

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

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

Parniyan Farvardin  
University of Miami  
pxf291@miami.edu

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

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

## Abstract

Discriminative deep learning models with a linear+softmax final layer have a problem: the latent space only predicts the conditional probabilities  $p(Y|X)$  but not the full joint distribution  $p(Y, X)$ , which necessitates a generative approach. The conditional probability cannot detect outliers, causing outlier sensitivity in softmax networks. This exacerbates model over-confidence impacting many problems, such as hallucinations, confounding biases, and dependence on large datasets. To address this 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 FlexMatch with our technique. We find our MoM constraint with the AAGMM layer is able to match the reported FlexMatch accuracy, while also modeling the joint distribution, thereby reducing outlier sensitivity. We also present a preliminary outlier detection strategy based on Mahalanobis distance and discuss future improvements to this strategy. Code is available at: <https://github.com/mmajurski/ssl-gmm>

## 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 pseudo-distribution of labels  $Y$  which argmax

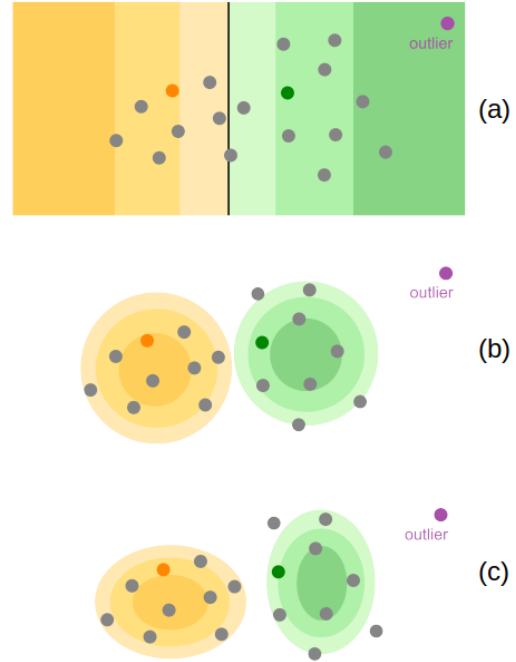


Figure 1. Schematic of the outlier problem, and how generative modeling of the joint probability can improve the situation. Prediction with (a) fully-supervised softmax, (b) semi-supervised KMeans, and (c) semi-supervised AAGMM.

can convert into a hard label. If  $X$  is distant from the decision boundary, then by definition, softmax assigns a prediction  $Y$  with high confidence. This works well for inlier samples, well represented by the training distribution. However, when presented with an outlier  $X$ , it is likely  $X$  will not be near the decision boundary (Figure 1, top). Therefore, softmax perceptrons, by definition, over-confidently hallucinate when given unexpected inputs [41]. Most deep

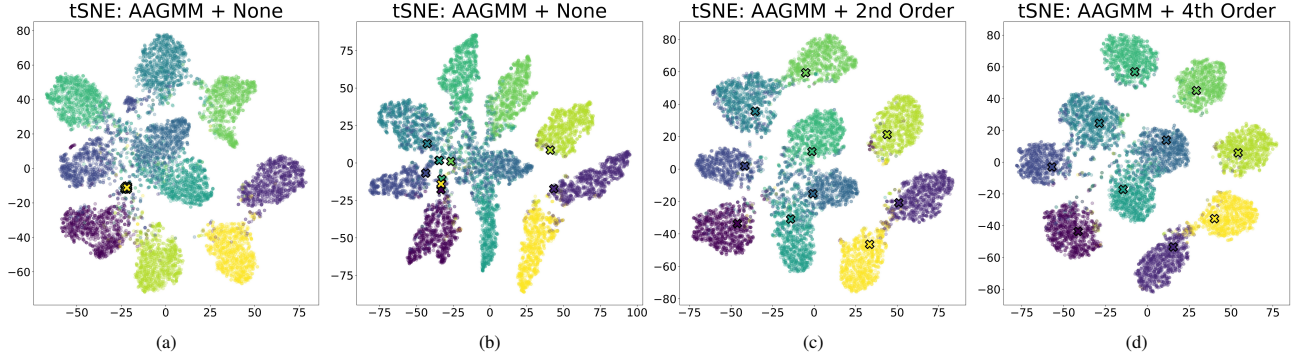


Figure 2. t-SNE[36] plot of the latent embedding space for various final layers with different MoM embedding constraints. AAGMM results include the explicitly modeled cluster centered marked with "X"s. Depending on the run, AAGMM cluster centers can be degenerate (a), non-degenerate but still mis-aligned with the clusters (b), acceptably aligned (c), or well aligned with the underlying clusters (d).

classifiers use softmax without a safety net and as such over-confidently predict  $Y$ . Ideally, when input  $X$  is distant from the decision boundary and training exemplars, the model should not be confident about the output class label  $Y$ .

