# AutoEmbedder: A semi-supervised DNN embedding system for clustering

**Abu Quwsar Ohi**
Department of Computer Science & Engineering
Bangladesh University of Business & Technology
Dhaka, Bangladesh
quwsarohi@gmail.com

**M. F. Mridha**
Department of Computer Science & Engineering
Bangladesh University of Business & Technology
Dhaka, Bangladesh
firoz@bubt.edu.bd

**Farisa Benta Safir**
Department of Computer Science & Engineering
Bangladesh University of Business & Technology
Dhaka, Bangladesh
farisabentasafir@gmail.com

**Md. Abdul Hamid**
Department of Information Technology
Faculty of Computing & Information Technology
King Abdulaziz University
Jeddah-21589, Kingdom of Saudi Arabia
mabdulhamid1@kau.edu.sa

**Muhammad Mostafa Monowar**
Department of Information Technology
Faculty of Computing & Information Technology
King Abdulaziz University
Jeddah-21589, Kingdom of Saudi Arabia
mmonowar@kau.edu.sa

## ABSTRACT

Clustering is widely used in unsupervised learning method that deals with unlabeled data. Deep clustering has become a popular study area that relates clustering with Deep Neural Network (DNN) architecture. Deep clustering method downsamples high dimensional data, which may also relate clustering loss. Deep clustering is also introduced in semi-supervised learning (SSL). Most SSL methods depend on pairwise constraint information, which is a matrix containing knowledge if data pairs can be in the same cluster or not. This paper introduces a novel embedding system named AutoEmbedder, that downsamples higher dimensional data to clusterable embedding points. To the best of our knowledge, this is the first research endeavor that relates to traditional classifier DNN architecture with a pairwise loss reduction technique. The training process is semi-supervised and uses Siamese network architecture to compute pairwise constraint loss in the feature learning phase. The AutoEmbedder outperforms most of the existing DNN based semi-supervised methods tested on famous datasets.

***Keywords*** Deep Neural Network · Unsupervised learning · Semi-supervised learning · Transfer learning · Embedding · Clustering · Dimensionality reduction · Pairwise constraint

# 1   Introduction

Clustering is a fundamental approach to perform unsupervised learning. It is a very widely studied topic and is applied on a wide range of applications, including image segmentation [1], image processing [2], network analysis [3], document analysis [4, 5], and so on. Clustering remains an active research area due to its simplicity and ability to find a pattern in unlabeled data. Although clustering is a broadly used method, the performance of clustering methods degrades when it is applied to high dimensional data. To overcome the limitation of higher-dimensional data, researchers perform feature reduction strategies to reduce the higher dimensional features while keeping the necessary features.

Principal Component Analysis (PCA), is a common method, used for data dimensionality reduction [6, 7]. Dimensionality reduction of data can be achieved through feature extraction or feature selection. However, in this paper, we attain dimensionality reduction of data through PCA, which reduces the dimension of data through feature extraction only. The critical part of this process lies in ignoring the relativity between clustering and feature learning procedure. To eliminate this issue, Discriminative Cluster Analysis (DCA) was introduced [8]. The process combines Linear Discriminant Analysis (LDA) and K-means into a joint framework. However, the aforementioned method fails to represent a better estimation.

Currently, due to the recent advancement, Deep Neural Network (DNN) has been widely applied in supervised learning as well as unsupervised learning. The usage of DNN on clustering methods is often derived as deep clustering. Almost all the deep clustering architectures contain two phases: feature transformation and clustering. However, some deep clustering methods learn feature transformation and clustering jointly [9]. Although most unsupervised deep clustering methods fail to generate appreciable performance on complex datasets, the current state of the art models generate promising results on simple datasets [9, 10].

Some studies also use a feature mapping network. The most used one is the basic Convolutional Deep Neural Network (CDNN), which is pre-trained on a bigger dataset. The trained CDNN is further used to generate embedding points from unseen data. The embedding points are used to perform cluster [11]. This type of learning is often interpreted as transfer learning method. However, until now, no studies attempted to improve the accuracy of CDNN networks by calculating cluster loss.

Although clustering is unsupervised, there may exist some pre-knowledge of the dataset that is to be used for a particular task. This pre-knowledge is often used in Semi-Supervised Clustering (SSC). SSC methods rely on a pairwise constrained matrix to gain better accuracy than unsupervised deep clustering [12]. A pairwise constraint matrix contains information if two instances are related or not. If they are related, then they must be in the same cluster, otherwise, in a different cluster. SSC can use this information to improve its learning.

This paper contributes to a semi-supervised clustering process via DNN architecture. We introduce a semi-supervised embedding system named AutoEmbedder that is aimed to generate clusterable embedding points based on pairwise constraints. The AutoEmbedder is built upon traditional DNN architecture. The training process of AutoEmbedder uses pairwise constraints and this type of training procedure is termed as an SSL process [12]. The AutoEmbedder is iteratively trained on a Siamese Neural Network (SNN) architecture. SNN architecture uses two same weighted AutoEmbedders in parallel. Therefore, the SNN receives a pair of input and generates a pair of output. From the SNN architecture, a pairwise loss is computed by calculating the pairwise distance of the SNN-AutoEmbedder generated embedding. This loss is further reduced using the traditional backpropagation technique along with an optimization function. The AutoEmbedder is extracted from the SNN, and the finally trained AutoEmbedder is further used to generate meaningful embeddings, on which clustering is performed.

Overall, our main contributions include:

- We develop a semi-supervised embedding system named AutoEmbedder that learns to produce cluster separable meaningful embedding points based on pairwise constraints.

- We introduce DNN as an embedding system.

- We introduce a procedure of using DNN architectures that links to the embedding system and clustering.

- We carry out experiments addressing unsupervised and semi-supervised DNN based algorithms and validate that the AutoEmbedder performs better in most of the complex datasets.

The rest of this paper is organized as follows: Section 2 presents the related work. Section 3 introduces the overall architecture and AutoEmbedder training procedure. Section 4 contains empirical results done on the architecture. Finally, Section 5 concludes the paper.

## 2   Related Work

Unsupervised learning is the process of identifying unknown patterns from an unlabeled dataset. Before the advancement of neural networks, clustering algorithms were the only methods used to perform unsupervised learning. Clustering defines the process of separating data points based on their dissimilarity. Before the full phase implementation of the deep neural network, unsupervised clustering was performed based on generalization techniques. Through a generalization technique, a machine learning model can identify specific unseen example classes based on some specific feature patterns. Nevertheless, these unseen examples must contain the appropriate feature patterns. Otherwise, machine learning models may generate inappropriate cluster regions. This process of causing a general feature/concept between seen and unseen data is often termed as a generalization process. There exist generalization techniques based on PCA as well [13]. However, down-sampling higher dimensional data with PCA does not greatly improve clustering accuracy.

