- [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. 549
- [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. 554
- [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. 558
- [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. 562
- [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* 566

- 570       on *Computer Vision and Pattern Recognition*, pages 3187–  
571       3197, 2023. 3
- 572 [12] Arthur Gretton, Kenji Fukumizu, Choon Teo, Le Song, Bern-  
573       hard Schölkopf, and Alex Smola. A kernel statistical test of  
574       independence. *Advances in neural information processing systems*, 20, 2007. 3
- 575 [13] Mohamed Farouk Abdel Hady and Friedhelm Schwenker.  
576       Semi-supervised learning. *Handbook on Neural Information  
577       Processing*, pages 215–239, 2013. 2
- 578 [14] Jiwon Kim, Youngjo Min, Daehwan Kim, Gyuseong Lee,  
579       Junyoung Seo, Kwangrok Ryoo, and Seungryong Kim. Comatch:  
580       Semi-supervised learning with confidence-guided  
581       consistency regularization. In *European Conference on  
582       Computer Vision*, pages 674–690. Springer, 2022. 2, 3
- 583 [15] Diederik P Kingma and Max Welling. Auto-encoding vari-  
584       ational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 2, 3
- 585 [16] Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen,  
586       Ilya Sutskever, and Max Welling. Improved variational in-  
587       ference with inverse autoregressive flow. *Advances in neural  
588       information processing systems*, 29, 2016. 3
- 589 [17] Diederik P Kingma, Max Welling, et al. An introduction to  
590       variational autoencoders. *Foundations and Trends® in Ma-  
591       chine Learning*, 12(4):307–392, 2019. 2
- 592 [18] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple  
593       layers of features from tiny images. Technical report, 2009.  
594       2, 5
- 595 [19] Doyup Lee, Sungwoong Kim, Ildoo Kim, Yeongjae Cheon,  
596       Minsu Cho, and Wook-Shin Han. Contrastive regulariza-  
597       tion for semi-supervised learning. In *Proceedings of the  
598       IEEE/CVF Conference on Computer Vision and Pattern  
599       Recognition*, pages 3911–3920, 2022. 2
- 600 [20] Junnan Li, Caiming Xiong, and Steven CH Hoi. Comatch:  
601       Semi-supervised learning with contrastive graph regulariza-  
602       tion. In *Proceedings of the IEEE/CVF International Confer-  
603       ence on Computer Vision*, pages 9475–9484, 2021. 2, 3
- 604 [21] Yujia Li, Kevin Swersky, and Rich Zemel. Generative mo-  
605       ment matching networks. In *International conference on ma-  
606       chine learning*, pages 1718–1727. PMLR, 2015. 3
- 607 [22] Yang Li, Quan Pan, Suhang Wang, Haiyun Peng, Tao Yang,  
608       and Erik Cambria. Disentangled variational auto-encoder for  
609       semi-supervised learning. *Information Sciences*, 482:73–85,  
610       2019. 2
- 611 [23] Yu-Feng Li and De-Ming Liang. Safe semi-supervised learn-  
612       ing: a brief introduction. *Frontiers of Computer Science*, 13:  
613       669–676, 2019. 2
- 614 [24] Ioannis E Livieris, Konstantina Drakopoulou, Vassilis T  
615       Tampakas, Tassos A Mikropoulos, and Panagiotis Pintelas.  
616       Predicting secondary school students’ performance utilizing  
617       a semi-supervised learning approach. *Journal of educational  
618       computing research*, 57(2):448–470, 2019. 3
- 619 [25] David McClosky, Eugene Charniak, and Mark Johnson.  
620       Reranking and self-training for parser adaptation. In *Pro-  
621       ceedings of the 21st International Conference on Compu-  
622       tational Linguistics and 44th Annual Meeting of the Associa-  
623       tion for Computational Linguistics*, pages 337–344, 2006. 3
- 624 [26] Sumeet Menon. *Semi-Supervised Expectation Maximiza-  
625       tion with Contrastive Outlier Removal*. PhD thesis, 2022.  