Replacing softmax with a generative method that models the joint probability  $p(Y, X)$  can improve the capability of deep classifiers. Models using a final layer capable of learning the joint probability  $p(Y, X)$  can infer the conditional  $p(Y|X)$ . More importantly, such a layer can also infer the prior probability  $p(X)$ . Thus, if  $X$  is an unexpected input, then such a layer can flag the input as a low-probability sample (equivalently a high-probability outlier), rather than simply confidently predicting a label.

Prior work has explored generative modeling for image classification [16, 18, 23]. How to best train and utilize generative modeling within a deep learning context remains an open question. The naive approach of minimizing cross entropy between  $Y$  and  $\hat{Y}$  will not work. Figure 2 (a & b) shows t-SNE plots of the model latent embedding space for 94% accurate semi-supervised CIFAR-10 [19] image classification models, demonstrating why the naive approach will not work as intended. 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 (c & d) demonstrates that with our proposed AAGMM final layer with either 2nd or 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 [14, 24, 51]. Contrastive learning methods leverage the intuition that similar instances should be close in the representation space, while different instances are farther apart [21, 43]. Consistency regularization borrows the intuition that modified views of the same instance should have similar representations and predictions [15, 20, 34, 47]. 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 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 preliminary outlier removal strategy based on Mahalanobis distance that is compatible with AAGMM and MoM techniques.

We apply this methodology to the task of semi-supervised image-classification using CIFAR-10 [19] and STL-10 [10] with 40 training labels. Only 40 labels were used to provide a sufficiently difficult SSL problem, as SOTA results with higher numbers of labeled samples are approaching fully supervised performance [50]. 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 FlexMatch [47]) 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 [47]. The goal of SSL is to produce a trained model of equivalent accuracy to fully supervised training, with vastly reduced data annotation requirements. Doing so relies on accurately characterizing inlier vs outlier unlabeled samples.

### 2.1. Pseudo-Labeling

Self-training was among the initial approaches employed in the context of SSL to annotate unlabeled images. This involves the initial training of a classifier with a limited set of labeled samples which incorporates pseudo-labels exceeding a predefined threshold into the gradient descent process [9, 25–27, 33, 44, 46]. A closely related method to self-training is co-training, where a given dataset is represented as two distinct feature sets [4]. These independent sample sets are subsequently trained separately using two distinct models, and the sample predictions surpassing predetermined thresholds are utilized in the final model training process [4, 30]. A notable approach to pseudo-labeling is the Mean Teacher algorithm [35], which leverages exponential moving averages of model parameters to acquire a notably more stable target prediction. This refinement substantially enhances the convergence of the algorithm.

Several papers have attempted to enhance the quality of pseudo-labels to either improve the final model accuracy, improve the rate of convergence, or avoid confirmation bias [1]. Rizve et al. [32] explore how uncertainty aware pseudo-label selection/filtering can be used to reduce the label noise. Incorrect pseudo-labels can be viewed as a network calibration issue [32] where better network logit calibration might improve results [42]. Improvements to the pseudo-labeling process have been demonstrated by imposing curriculum [47] or by including a class-aware contrastive term [43]. Leveraging the concept of explicit class cluster centers for conditioning semantic similarity improves final model accuracy [49]. Additionally, improvements have been found in incorporating purely clustering based methods like DINO [7] into semi-supervised methods [12].

### 2.2. Consistency Regularization

Consistency Regularization is a branch of techniques that have been instrumental toward many of the state of the art techniques in semi-supervised learning within the last several years [2, 3, 15, 21, 34, 47, 49, 50]. The idea being that augmentation does not typically change the meaning of images. MixMatch is a semi-supervised pseudolabeling that greatly popularized the use of consistency regularization to ensure that augmentation does not affect the predicted label [2, 3]. FixMatch further extended this method

by introducing the notion of weak and strong augmentations including the cutout operator to increase the robustness of the regularization [34]. FlexMatch is a further improvement that introduces a curriculum pseudo-labeling strategy for flexible threshold values [47]. Co-Match made use of a form of consistency regularization to ensure that strong augmentations shared not only a similar pseudolabel, but furthermore a similar embedding space. Moreover, a neighborhood graph was constructed for embeddings and pseudolabels and refined via co-learning [21]. Con-match introduced a confidence metric based on the similarity of a basket of augmented embeddings [15]. SimMatch introduced a graph-based label propagation algorithm through a low-dimensional latent projection, and utilized multiple forms of consistency regularization including both semantic-level and instance-level consistency terms [49, 50].