After the advancement of neural network architectures, researchers have introduced many unsupervised learning methods that are based on neural networks. Autoencoders (AE), Variational Autoencoders (VAE), Generative Adversarial Networks (GAN), Deep Belief Net (DBN), etc. are common architectures that are used to create the recent state of the art unsupervised learning architectures [14]. Most of the methods relying on autoencoders require a pre-training process for calibration. Deep Clustering Network (DCN) represents a combined approach of autoencoders and k-means algorithm [15]. Deep Embedding Network (DEN) only rely on reconstruction loss of autoencoders and converges to a cluster friendly representation [16]. Deep Continuous Clustering (DCC), Deep Embedded Regularized Clustering (DEPICT), Deep Multi-Manifold Clustering (DMC), and Deep Subspace Clustering Networks (DSC-Nets) perform similar reconstruction loss of autoencoders to perform clustering [17, 18, 19, 20]. The Deep Embedding Clustering (DEC) is a renowned model, which uses a pre-trained autoencoder [9]. Later, the model is fine-tuned using the combination of cluster loss and autoencoder loss, which is described as hardening loss. Unsupervised CDNN depends on features extracted from deep neural networks, and fine-tunes the network based on clustering loss, where clustering loss defines the error of combining different cluster data points into the same cluster and vice versa. It is also proven that a supervised CDNN architecture, pre-trained on a large dataset, acquire better accuracy than traditional unsupervised pre-trained CDNN based methods [21]. On the contrary, Joint Unsupervised Learning (JULE), and Deep Adaptive Image Clustering (DAC) are non-pre-trained CDNN architectures [22, 23]. The drawback of JULE is the computational cost and memory complexity of the learning process is greater than average for large datasets. Unlike JULE, DAC achieves better performance in some challenging datasets. Compared to the vast implementation of AE and CDNN, there are quite a few implementations based on VAE and GAN architectures. Most VAE and GAN based architectures fall behind due to their high complexity and hard to converge architectures. Gaussian Mixture VAE (GMVAE), and Variational Deep Embedding (VaDE) uses VAE architecture [24, 25]. Categorical Generative Adversarial (CatGAN), and Information Maximizing Generative Adversarial Network (InfoGAN) are well-recognized architectures that are based on GAN. Although these procedures present outstanding performance on simple datasets, they fail to produce a reasonable performance on complex datasets. To improve performance measures, semi-supervised training architecture is used.

SSL process contains a small portion of labeled data with a large portion of unlabeled data. In the aspect of clustering, there are some variants of semi-supervised methods [26]. In SSL, some data may contain information suitable to train the unsupervised architecture. This information might be the label of data or cluster linkage. Cluster linkage defines pairwise information by establishing connections among data pairs [27, 28, 29, 30]. Cluster linkage is also named as a pairwise constraint because it contains pairwise data linkage information. Most SSL use pairwise constraints as their basic training information. The main point of improving the overall process is selecting the appropriate data dimensional reduction technique or objective function, the pairwise constraints as well as the pairwise loss calculation. Gang Chen used a Restricted Boltzmann Machine (RBM) as the objective function [31]. Although the approach is promising, it fails to generate better accuracy with lesser data. A semi-supervised implementation of the renowned DEC method, named Semi-supervised Deep Embedded Clustering (SDEC), proved to give better accuracy than traditional DEC architecture [12]. However, the method fails to give better estimations on famous datasets.

To further strengthen the position, many neural network architectures on unsupervised learning that are pre-trained on complex datasets are used to perform clustering on the different datasets. This type of training is often referred to as transfer learning. Transfer learning relies on data dependency [32] and is actively used in many domains, including unsupervised learning [11]. Convolutional Neural Networks (CNN) are also proven to perform better on complex image datasets [21]. Therefore, feed-forward CNN is used to perform semi-supervised classification that is referred to as Semi-Supervised Feed-Forward CNN (SSFF-CNN) [33]. The SSFF-CNN architecture relates Feed-Forward CNN (FF-CNN) with SSL [34]. The parameters of FF-CNN target layers are generated by the statistics of the previous layers. Again, the CNN methods that learn from backpropagation are referred to as BP-CNN methods [34] which are used in SSL. Most promising CNN based semi-supervised network architectures are built using an ensemble system,

where multiple weak architectures are connected to obtain a stronger one [35]. Nevertheless, higher time and memory complexity is the burden of these ensemble systems.

SSL is mostly suitable for large datasets in which no sufficient information is available. SSL becomes essential when there exists a large unlabeled dataset. The challenge of SSL is to acquire higher accuracy when it is trained with a small amount of labeled dataset. Most semi-supervised and unsupervised learning architecture relates AE for data dimensionality reduction. In most cases, the AE is to be pre-trained using similar types of datasets to achieve a good performance. On the contrary, AE stores most of the information of data, which in some cases, may or may not be used as a feature. Instead of relying on AE, we propose a DNN architecture that performs data dimensionality reduction and can be used as an embedding system. The recent CDNN semi-supervised architectures rely on BP-CNN and FF-CNN architectures. But both of the architectures exhibit poor accuracy while they are trained with fewer data. Although ensemble-based CDNN architectures show better performance on complex datasets, they have higher time and memory complexity due to the overall fusion of multiple DNNs in the architecture. We solve the difficulty by implementing a better training method that requires a CDNN architecture for producing clusterable embedding points.

This paper refers to DNN as an embedding system on which clustering is applied. Although this architecture can be applied to most types of datasets, this paper specifically relates the AutoEmbedder to image-based datasets and most evaluations are performed on image-based datasets. The AutoEmbedder extracts features from higher dimensional data and compresses the features into a lower dimension embedding point, in which clustering is performed. The DCNN network performs the dimensionality reduction based on backpropagation distance loss calculation that is generated from Siamese network architecture.

## 3    Methodology

The proposed approach of this paper relates the dimension reduction technique of DNN networks [21] with a semi-supervised deep embedding system. Firstly, an embedding function is generated using DNN architecture, and it is trained using SNN architecture. Finally, the trained embedding function is used to transform higher dimensional data into meaningful low dimensional embedding points, on which clustering can be performed. The algorithm of the AutoEmbedder training process is presented in Algorithm 1. In the following subsections, we first derive the properties of the AutoEmbedder along with its training process. In the subsequent section 3.6, we emphasize the intuition and the proof of the overall training architecture of the AutoEmbedder. Finally, in the last subsection 3.7, we present a randomized training data selection scheme that boosts the accuracy of the AutoEmbedder.
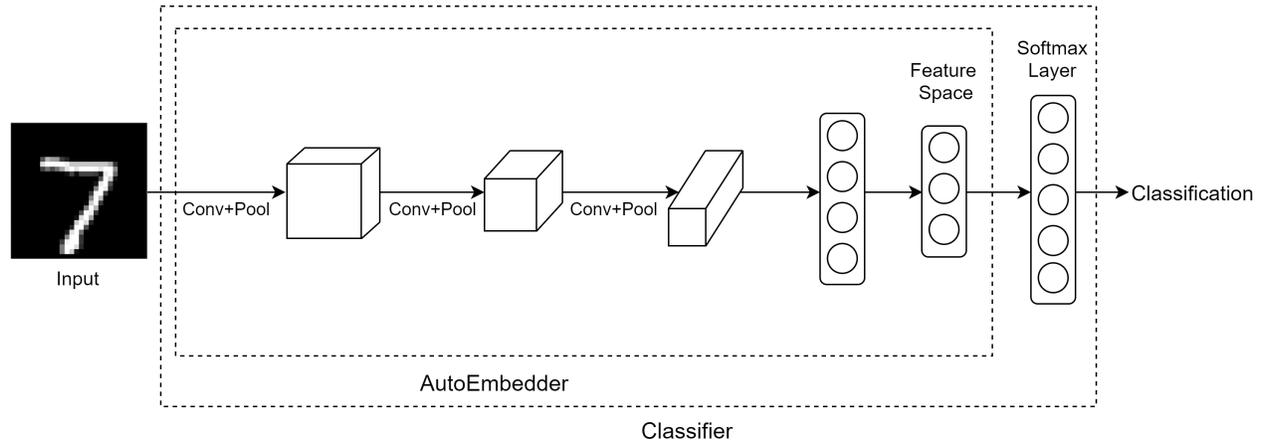


Figure 1: AutoEmbedder extraction from classifier CNN.