626       AAI29167191. 3
- 627 [27] Sumeet Menon, David Chapman, Phuong Nguyen, Ye-  
628       lena Yesha, Michael Morris, and Babak Saboury. Deep  
629       expectation-maximization for semi-supervised lung cancer  
630       screening. 2019. 3
- 631 [28] Aamir Mustafa and Rafał K Mantiuk. Transformation  
632       consistency regularization—a semi-supervised paradigm for  
633       image-to-image translation. In *European Conference on  
634       Computer Vision*, pages 599–615. Springer, 2020. 3
- 635 [29] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer,  
636       James Bradbury, Gregory Chanan, Trevor Killeen, Zeming  
637       Lin, Natalia Gimelshein, Luca Antiga, Alban Desma-  
638       aison, Andreas Kopf, Edward Yang, Zachary DeVito, Mar-  
639       tin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit  
640       Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch:  
641       An imperative style, high-performance deep learning library.  
642       In *Advances in Neural Information Processing Systems* 32,  
643       pages 8024–8035. Curran Associates, Inc., 2019. 4
- 644 [30] Karl Pearson. Method of moments and method of maximum  
645       likelihood. *Biometrika*, 28(1/2):34–59, 1936. 4
- 646 [31] V Jothi Prakash and Dr LM Nithya. A survey on  
647       semi-supervised learning techniques. *arXiv preprint  
648       arXiv:1402.4645*, 2014. 3
- 649 [32] Danilo Rezende and Shakir Mohamed. Variational inference  
650       with normalizing flows. In *International conference on ma-  
651       chine learning*, pages 1530–1538. PMLR, 2015. 3
- 652 [33] Mamshad Nayeem Rizve, Kevin Duarte, Yogesh S Rawat,  
653       and Mubarak Shah. In defense of pseudo-labeling: An  
654       uncertainty-aware pseudo-label selection framework for  
655       semi-supervised learning. In *International Conference on  
656       Learning Representations*, 2021. 3
- 657 [34] Chuck Rosenberg, Martial Hebert, and Henry Schneiderman.  
658       Semi-supervised self-training of object detection models. In  
659       *2005 Seventh IEEE Workshops on Applications of Computer  
660       Vision (WACV/MOTION'05) - Volume 1*, pages 29–36, 2005.  
661       3
- 662 [35] Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao  
663       Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk,  
664       Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying  
665       semi-supervised learning with consistency and confidence.  
666       *Advances in neural information processing systems*, 33:596–  
667       608, 2020. 2, 3, 5, 6, 7, 8
- 668 [36] Antti Tarvainen and Harri Valpola. Mean teachers are better  
669       role models: Weight-averaged consistency targets improve  
670       semi-supervised deep learning results. *Advances in neural  
671       information processing systems*, 30, 2017. 3
- 672 [37] Laurens Van der Maaten and Geoffrey Hinton. Visualizing  
673       data using t-sne. *Journal of machine learning research*, 9  
674       (11), 2008. 2
- 675 [38] Mei Wang and Weihong Deng. Deep visual domain adap-  
676       tation: A survey. *Neurocomputing*, 312:135–153, 2018. 3
- 677 [39] Wei Wang, Haojie Li, Zhengming Ding, Feiping Nie, Jun-  
678       yang Chen, Xiao Dong, and Zhihui Wang. Rethinking maxi-  
679       mum mean discrepancy for visual domain adaptation. *IEEE  
680       Transactions on Neural Networks and Learning Systems*, 34  
681       (1):264–277, 2023. 3
- 682 [40] Yidong Wang, Hao Chen, Yue Fan, Wang Sun, Ran Tao,  
683       Wenxin Hou, Renjie Wang, Linyi Yang, Zhi Zhou, Lan-Zhe  
684

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