### 2.3. Latent Embedding Constraints

A notable latent embedding constraint that is related, yet substantially different from our approach is the Evidence Lower Bound (ELBO) [16]. ELBO approximates a latent sample with a variational distribution and constrains the KL-divergence between the variational distribution and a target shape which is typically a multivariate standard normal distribution [16]. The main drawback of this approach is that the true KL-divergence is intractable to calculate. As such, the posterior must take on a simplified form. Most practical implementations use a diagonal posterior which can only penalize simple differences in shape such as mean and standard deviation. Arbitrarily complex posteriors are nevertheless possible using the method of Normalizing flows [8, 17, 31], which provides an iterative framework based on change of variables although this method is quite involved. Our MoM constraint is relatively simple but can also penalize complex differences in shape by constraining 2nd, 3rd, and 4th order hyper-covariance matrices, although we do so by comparing the moments directly. This greatly simplifies implementation as we do not need to explicitly construct a posterior distribution.

Another notable embedding constraint is the Maximum Mean Discrepancy (MMD), also known as the two-sample test [13]. MMD was used in the Generative Moment Matching Network [22] and has since been used extensively for the problem of domain adaptation [37, 38], in order to constrain the latent projections of the source and target distributions to follow the same distribution. MMD is a moment-matching constraint based on the kernel trick and can therefore constrain any difference in shape between two samples including very high order moments. Due to the kernel trick requiring proper inner products, MMD can only be used to constrain one sample to another sample. It cannot directly constrain sample statistics to population statistics, although it is possible to approximate populations nu-

merically via monte-carlo sampling [48]. Like MMD, our method is based on MoM, but it does not involve the kernel trick, and instead penalizes polynomial moments explicitly thereby enabling the sample embedding to be constrained to an exact target distribution.

### 3. Methodology

In this section, we explore our proposed replacement final activation layers and our embedding space constraints. Our methodology is based upon the published FlexMatch [47] algorithm as implemented in the USB framework [39], with identical hyper-parameters unless otherwise stated. We extend FlexMatch with a few minor training algorithm modifications explored in Section 4.1. FlexMatch [47] is a simple, well performing SSL algorithm. As such, it serves as a good comparison point for exploring the effect of our contributions.

Both the linear layer replacements and the embedding constraints explored herein represent increasing levels of prescription about how the final latent embedding space should be arranged compared to a traditional linear layer. The idea of leveraging clusters in embedding space is not new [5, 6, 11], 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

As we discussed in the introduction traditional final activation layers such as linear+softmax are fully discriminative in that they directly estimate the conditional probability  $p(Y|X)$ . These layers do not estimate  $p(X)$  or the joint probabilities  $p(Y, X)$ . To overcome this limitation, we present two generative 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  $\mu_k$  and cluster covariance  $\Sigma_k$ . The prior probability of any given sample  $X_i$  is defined by the mixture of cluster probability densities over the  $D$ -dimensional latent representation  $Z_i$  as follows,

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

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

Where  $\mathcal{N}(Z_i, \mu_k, \Sigma_k)$  represents the multivariate Gaussian pdf with centroid  $\mu_k$  and covariance  $\Sigma_k$ . AAGMM is axis aligned because  $\Sigma_k$  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 conditional 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 [28] module. It has two parameters updated by backprop. (1) the explicit cluster centers, a matrix  $K \times D$  initialized randomly, and (2) the diagonal elements of the  $D \times D$  matrix  $\Sigma_k$  are 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  $\Sigma_k$  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. See the published codebase for implementation details about the AAGMM and KMeans layers.

### 3.1.3 Relation between K-means, AAGMM and Softmax layers

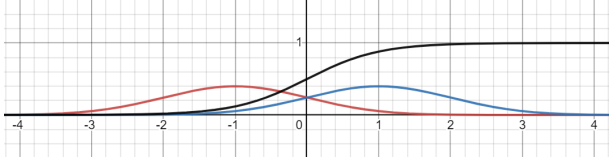


Figure 3. Illustrated relationship between K-means and Softmax Layers for 1D binary classification. If the joint distributions (blue curve and red curve) follow equivariate normal probability densities, then the conditional distribution (black curve) is softmax.

The AAGMM, KMeans and Softmax layers are theoretically related because Softmax is the conditional distribution  $p(Y|X)$  that arises when the joint distributions  $p(Y, X)$  follow the equivariate normal distributions as modeled by KMeans. Figure 3 shows a simple example of this relationship in one dimension with two labels  $A$  (blue curve) and  $B$  (red curve), with joint distributions as follows,

$$\begin{aligned} p(Y = A, X) &= \frac{1}{2} \mathcal{N}(X, \mu_A, \sigma) \\ p(Y = B, X) &= \frac{1}{2} \mathcal{N}(X, \mu_B, \sigma) \end{aligned} \quad (5)$$

In this case the conditional distribution can be described by sigmoid which is a special case of Softmax,

$$\begin{aligned} p(Y = A | X) &= \text{sigmoid}(mX + b) \\ \text{where} \\ m &= \frac{\mu_A - \mu_B}{\sigma^2} \quad \text{and} \quad b = \frac{\mu_A^2 + \mu_B^2}{2\sigma^2} \end{aligned} \quad (6)$$

The AAGMM layer is a generalization of the KMeans layer to allow different diagonal covariance matrices for each cluster. This gives the AAGMM somewhat greater expressive power than Softmax and KMeans as it can model the joint distribution of latent spaces with different cluster sizes and non-spherical shapes.