### 3.1    AutoEmbedder Architecture

A classifier neural network classifies inputs based on the activation node, which resides on the output layer. The output layer is also referred to as a softmax layer of a DNN architecture. In this paper, the previous layer of the last softmax layer of any classifier architecture is represented as the 'feature space' layer, which is illustrated in Figure 1. The feature space layer extracts final features from which the output layer performs the classification task. The feature space layer will work as an output layer of the AutoEmbedder after completing the AutoEmbedder training process. The number of nodes that reside in the feature space layer denotes the dimension of embedding points. The existing DNN architecture

is defined as AutoEmbedder. Mathematically, the AutoEmbedder can be represented as $f(x) = \mathbb{R}^k$, such that $k$ being the dimension of the embedding subspace. Like the autoencoder, the DNN architecture of the AutoEmbedder performs better in downsampling higher dimensional data, while keeping the required clusterable properties.
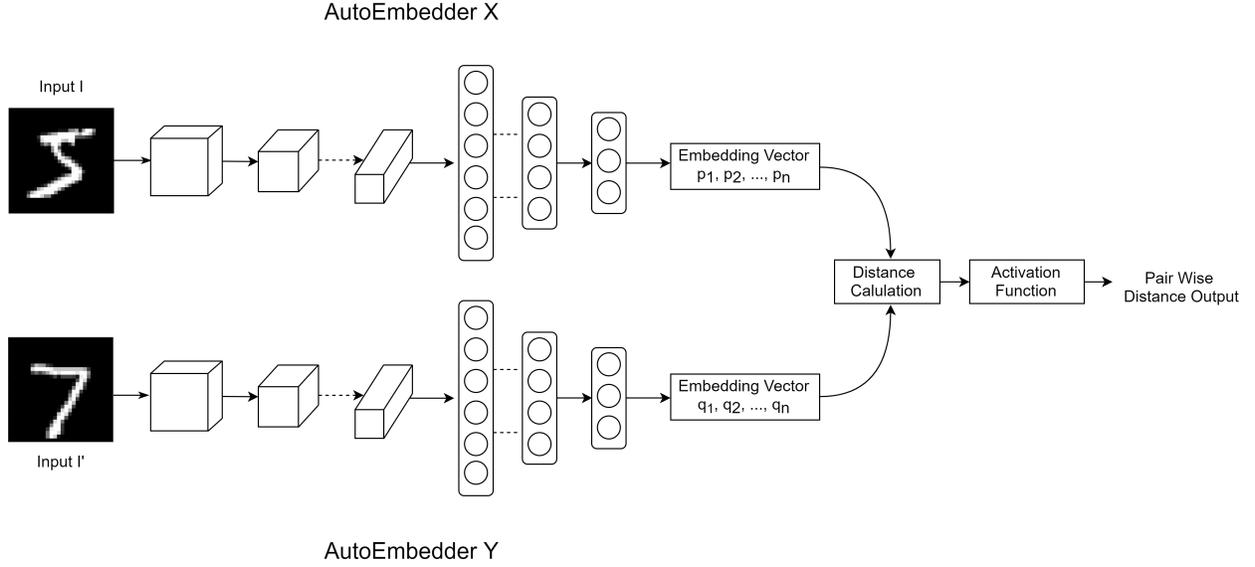


Figure 2: The SNN training architecture of AutoEmbedder.

## 3.2 AutoEmbedder Training Network

The training process of AutoEmbedder uses an SNN architecture as shown in Figure 2. In Siamese network architecture, a pair of AutoEmbedders with the same initial weights are placed parallelly. However, it is to be noted that the pair of networks does not share weights. The output of the AutoEmbedder pair is connected to a Euclidian distance calculation function as,

$$d(p, q) = d(q, p) = \sqrt{\sum_{i=1}^{k} (p_i - q_i)^2} \tag{1}$$

Where $f(x) = p = [p_1, p_2, ..., p_k]^T$ is the output of the first AutoEmbedder, and $g(x) = q = [q_1, q_2, ..., q_k]^T$ is the output of the second AutoEmbedder. The functions $f(x)$ and $g(x)$ are the AutoEmbedder pairs. The output of the distance calculation layer is further passed through a Rectified Linear Unit (ReLU) activation function with an upper bound value $\alpha$, defined as,

$$ReLU(x) = \begin{cases} x & \text{if } 0 \leq x < \alpha \\ \alpha & \text{if } x \geq \alpha \end{cases} \tag{2}$$

Where $x$ the input of the ReLU activation function and $\alpha$ is a hyperparameter set on the AutoEmbedder training phase. The hyperparameter $\alpha$ is also used in pairwise constraints. Due to the upper bound value set on the ReLU activation layer, the output of the SNN architecture will be in the range $[0, \alpha]$. The defined SNN architecture receives a pair of data and outputs the Euclidean distance of the embedding vector pairs. Combining equation 1 and 2, the overall training framework of the SNN can be represented as,

$$S(x_1, x_2) = ReLU(d(f(x_1), g(x_2))) = \mathbb{R}^+_{\leq \alpha} \tag{3}$$

Where $S(x_1, x_2)$ represents the SNN architecture function.

## 3.3 AutoEmbedder Pairwise Constraints

To train the SNN architecture of the AutoEmbedder with a precise target value, a distance hyperparameter $\alpha$ is to be decided. For any input data pair $x_i$ and $x_j$, the pairwise constraint is rated to be $\alpha$ if there exists a can-not-link constraint. Otherwise, the distance value is estimated to be 0. This can be mathematically stated as,

$$c_{ij} = \begin{cases} 0 & \text{if } x_i \text{ and } x_j \text{ must link} \\ \alpha & \text{if } x_i \text{ and } x_j \text{ can not link} \end{cases} \tag{4}$$

The pairwise constraints instruct the AutoEmbedder to create clusterable points. By equation 4, the AutoEmbedder pair is instructed to generate embedding vectors closer to zero if the input pair refer to the same class. Otherwise, it is instructed to generate embedding vectors greater or equal to $\alpha$, if the input pair refer to mixed classes.
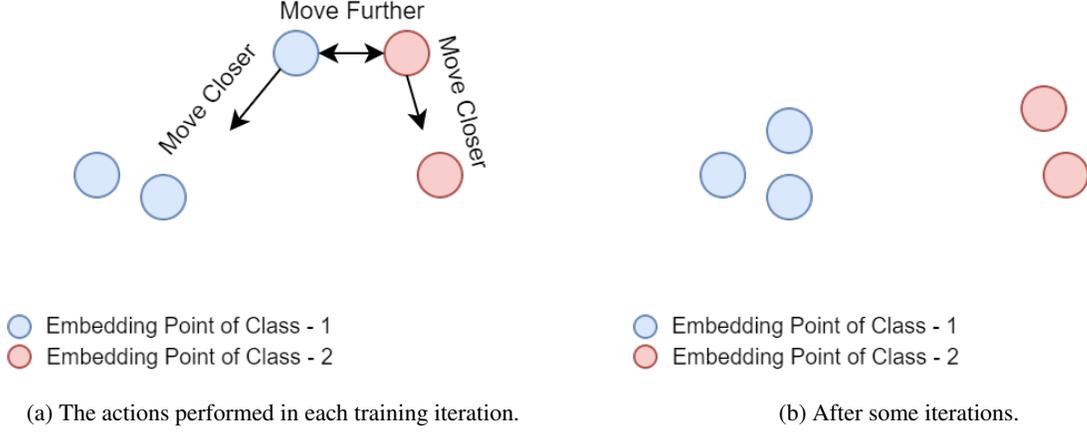


(a) The actions performed in each training iteration.    (b) After some iterations.