## 3.2. Method of Moments Embedding Constraints

We introduce and evaluate a series of embedding constraints based on the Method of Moments (MoM) [29]. For each sample  $i$  and each cluster  $k$ , the joint  $p(Y_{i,k}, X_i)$  is calculated as in Eq. 3 and then used to infer the prior  $p(X_i)$  and the conditional  $p(Y_{i,k}|X_i)$ . As usual, the conditional

probability 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 conditional probability without modeling the latent joint distribution. MoM solves these problems and is an appropriate strategy for semi-parametric models.

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  $P$  population statistic as follows,

$$\forall p = 1 \dots P \quad \frac{1}{n} \sum_{i=1}^n g_p(z_i) \approx E(g_p(Z)) \quad (7)$$

We can therefore constrain the latent representation of our model to approximate a multivariate standard normal distribution.

The centralized moments are a classical choice for the consistent estimator  $g_p$  representing the terms of a power series around the mean  $\mu$

$$g_p(Z) = (Z - \mu)^p \quad (8)$$

In the univariate standard normal case, the  $p^{th}$  order centralized moment constraint is the following.

$$E[(Z - \mu)^p] = \begin{cases} 0 & \text{if } p \text{ is odd} \\ \sigma^p (p-1)!! & \text{if } p \text{ is even} \end{cases} \quad (9)$$

Where '!!' represents the double factorial operator. By this formula, the univariate unit Gaussian has mean 0, standard deviation 1, skew 0, and kurtosis 3.

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

$$E[g_p(Z - \mu)] = E \left[ \prod_{d=1}^D (Z_d - \mu_d)^{p_d} \right] \quad (10)$$

Due to independence of the axes, this multivariate population moment can be represented as a product of univariate moments of the individual standard normal distributions as follows,

$$E \left[ \prod_{d=1}^D (Z_d - \mu_d)^{p_d} \right] = \prod_{d=1}^D E[(Z_d - \mu_d)^{p_d}] \quad (11)$$

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

$$\varepsilon_p = \left( \frac{1}{n} \sum_{i=1}^n g_p(z_i) - E(g_p(Z)) \right)^2 \quad (12)$$

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 of 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.

### 3.3. Mahalanobis Outlier Removal

The AAGMM layer allows us to detect and remove outliers based on Mahalanobis distance in the latent feature space. By *outlier*, we are referring to the problem that the pseudolabel learner (i.e. FlexMatch) is simply not yet ready to learn a given unlabeled sample, because the model has only attempted to learn the distribution of labeled and previously pseudolabeled samples up until that point. Due to small labeled sample size, the labeled and pseudolabeled samples do not fully represent the distribution of the unlabeled samples in early iterations. Thus, unlabeled samples far from the learned distribution are considered to be outliers in a given iteration.

In order to implement *outlier removal*, in the context of pseudolabeling, we exclude the unlabeled *detected outlier* sample from gradient updates for any the given iteration of the semi-supervised procedure. For FlexMatch in particular, all unlabeled samples are augmented with weak and strong RandAugment. As such we use the weakly augmented samples as input to outlier detection.

Mathematically, we consider an unlabeled sample  $x$  to be an outlier if it is distance from all cluster centers  $\mu_1 \dots \mu_K$

in terms of Mahalanobis distance based on the cluster covariances  $\Sigma_1 \dots \Sigma_K$ . Mahalanobis distance of a given point  $x$  to cluster  $k$  is defined as follows,

$$d_{M,k}(x) = \sqrt{(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)} \quad (13)$$

Define  $X_L$  as the labeled population and  $X_U$  as the unlabeled population, with  $x \in X_U$  as an unlabeled sample, and  $P_{90}$  as the 90th percentile. As such, we detect outliers as follows,

$$\begin{aligned} x \in X_U \text{ is an outlier iff.} \\ \max_k d_{M,k}(x) > \tau \\ \text{where } \tau = P_{90}(\max_k d_{M,k}(X_L)) \end{aligned} \quad (14)$$

## 4. Experiments

We evaluate both AAGMM and KMeans linear layer replacements and the embedding space constraints using our modified FlexMatch [47] on the common SSL benchmarks CIFAR-10 [19] at 40 labels (4 labels per class) and STL-10 [10] at 40 labels. We randomly selected 5 seeds 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 are used as unlabeled data (i.e., the labels are discarded).