Figure 3: The AutoEmbedder training process moves can-not-link embedding pairs at a distance greater than or equal to $\alpha$ and moves embedding points of the must link pairs to a closer distance. Subfigures 3a and 3b illustrates the scenario.

## 3.4 AutoEmbedder Training

In each iteration of the AutoEmbedder training, the pair of AutoEmbedder is trained twice. Let the SNN architecture is trained though the function $S.train()$, which receives three parameters. The first two parameters are the inputs of the first and second AutoEmbedder, respectively, and the third parameter is the ground pairwise constraint output. At each iteration, the model is trained as,

$$S.train(I, I', Y)$$
$$S.train(I', I, Y) \tag{5}$$

Here $I$ and $I'$ define two subsets of data, and $Y$ defines the pairwise constraints. After completing the training of the AutoEmbedder pairs, any one of the two trained AutoEmbedders can be used from the SNN architecture. The backpropagation phase of each training iteration reduces the AutoEmbedder pairwise loss by moving the can-not-link pairs (mixed-class input pairs) at a distance and must-link pairs (same-class input pairs) closer as illustrated in Figure 3. This is the basic requirement of a data embedding for being clusterable.

## 3.5 AutoEmbedder Pairwise Loss Calculation

The output of the SNN architecture is a thresholded pairwise input distance, which is a continuous output in the range $[0, \alpha]$. Due to the criteria, the SNN architecture is trained based on regression. Most SNN based architectures often implement contrastive loss function [36]. However, due to the threshold of the final ReLU activation function derived in equation 2, the contrastive loss function may generate improper results. Table 1 illustrates a comparison of accuracy, while the AutoEmbedder is trained based on mean squared error (MSE) and contrastive loss. The training parameters used in the comparison are reported in Table 4. The comparison apprises that MSE is the most suitable for SNN training architecture. Hence, the pairwise loss is calculated using MSE.
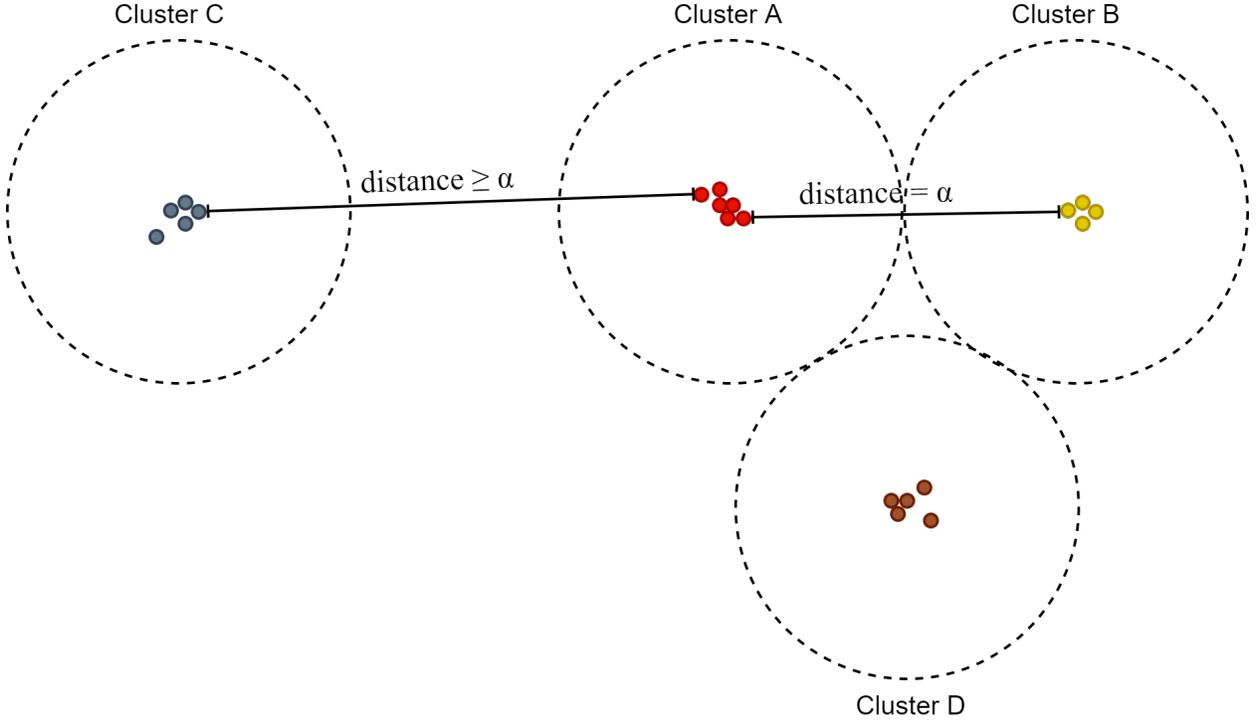
Let the pairwise ground truth vector be $Y$, and the predicted pairwise vector be $\hat{Y}$. The MSE for each iteration batch is calculated as,

$$MSE(Y, \hat{Y}) = \frac{1}{batchSize} \sum_{i=1}^{batchSize} (y_i - \hat{y}_i)^2 \tag{6}$$

Adam optimizer [37] and backpropagation is used to reduce the MSE.

Table 1: A comparison between mean square error loss and contrastive loss based on ACC, NMI, and ARI metrics.

| Dataset | Contrastive Loss ACC | MSE Loss ACC | Contrastive Loss NMI | MSE Loss NMI | Contrastive Loss ARI | MSE Loss ARI |
|---|---|---|---|---|---|---|
| MNIST | $13 \pm 4.1$ | $98.4 \pm 0.5$ | $0.95 \pm 0.02$ | $0.02 \pm 0.01$ | $0.01 \pm 0.01$ | $0.98 \pm 0.01$ |
| Fashion-MNIST | $11.07 \pm 1.2$ | $91.9 \pm 1.7$ | $0.0 \pm 0.03$ | $0.87 \pm 0.03$ | $0.0 \pm 0.03$ | $0.84 \pm 0.03$ |
| CIFAR10 | $4.7 \pm 2.1$ | $87.1 \pm 2.7$ | $0 \pm 0$ | $0.78 \pm 0.04$ | $0 \pm 0$ | $0.72 \pm 0.05$ |
| REUTERS | $6.7 \pm 1.2$ | $72.2 \pm 1.2$ | $0 \pm 0$ | $0.80 \pm 0.03$ | $0 \pm 0$ | $0.72 \pm 0.06$ |



Figure 4: An illustration of an optimal clusterable data points. Each pair of clusters tries to maintain a distance of $\alpha$, which is maintained by the SNN architecture through the pairwise constraints.

### 3.6 Proof of AutoEmbedder Training Architecture

The significance of SNN architecture is to train the AutoEmbedder so that, a) must link embedding vector pairs remains as close as possible, and b) the can-not-link embedding vector pairs obtain at least $\alpha$ distance from each other. If it is possible to construct a function that maintains the aforementioned properties, it can be concluded that the function generates clusterable embedding vectors. Also, the role of hyperparameter $\alpha$ is to construct a minimum margin or distance between the dissimilar data embeddings. Let us consider the two AutoEmbedders of the SNN architecture as two functions, $f(x)$ and $g(x)$. Let us also consider that $u$ and $v$ are two different data class sets, and they contain a can-not-link constraint for each other ($c_{uv} = \alpha$), and a must link constraint for themselves ($c_{uu} = 0, c_{vv} = 0$).

By considering the aforementioned cluster characteristics, it can be assumed that a function $f(x)$ generates clusterable embedding points if it maintains the following properties,

$$\|f(u_i) - f(u_j)\| \approx 0 \qquad [\forall_i \ u_i \in u, \forall_j \ u_j \in u] \tag{7}$$
$$and, \quad \|f(u_i) - f(v_j)\| \geq \alpha \qquad [\forall_i \ u_i \in u, \forall_j \ v_j \in v] \tag{8}$$

The abovementioned property must satisfy for function $g(x)$, since in subsection 3.4, we discussed any of the functions can be used as AutoEmbedder.

Through the training process of the SNN architecture, let us consider that both functions converge to an optimal state such that both functions optimally maintain the pairwise constraints. Therefore, the functions hold the following properties due to the must-link constraints,

$$f(u_i) \approx g(u_j) \tag{9}$$

Also, due to the can-not-link constraint, the functions hold the following properties,

$$\|f(v_i) - g(u_j)\| \geq \alpha \tag{10}$$

By placing the approximate value from equation 9,

$$\|f(v_i) - f(u_j)\| \geq \alpha$$
$$or, \quad \|f(u_i) - f(v_j)\| \geq \alpha \tag{11}$$

The above equation proves that the SNN architecture maintains the property described in equation, 8.

Furthermore, let us consider that, the function $f(x)$ produces embedding points for all input pairs $u_i$ and $w$, which are in the distance greater than zero. Consider $w$ as random data input. This can be stated as,

$$\|f(u_i) - f(w)\| > 0 \tag{12}$$

However, if we consider $w \in u$ and $w$ is equal to the $k^{th}$ input of the class $u$. Mathematically this can be formed as,

$$\|f(u_i) - f(u_k)\| > 0 \tag{13}$$

If we place relevant value from equation 9, we get,

$$\|f(u_i) - g(u_k)\| > 0 \tag{14}$$

The above equations 13 and 14 are contradictory to equation 9. Because, it is considered that both functions satisfies equation 9 and 10 after they are trained through SNN. As the distance can not be negative, we can reform equation 13 as,

$$\|f(u_i) - f(u_k)\| = 0$$
$$or, \quad \|f(u_i) - f(u_k)\| \approx 0 \tag{15}$$
$$similarly, \quad \|f(u_i) - f(u_j)\| \approx 0 \tag{16}$$

The above equation proves that SNN architecture maintains the property of equation 7. We can also construct a similar proof for function $g(x)$. Hence, it can be concluded that the SNN architecture can train two functions such that they generate clusterable embedding points based on the pairwise constraints.

From the general implementation of the AutoEmbedder training architecture, it can be understood that the hyperparameter $\alpha$ works as a cluster margin. The margin states that two inputs are considered to be in separate clusters if their distance is greater than or equal to the margin value. Maintaining the pairwise constraints, the AutoEmbedder can obtain an optimal knowledge to downsample higher dimensional data into clusterable points, as illustrated in Figure 4. Also, the threshold of the ReLU activation function (equation 2) serves to reduce the distance of the can-not-link pairs if they are farther than $\alpha$. This scenario is illustrated in Figure 4 for cluster pair A and C.

### 3.7 AutoEmbedder Random Train Data Selection

The AutoEmbedder is trained by selecting a fixed number of input pairs per training iteration. At each training iteration of the AutoEmbedder, two data-subset $I$ and $I'$ are created where, $|I| = |I'| = batchSize$, $I \in X$, and $I' \in X$. Here, $batchSize$ is the number of inputs in each iteration. The elements of the respective sets $I$ and $I'$ are randomly selected from the dataset $X$. The two data-subsets $I$ and $I'$ are used to train the AutoEmbedder at each iteration.

Let $P$ be the approximate probability of randomly selecting a data pair of same class, $s$ be the number of data of the same class, and $c$ be the total number of classes. Then $P$ can be defined as,

$$P = \binom{s}{1} \Big/ \binom{s*c}{1} \tag{17}$$

As $s < s*c$, the approximate probability of randomly selecting a pair of input data of different class $P'$ is greater than $P$. Due to the fact, the randomized selection of data pairs will contain more can-not-link constraints than must-link
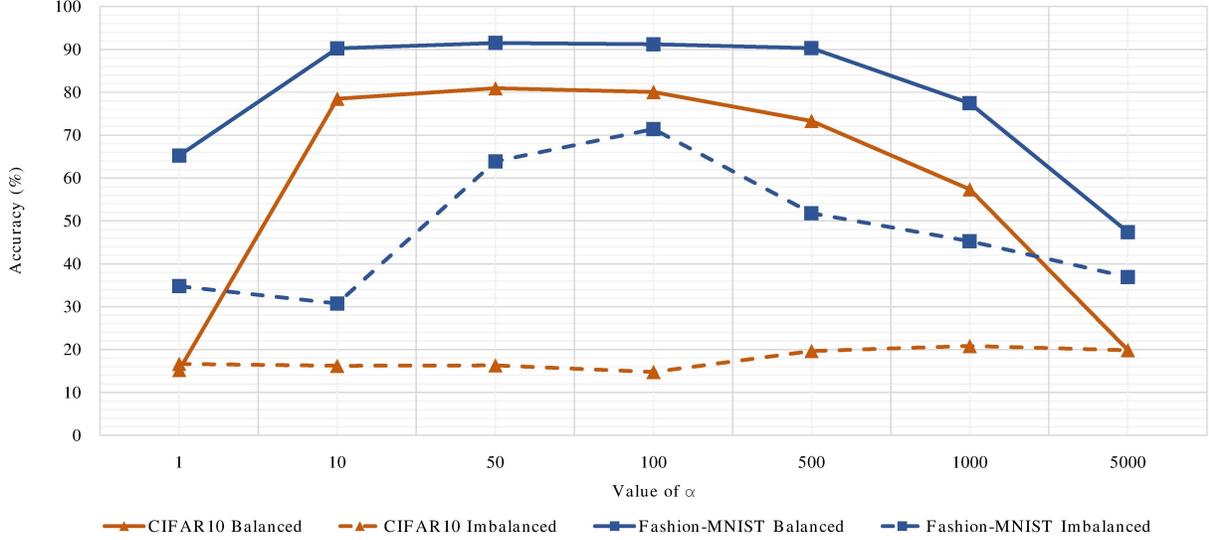
Figure 5: Accuracy measurement of two datasets by selecting different $\alpha$ values and containing balanced and imbalanced pairwise constraints while training. The illustrated result is obtained based on 300 epochs.

constraints. This leads to biased training data. If these link constraints are kept fully random or imbalanced, it is observed that the accuracy of AutoEmbedder is much lower after completing the training process of the AutoEmbedder, as depicted in Figure 5. Hence, for each iteration, half of the data pair is randomly selected, which contains must-link pairs, and the other half is randomly selected containing can-not-link pairs. Mathematically it can be viewed as,

$$nML = nCL \tag{18}$$

Where $nML$ is the number of pairs containing must-link constraint, $nCL$ is the number of pairs containing can-not-link constraint, and $nML + nCL = batchSize$. By maintaining this criterion, a massive accuracy improvement is observed that is illustrated in Figure 5.

## 4    Empirical Results

In this section, we describe evaluation metrics along with the experimental setup of the AutoEmbedder architecture, and the information of the datasets on which the tests were performed. Finally, we demonstrate the result of the AutoEmbedder.

### 4.1    Evaluation Metric

To evaluate the correctness of the proposed method, three well known and standard accuracy metrics are used. The evaluation metrics are presented below.