### 4.1. Hyper-Parameters

All CIFAR-10 models were trained with the standard benchmark WideResNet28-2 architecture [45]. All STL-10 models were trained with the standard benchmark WideResNet37-2 architecture [45]. This work leverages the published FlexMatch [47] code, hyper-parameters, and training configurations within the USB Framework [39]. However, due to the higher training instability of the AAGMM layers compared to a linear layer, the gradient norm was clipped to 1.0. Despite specific attention to computing the AAGMM and embedding constraints in a numerically stable manner, they are still less stable during backprop than a simple linear layer. Gradient clipping was especially important when there were latent points that were multiple standard deviations away from the cluster centers. In this case, the gradient of the Gaussian probability density function converges rapidly toward zero, which can affect the division step of equation 4. We have almost entirely overcome this issue in the code by using laws of exponents in order to normalize the denominator to be greater or equal to 1 prior to division, but in extreme cases of latent points distant from cluster centers, the gradient clipping is still necessary to achieve stable gradient descent.

Mean Test Accuracy Per Method and Dataset

Last Layer	Emb Dim	CIFAR-10 at 40 Labels (5 trials)		
Embedding Constraint		None	1st Order	2nd Order
Linear (unmodified FlexMatch)	128	95.03 $\pm$ 0.06		
Linear (FlexMatch w/ Reduced Emb Dim)	8	90.89 $\pm$ 3.24		
AAGMM	128	94.66 $\pm$ 0.32	94.80 $\pm$ 0.14	94.83 $\pm$ 0.07
AAGMM	8	94.98 $\pm$ 0.07	94.64 $\pm$ 0.27	93.58 $\pm$ 2.74
KMeans	8	90.15 $\pm$ 6.40	93.50 $\pm$ 0.62	93.44 $\pm$ 0.50
Last Layer	Emb Dim	STL-10 at 40 Labels (3 trials)		
Embedding Constraint		None	1st Order	2nd Order
Linear (unmodified FlexMatch)	128	70.85 $\pm$ 4.16		
AAGMM	8	58.79 $\pm$ 11.00	71.11 $\pm$ 7.60	70.40 $\pm$ 6.39

Table 1. Mean test accuracy % for CIFAR-10 and STL-10 SSL benchmarks comparing various configurations of our method. The FlexMatch results in the table is drawn from the publication. For CIFAR-10, the WideResNet model used by FlexMatch has an embedding size of 128 dimension. Results for a given order of embedding constraint include all lower constraints.

We believe that the rapid decrease in slope of the Gaussian distribution pdf for points that are multiple standard deviations away from the mean. As the conditional distribution involves calculating a division step, it is very likely that such a division may become less stable for points that are not near the mean.

This work includes an exploration of how various latent embedding dimensionalities affect 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 WideResNet28-2) down to a reduced 8 dimensional space. Results listed with an embedding dimensionality of 128 do not include the additional linear layer which reduces the latent dimensionality.

Due to exponential GPU memory requirements with each successive MoM moment, only the 8D embedding can operate with higher order MoM embedding constraints. While our method can place constraints on any number of moments, we only explored MoM constraints up to the second order in this paper. Results for any given order of embedding constraint include all lower constraints.

## 4.2. CIFAR-10 and STL-10 Results

While current SOTA on CIFAR-10 at 250 labels is close to fully supervised accuracy, the 40 label case provides a far more challenging task. Table 1 summarizes the relative performance of our various configurations. As we simply extended FlexMatch, our hyper-parameter selection is likely sub-optimal, and further experimentation might yield improvements. The 'Linear (unmodified FlexMatch)' row in Table 1 represents the baseline fully connected linear last layer exactly as FlexMatch published. The 'Linear (FlexMatch w/ Reduced Emb Dim)' row in Table 1 represents FlexMatch performance where its final Linear layer pro-

duced a latent embedding dimensionality of 8 (instead of the stock model latent dimensionality of 128). Its worth noting that the significantly smaller embedding dimensionality reduced average model accuracy by 5%. That accuracy is restored by replacing the Linear layer with an AAGMM layer; despite keeping the reduced latent embedding dimensionality. That trend is mirrored with STL-10, where the 70.85% FlexMatch accuracy drops to 58.79% using a latent embedding dimensionality of 8 (and the unconstrained AAGMM layer), which the 1st order embedding constrained version of AAGMM restores to 71.11% accuracy.

We see that the AAGMM final layer consistently outperforms the KMeans final layer for the CIFAR-10 Test Accuracy with 40 labels. We furthermore see that the KMeans layer performs significantly better with the 1st and 2nd order MoM constraint, as compared with no constraints.

For STL-10 the AAGMM final layer (71.11%) improved upon the FlexMatch [47] result (70.85%), though within the margin of error. Additionally, both 1st and 2nd order MoM constraints significantly improved the upon the AAGMM with no embedding constraints.

The modeled cluster centers vary in quality between individual model runs of the AAGMM layer due to the stochasticity of the training process. Figure 2 (a & b) show cases degenerate cluster centers. The Figure 2 (c & d) AAGMM model learned cluster centers that are an adequate approximation of the underlying data, where the embedding constraints encourage cluster centers which are better aligned with the underlying data. It is worth noting that we did not observe 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. To quantify the modeled cluster compactness, we measure the average L2 distance from each test data point to its assigned cluster as shown in Table

2.

Latent Embedding Space Cluster Compactness

Dataset	Last Layer	None	1st MoM	2nd MoM
CIFAR-10	AAGMM	1.03	0.58	0.53
CIFAR-10	KMeans	18.41	0.96	1.06
STL-10	AAGMM	2.32	0.79	0.76

Table 2. Average L2 distance from each test data point to its assigned cluster center.

To place our results in context with the current SSL SOTA for CIFAR-10 and STL-10 at 40 labels Table 3 compares against the current best methods, demonstrating that this methodology is nearly competitive with for CIFAR-10, but still requires improvement and tuning for STL-10.

Method	CIFAR-10 (40 Labels)	STL-10 (40 Labels)
FixMatch[34]	13.81 $\pm$ 3.37	35.97 $\pm$ 4.14
FlexMatch[47]	4.97 $\pm$ 0.06	29.15 $\pm$ 4.16
FreeMatch[40]	4.90 $\pm$ 0.29	15.56 $\pm$ 0.55
SimMatchV2[50]	4.90 $\pm$ 0.04	15.85 $\pm$ 2.62
<b>AAGMM+None</b>	5.02 $\pm$ 0.07	41.21 $\pm$ 11.00
<b>AAGMM+1stOrder</b>	5.36 $\pm$ 0.27	28.89 $\pm$ 6.0

Table 3. Error rate % for CIFAR-10 and STL-10 SSL benchmarks with 40 labels, comparing to state of the art results. Results for previously published methods are drawn from USB [39] except for FreeMatch [40] and SimMatchV2 [50] publications.

## 5. Discussion

Although our preliminary results with the proposed AAGMM and MoM achieve high accuracy relative to SOTA, this was achieved without Mahalanobis outlier detection as documented in Table 4. When the outlier detection was enabled, we believe the reduced accuracy is due to: 1. the 90th percentile distance threshold being too aggressive and filtering too much signal relative to noise, and 2. the need for an adaptive outlier detection threshold. In early epochs, an aggressive outlier detection threshold is viable because the model will not adequately fit many of the unlabeled samples, but as the model converges the fit improves reducing the need for outlier removal. As such, we believe that the aggressive outlier filtering is removing too many inliers, particularly in later epochs.

The proposed MoM embedding constraint has at least one significant downside, by requiring exponentially increasing amounts of GPU memory for each successive moment penalty included as shown in Table 5.

This limits the current practicality of these MoM constraints. Additional optimization and/or avoiding the ex-

AAGMM (8D) Mean Test Accuracy With Outlier Removal

CIFAR-10		
Outlier Threshold	None	90th Percentile
MoM: None	94.98 $\pm$ 0.07	94.9 $\pm$ 0.125
MoM: 1st Order	94.64 $\pm$ 0.27	87.70 $\pm$ 2.96
MoM: 2nd Order	93.58 $\pm$ 2.74	87.25 $\pm$ 2.51
STL-10		
Outlier Threshold	None	90th Percentile
MoM: None	58.79 $\pm$ 11.00	57.50 $\pm$ 12.75
MoM: 1st Order	71.11 $\pm$ 7.60	64.18 $\pm$ 3.82
MoM: 2nd Order	70.40 $\pm$ 6.39	65.90 $\pm$ 4.09

Table 4. Mean test accuracy % for CIFAR-10 and STL-10 SSL benchmarks comparing an 8D embedding AAGMM final layer with and without outlier removal during training. Results for a given order of embedding constraint include all lower constraints.

AAGMM Training GPU Memory Requirements (in GiB)

Emb Dim	None	1st	2nd	3rd	4th
8	7.72	7.71	7.70	7.76	8.76
32	7.71	7.71	7.79	13.15	> 20.47

Table 5. GPU RAM utilization in GiB on Nvidia RTX A4500 with 20.47GiB of VRAM evaluated on CIFAR10 using WRN28-2 with a batch size of 64. Each row shows the results for a given embedding dimensionality, and each row a MoM embedding constraint order, where "None" indicates a stock linear layer.

plicit creation of both the  $n^{th}$  order moment and its target value on GPUs would likely improve the usability.