**Accuracy:** Accuracy (ACC) defines the unsupervised clustering accuracy, stated as,

$$ACC(c, c') = \left( \max_m \frac{\sum_{i=1}^{n} l\{c_i = m(c_i')\}}{2} \right) \times 100\% \tag{19}$$

Where $l_i$ defines the ground-truth label, $c_i$ defines the cluster assignment produced by the algorithm, and $m(.)$ ranges over all possible one-to-one mapping of the labels and clusters, from which the best mapping is taken. The mapping using the Hungarian algorithm is found to be efficient.

**Normalized Mutual Information:** The normalized mutual information (NMI) is defined as,

$$NMI(c, c') = \frac{I(c; c')}{\max(H(c), H(c'))} \tag{20}$$

9

---

**Algorithm 1:** AutoEmbedder Training Algorithm

---

**Input:** Subset of the dataset to be clustered $X$, AutoEmbedder model $M$, Number of iterations $e$, Training batch per
      iteration $batchSize$, Distance hyperparameter $\alpha$

**Result:** Trained AutoEmbedder

initialize two AutoEmbedder models $m_1$, $m_2$ with the same architecture and weight as $M$;

build a siamese network $S$ with AutoEmbedders $m_1$ and $m_2$;

$iter \leftarrow 0$;

**while** $iter < e$ **do**

    initialize two empty input data set $I$ and $I'$;

    initialize an empty target output set $Y$;

    $b \leftarrow 0$;

    **while** $b < \frac{batchSize}{2}$ **do**

        select two random data input $x_i$ and $x_j$ containing must link constraint $c_{ij} == 0$;

        append $x_i$ to $I$, $x_j$ to $I'$ and $c_{ij}$ to $Y$;

        $b = b + 1$;

    $b \leftarrow 0$;

    **while** $b < \frac{batchSize}{2}$ **do**

        select two random data input $x_i$ and $x_j$ containing must link constraint $c_{ij} == \alpha$;

        append $x_i$ to $I$, $x_j$ to $I'$ and $c_{ij}$ to $Y$;

        $b = b + 1$;

    train siamese network S with inputs $I$, $I'$ and $Y$;

    train siamese network S with inputs $I'$, $I$ and $Y$;

    $iter = iter + 1$;

---

Where $c$ is the ground truth and $c'$ is the predicted cluster. $I(.)$ refers to the mutual information between $c$ and $c'$. $H(.)$ denotes the entropy.

**Adjusted Rand Index:** The adjusted random index (ARI) is calculated using the contingency table [38]. The ARI can be derived as,

$$ARI = \frac{\sum_{ij}\binom{n_{ij}}{2} - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}\right] \Big/ \binom{n}{2}}{\frac{1}{2}\left[\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2}\right] - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}\right] \Big/ \binom{n}{2}} \tag{21}$$

Here, $n_{ij}$, $a_i$, and $b_j$ are the values from the contingency table.

Both NMI and ARI produce a result in the range $[0, 1]$ whereas, ACC produces a result in the scale $[0, 100]$. The higher value of these indices indicates a better correspondence between the cluster and the ground truth.

## 4.2 Experimental Setup

The neural network architecture is implemented using *Keras* [39]. To perform mathematical operations *numpy* is used [40]. The overall clustering is performed using K-means via *scikit-learn* library [41] with the default parameters. As an AutoEmbedder, *Keras* implemented MobileNet architecture is used. The parameters of the MobileNet architecture is defned in Table 2.

The 'input_shape' parameter of Table 2 is a variable that is based on the applied datasets. Also, the last layer of the MobileNet is 'conv_pw_13_relu' with a shape of $(1, 1, 1024)$. To reduce the dimension for each dataset, a feature space layer is added. The shape of the feature space layer is reported in column 'Embedding Dimensions' of Table 4. The number of nodes residing in the feature space layer denotes the embedding dimension of the AutoEmbedder. As the MobileNet can handle a minimum image shape of $(32, 32, 3)$, the datasets with lesser image shapes are zero-padded. To convert a single channel of a black and white image to a three-channel image, a copy of the black and white image is generated into each of the three channels. A graphical process of this image conversion is shown in Figure 6. The

Table 2: The MobileNet parameters set for AutoEmbedder.

| Argument | Value |
|---|---|
| alpha | 1 |
| depth_multiplier | 1 |
| dropout | 0.2 |
| include_top | False |
| pooling | None |
| input_tensor | None |
| weights | None |



Single Channel Image

Single Channel Zero Padded Image

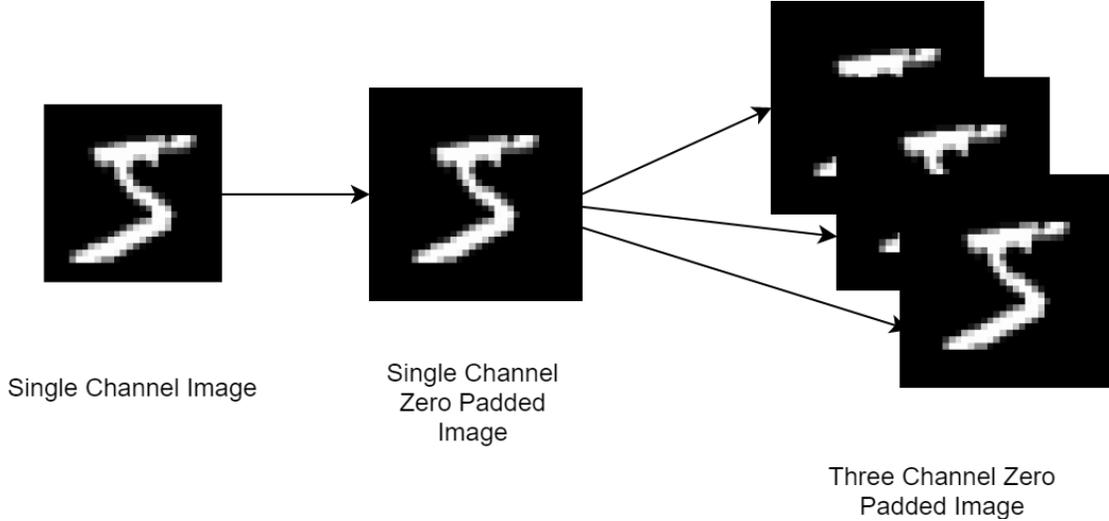Three Channel Zero Padded Image

Figure 6: Conversion process of MNIST $(28, 28, 1)$ input image shape to $(32, 32, 3)$.

converted dimension is reported in Table 3. The must-link and can-not link constraints are defined using the available data labels of the tested datasets.

The AutoEmbedder for REUTERS dataset contains four dense layers of nodes 512, 256, 128, and 64, respectively, with a default ReLU activation function. Keras implemented mean squared error and Adam optimization function is used to train the AutoEmbedder for all the datasets.

Table 3: Information on the datasets on which AutoEmbedder is evaluated. The reported image dimension is given by the format $(height, width, channel)$.

| Dataset | Classes | Size | Description | Input Dimension | Converted Input Dimension |
|---|---|---|---|---|---|
| MNIST | 10 | 60,000 | Handwritten Digits | (28, 28, 1) | (32, 32, 3) |
| Fashion-MNIST | 10 | 60,000 | Shoes/Clothing | (28, 28, 1) | (32, 32, 3) |
| CIFAR10 | 10 | 50,000 | Vehicle/Animals | (32, 32, 3) | (32, 32, 3) |
| COIL20 | 20 | 1,440 | Objects | (128, 128, 1) | (128, 128, 3) |
| SVHN | 10 | 99,289 | Street View House Number | (32, 32, 3) | (32, 32, 3) |
| REUTERS | 46 | 11,228 | Word Sequence | - | - |

## 4.3   Datasets

The AutoEmbedder is tested on well-known datasets. A short description of the datasets on which the AutoEmbedder is evaluated is presented in Table 3.

Table 4: The parameters used to train the AutoEmbedder on different datasets.

| Dataset | Embedding Dimension | Batch Size | Epochs | Distance ($\alpha$) |
|---|---|---|---|---|
| MNIST | 2 | 128 | 3,000 | 100 |
| Fashion-MNIST | 2 | 128 | 3,000 | 100 |
| CIFAR10 | 3 | 128 | 3,000 | 100 |
| COIL20 | 2 | 128 | 3,000 | 100 |
| SVHN | 3 | 128 | 3,000 | 100 |
| REUTERS | 16 | 128 | 3,000 | 100 |

Table 5: Accuracy, NMI, and ARI score comparison of different unsupervised and semi-supervised architectures tested on MNIST, Fashion-MNIST, CIFAR10, and COIL20 dataset. The semi-supervised architectures are marked with an asterisk (*). Methods that are unsuitable for calculating NMI and ARI are kept blank.
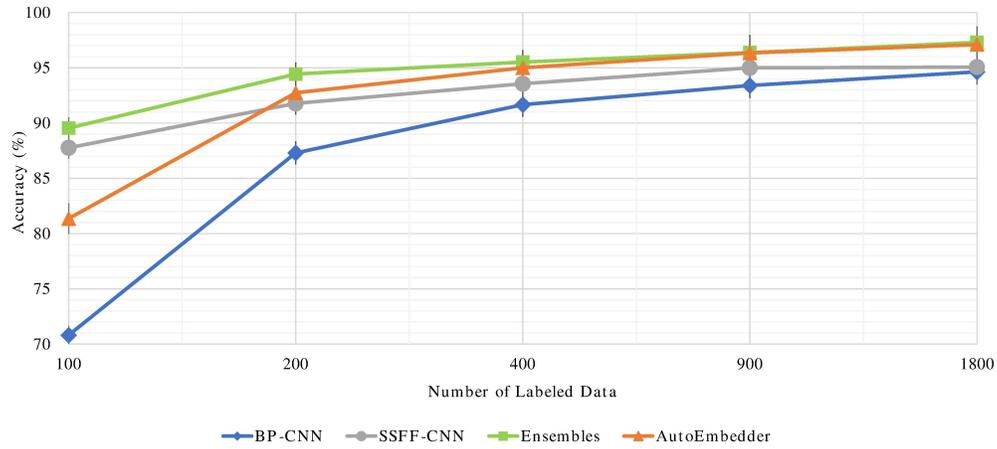
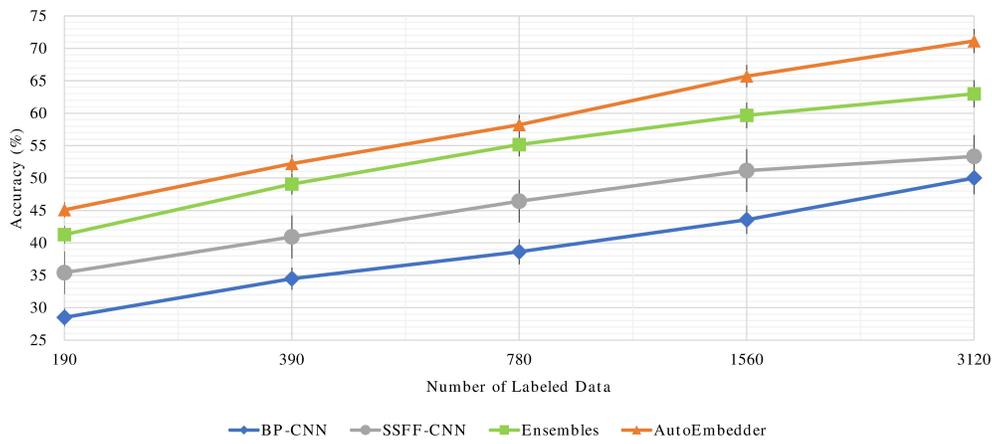| Model | MNIST | | | Fashion-MNIST | | | CIFAR10 | | | COIL20 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | ACC | NMI | ARI | ACC | NMI | ARI | ACC | NMI | ARI | ACC | NMI | ARI |
| ADC [42] | $98.2 \pm 0.6$ | - | - | $56.1 \pm 2.9$ | - | - | $28.1 \pm 2.9$ | - | - | $61.2 \pm 2.1$ | - | - |
| CatGAN [43] | $97 \pm 0.8$ | - | - | $87.49 \pm 1.4$ | - | - | $80.42 \pm 1.5$ | - | - | $92.38 \pm 0.7$ | - | - |
| ClusterGAN [45] | $92.4 \pm 1.1$ | $0.86 \pm 0.06$ | $0.89 \pm 0.04$ | $62 \pm 1.7$ | $0.65 \pm 0.07$ | $0.50 \pm 0.08$ | $41 \pm 2.1$ | $0.16 \pm 0.04$ | $0.30 \pm 0.02$ | $69 \pm 1.1$ | $0.68 \pm 0.07$ | $0.59 \pm 0.2$ |
| DAE Network [46] | $95.4 \pm 0.7$ | $0.93 \pm 0.04$ | $0.91 \pm 0.03$ | $63.8 \pm 1.4$ | $0.61 \pm 0.07$ | $0.53 \pm 0.04$ | $42.1 \pm 2.3$ | $0.28 \pm 0.04$ | $0.31 \pm 0.03$ | $69.4 \pm 1.3$ | $0.72 \pm 0.04$ | $0.59 \pm 0.3$ |
| DBC [47] | $94.7 \pm 1.1$ | $0.90 \pm 0.03$ | $0.93 \pm 0.02$ | $63.1 \pm 1.9$ | $0.65 \pm 0.06$ | $0.52 \pm 0.03$ | $34.2 \pm 1.4$ | $0.33 \pm 0.04$ | $0.22 \pm 0.03$ | $77.1 \pm 1.6$ | $0.87 \pm 0.05$ | $0.63 \pm 0.2$ |
| DEC-DA [48] | $95 \pm 1.6$ | $0.95 \pm 0.07$ | $0.94 \pm 0.02$ | $55.2 \pm 1.1$ | $0.64 \pm 0.04$ | $0.46 \pm 0.04$ | $18.9 \pm 3.1$ | $0.35 \pm 0.08$ | $0.07 \pm 0.03$ | $59.7 \pm 2.1$ | $0.82 \pm 0.08$ | $0.47 \pm 0.04$ |
| DEN [16] | $96.2 \pm 1.4$ | $0.91 \pm 0.03$ | $0.94 \pm 0.03$ | $78.5 \pm 2.1$ | $0.47 \pm 0.07$ | $0.67 \pm 0.03$ | $38.3 \pm 2.7$ | $0.22 \pm 0.04$ | $0.24 \pm 0.02$ | $86 \pm 1.8$ | $0.82 \pm 0.07$ | $0.71 \pm 0.2$ |
| DEPICT [18] | $95.5 \pm 0.6$ | $0.90 \pm 0.07$ | $0.93 \pm 0.02$ | $39.1 \pm 1.1$ | $0.37 \pm 0.04$ | $0.30 \pm 0.04$ | $43.7 \pm 1.3$ | $0.18 \pm 0.05$ | $0.31 \pm 0.03$ | $67.2 \pm 1.8$ | $0.43 \pm 0.03$ | $0.56 \pm 0.2$ |
| JULE [22] | $95.1 \pm 0.7$ | $0.91 \pm 0.02$ | $0.92 \pm 0.03$ | $54.1 \pm 1.9$ | $0.58 \pm 0.08$ | $0.39 \pm 0.06$ | $26.5 \pm 2$ | $0.34 \pm 0.07$ | $0.12 \pm 0.01$ | $89.8 \pm 1.2$ | $0.96 \pm 0.04$ | $0.88 \pm 0.05$ |
| RTM [49] | $94.8 \pm 0.9$ | $0.93 \pm 0.03$ | $0.92 \pm 0.02$ | $69.1 \pm 1.5$ | $0.68 \pm 0.6$ | $0.58 \pm 0.02$ | $28.5 \pm 2.3$ | $0.19 \pm 0.08$ | $0.13 \pm 0.03$ | $82 \pm 2.1$ | $0.70 \pm 0.09$ | $0.68 \pm 0.03$ |
| SpectralNet [10] | $96.1 \pm 0.8$ | $0.91 \pm 0.01$ | $0.95 \pm 0.02$ | $53.3 \pm 2.1$ | $0.53 \pm 0.04$ | $0.47 \pm 0.04$ | $20.8 \pm 2.5$ | $0.32 \pm 0.04$ | $0.10 \pm 0.02$ | $69 \pm 1.6$ | $0.62 \pm 0.07$ | $0.59 \pm 0.02$ |
| TAGnet [50] | $68 \pm 0.8$ | $0.64 \pm 0.09$ | $0.58 \pm 0.05$ | $58.7 \pm 1.4$ | $0.28 \pm 0.07$ | $0.46 \pm 0.02$ | $24.1 \pm 2.9$ | $0.12 \pm 0.02$ | $0.11 \pm 0.04$ | $89.9 \pm 1.1$ | $0.91 \pm 0.05$ | $0.72 \pm 0.04$ |
| VaDE [25] | $93.9 \pm 0.6$ | $0.88 \pm 0.06$ | $0.91 \pm 0.03$ | $56.1 \pm 1.7$ | $0.47 \pm 0.02$ | $0.42 \pm 0.03$ | $31.7 \pm 2.1$ | $0.27 \pm 0.07$ | $0.21 \pm 0.03$ | $88 \pm 1.4$ | $0.67 \pm 0.05$ | $0.72 \pm 0.03$ |
| SS-KM* [44] | $52.6 \pm 3.1$ | $0.48 \pm 0.09$ | $0.42 \pm 0.08$ | $22.3 \pm 1.4$ | $0.17 \pm 0.04$ | $0.16 \pm 0.04$ | $20.7 \pm 2.4$ | $0.14 \pm 0.08$ | $0.11 \pm 0.09$ | $21.5 \pm 1.7$ | $0.17 \pm 0.5$ | $0.18 \pm 0.06$ |
| SS-DEC* [12] | $84.7 \pm 0.7$ | $0.83 \pm 0.04$ | $0.81 \pm 0.03$ | $44.2 \pm 1.1$ | $0.39 \pm 0.05$ | $0.41 \pm 0.03$ | $26.2 \pm 1.6$ | $0.19 \pm 0.06$ | $0.20 \pm 0.09$ | $78.4 \pm 1.3$ | $0.72 \pm 0.07$ | $0.71 \pm 0.08$ |
| CatGAN* [43] | $98.4 \pm 0.3$ | - | - | $89.2 \pm 1.2$ | - | - | $86.4 \pm 1.7$ | - | - | $98 \pm 1.3$ | - | - |
| **AutoEmbedder*** | $\mathbf{98.4 \pm 0.5}$ | $\mathbf{0.95 \pm 0.02}$ | $\mathbf{0.98 \pm 0.01}$ | $\mathbf{91.9 \pm 1.7}$ | $\mathbf{0.87 \pm 0.03}$ | $\mathbf{0.84 \pm 0.03}$ | $\mathbf{87.1 \pm 2.7}$ | $\mathbf{0.78 \pm 0.04}$ | $\mathbf{0.72 \pm 0.05}$ | $\mathbf{100 \pm 0}$ | $\mathbf{0.98 \pm 0.01}$ | $\mathbf{0.89 \pm 0.02}$ |