Semi-supervised learning is highly sensitive to both which samples are selected from the labeled population [34] and the stochasticity of the training process itself.

Future work in this area will explore alternate outlier removal strategies, including thresholding the unlabeled samples based on their latent sample probability  $p(X)$  as opposed to latent Mahalanobis distance. The Mahalanobis distance is part of the exponential term in the calculation of the multivariate Gaussian PDF. As such, we expect that a removing outliers based on low  $p(X)$  is likely to perform comparably to removing samples based on far Mahalanobis distance, in the special case that all clusters share a similar determinant  $\det(\Sigma_k)$ . Additionally, we plan to 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 can be applied to semi-supervised learning. The combination of these techniques enables outlier detection strategies that would otherwise not be possible with a traditional softmax discriminator approach. This preliminary work constructs these novel layers with the associated constraints and demonstrates



reasonable performance on challenging benchmark semi-supervised learning tasks, while opening the door for future outlier detection strategies that can make semi-supervised learning more robust to large and diverse unlabeled sample distributions.

**Acknowledgements.** We would like to thank the Frost Institute for Data Science and Computing for their support.

## References

- [1] Eric Arazo, Diego Ortego, Paul Albert, Noel E O'Connor, and Kevin McGuinness. Pseudo-labeling and confirmation bias in deep semi-supervised learning. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2020. [2](#), [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. Remix-match: 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] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 215–223. JMLR Workshop and Conference Proceedings, 2011. [2](#), [6](#)
- [11] 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](#)
- [12] Enrico Fini, Pietro Astolfi, Karteek Alahari, Xavier Alameda-Pineda, Julien Mairal, Moin Nabi, and Elisa Ricci. Semi-supervised learning made simple with self-supervised clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3187–3197, 2023. [3](#)
- [13] Arthur Gretton, Kenji Fukumizu, Choon Teo, Le Song, Bernhard Schölkopf, and Alex Smola. A kernel statistical test of independence. *Advances in neural information processing systems*, 20, 2007. [3](#)
- [14] Mohamed Farouk Abdel Hady and Friedhelm Schwenker. Semi-supervised learning. *Handbook on Neural Information Processing*, pages 215–239, 2013. [2](#)
- [15] Jiwon Kim, Youngjo Min, Daehwan Kim, Gyuseong Lee, Junyoung Seo, Kwangrok Ryoo, and Seungryong Kim. Con-match: Semi-supervised learning with confidence-guided consistency regularization. In *European Conference on Computer Vision*, pages 674–690. Springer, 2022. [2](#), [3](#)
- [16] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. [2](#), [3](#)
- [17] Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improved variational inference with inverse autoregressive flow. *Advances in neural information processing systems*, 29, 2016. [3](#)
- [18] Diederik P Kingma, Max Welling, et al. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392, 2019. [2](#)
- [19] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, 2009. [2](#), [6](#)
- [20] Doyup Lee, Sungwoong Kim, Ildoo Kim, Yeongjae Cheon, Minsu Cho, and Wook-Shin Han. Contrastive regularization for semi-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3911–3920, 2022. [2](#)
- [21] Junnan Li, Caiming Xiong, and Steven CH Hoi. Comatch: Semi-supervised learning with contrastive graph regularization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9475–9484, 2021. [2](#), [3](#)
- [22] Yujia Li, Kevin Swersky, and Rich Zemel. Generative moment matching networks. In *International conference on machine learning*, pages 1718–1727. PMLR, 2015. [3](#)
- [23] Yang Li, Quan Pan, Suhang Wang, Haiyun Peng, Tao Yang, and Erik Cambria. Disentangled variational auto-encoder for semi-supervised learning. *Information Sciences*, 482:73–85, 2019. [2](#)
- [24] Yu-Feng Li and De-Ming Liang. Safe semi-supervised learning: a brief introduction. *Frontiers of Computer Science*, 13: 669–676, 2019. [2](#)
- [25] Ioannis E Livieris, Konstantina Drakopoulou, Vassilis T Tampakas, Tassos A Mikropoulos, and Panagiotis Pintelas. Predicting secondary school students' performance utilizing a semi-supervised learning approach. *Journal of educational computing research*, 57(2):448–470, 2019. [3](#)

- [26] David McClosky, Eugene Charniak, and Mark Johnson. Reranking and self-training for parser adaptation. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pages 337–344, 2006.
- [27] Sumeet Menon, David Chapman, Phuong Nguyen, Yelena Yesha, Michael Morris, and Babak Saboury. Deep expectation-maximization for semi-supervised lung cancer screening. 2019. **3**
- [28] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. **4**
- [29] Karl Pearson. Method of moments and method of maximum likelihood. *Biometrika*, 28(1/2):34–59, 1936. **5**
- [30] V Jothi Prakash and Dr LM Nithya. A survey on semi-supervised learning techniques. *arXiv preprint arXiv:1402.4645*, 2014. **3**
- [31] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pages 1530–1538. PMLR, 2015. **3**
- [32] Mamshad Nayeem Rizve, Kevin Duarte, Yogesh S Rawat, and Mubarak Shah. In defense of pseudo-labeling: An uncertainty-aware pseudo-label selection framework for semi-supervised learning. In *International Conference on Learning Representations*, 2021. **3**
- [33] Chuck Rosenberg, Martial Hebert, and Henry Schneiderman. Semi-supervised self-training of object detection models. In *2005 Seventh IEEE Workshops on Applications of Computer Vision (WACV/MOTION’05) - Volume 1*, pages 29–36, 2005. **3**
- [34] Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in neural information processing systems*, 33:596–608, 2020. **2, 3, 8**
- [35] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in neural information processing systems*, 30, 2017. **3**
- [36] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9 (11), 2008. **2**
- [37] Mei Wang and Weihong Deng. Deep visual domain adaptation: A survey. *Neurocomputing*, 312:135–153, 2018. **3**
- [38] Wei Wang, Haojie Li, Zhengming Ding, Feiping Nie, Junyang Chen, Xiao Dong, and Zhihui Wang. Rethinking maximum mean discrepancy for visual domain adaptation. *IEEE Transactions on Neural Networks and Learning Systems*, 34 (1):264–277, 2023. **3**
- [39] Yidong Wang, Hao Chen, Yue Fan, Wang Sun, Ran Tao, Wenxin Hou, Renjie Wang, Linyi Yang, Zhi Zhou, Lan-Zhe Guo, et al. Usb: A unified semi-supervised learning benchmark for classification. *Advances in Neural Information Processing Systems*, 35:3938–3961, 2022. **4, 6, 8**
- [40] Yidong Wang, Hao Chen, Qiang Heng, Wenxin Hou, Yue Fan, Zhen Wu, Jindong Wang, Marios Savvides, Takahiro Shinozaki, Bhiksha Raj, Bernt Schiele, and Xing Xie. Freematch: Self-adaptive thresholding for semi-supervised learning. In *The Eleventh International Conference on Learning Representations*, 2023. **8**
- [41] Hongxin Wei, Renchunzi Xie, Hao Cheng, Lei Feng, Bo An, and Yixuan Li. Mitigating neural network overconfidence with logit normalization. In *International Conference on Machine Learning*, pages 23631–23644. PMLR, 2022. **1**
- [42] Chen Xing, Sercan Arik, Zizhao Zhang, and Tomas Pfister. Distance-based learning from errors for confidence calibration. In *International Conference on Learning Representations*, 2020. **3**
- [43] Fan Yang, Kai Wu, Shuyi Zhang, Guannan Jiang, Yong Liu, Feng Zheng, Wei Zhang, Chengjie Wang, and Long Zeng. Class-aware contrastive semi-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14421–14430, 2022. **2, 3**
- [44] David Yarowsky. Unsupervised word sense disambiguation rivaling supervised methods. In *33rd annual meeting of the association for computational linguistics*, pages 189–196, 1995. **3**
- [45] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. 2016. **6**
- [46] Xiaohua Zhai, Avital Oliver, Alexander Kolesnikov, and Lucas Beyer. S4l: Self-supervised semi-supervised learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1476–1485, 2019. **3**
- [47] Bowen Zhang, Yidong Wang, Wenxin Hou, Hao Wu, Jindong Wang, Manabu Okumura, and Takahiro Shinozaki. Flexmatch: Boosting semi-supervised learning with curriculum pseudo labeling. *Advances in Neural Information Processing Systems*, 34:18408–18419, 2021. **2, 3, 4, 6, 7, 8**
- [48] Shengjia Zhao, Jiaming Song, and Stefano Ermon. Infovae: Balancing learning and inference in variational autoencoders. In *Proceedings of the aaai conference on artificial intelligence*, pages 5885–5892, 2019. **4**
- [49] Mingkai Zheng, Shan You, Lang Huang, Fei Wang, Chen Qian, and Chang Xu. Simmatch: Semi-supervised learning with similarity matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14471–14481, 2022. **3**
- [50] Mingkai Zheng, Shan You, Lang Huang, Chen Luo, Fei Wang, Chen Qian, and Chang Xu. Simmatchv2: Semi-supervised learning with graph consistency. pages 16432–16442, 2023. **2, 3, 8**
- [51] Xiaojin Zhu and Andrew B Goldberg. *Introduction to semi-supervised learning*. Springer Nature, 2022. **2**