## 4.4 Results

The test results are obtained by calculating mean of the maximum ACC, NMI, and ARI scores for six runs. The results are reported in the mean±std format. The batch size, the epochs, the embedding dimension used for each dataset, and the distance hyperparameter $\alpha$ is reported in Table 4. The training iteration ($e$) stated in Algorithm 1 can be calculated as, $e = \frac{Epochs \times DatasetSize}{BatchSize}$. The AutoEmbedder is compared with both unsupervised and semi-supervised methods. All of the comparisons were performed in the same environment, and the other methods were implemented by maintaining the optimal hyperparameters.

Table 5 illustrates a comparison based on the ACC, NMI, and ARI scores tested over different image datasets. The table contains both unsupervised and semi-supervised methods, whereas the semi-supervised methods are marked with an asterisk (*). Furthermore, the highest scores are marked bold. In this comparison, some implemented methods generate pseudo labeling instead of generating embeddings [42, 43]. Therefore it is not possible to calculate NMI and ARI scores for these methods, and they are kept blank. From the comparison of Table 5, it can be observed that some semi-supervised methods perform less accurately than some unsupervised methods [44, 12]. However, it can also be witnessed that unsupervised methods fail to generate better accuracy in complex datasets, such as Fashion-MNIST and CIFAR-10. Yet, GAN based semi-supervised and unsupervised methods [43, 45] perform better than most other methods. However, GAN based methods suffer from some unfortunate circumstances due to the unstable learning of generators and discriminators [43]. On the contrary, the AutoEmbedder does not suffer from any instability and outperforms all of the other implemented semi-supervised and unsupervised methods.
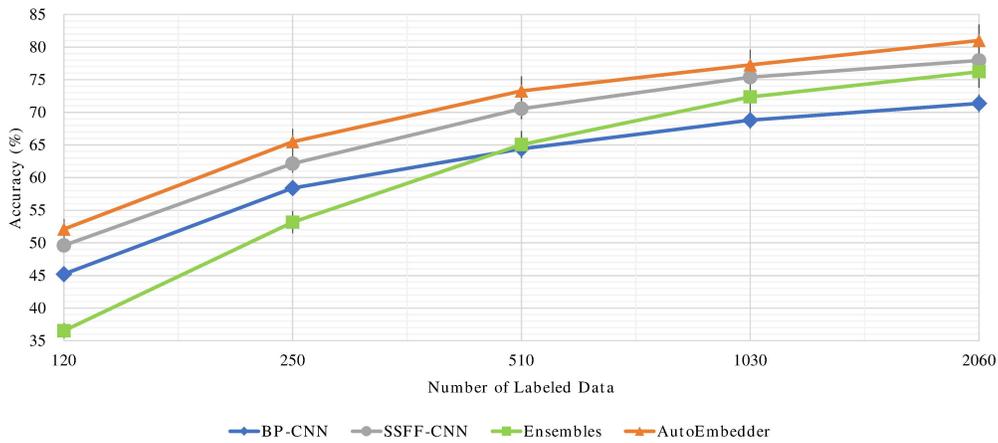
The AutoEmbedder architecture also outperforms other unsupervised methods that are applied to textual data. It is tested on the REUTERS dataset, and the evaluation report is presented in Table 6.

(a) Accuracy on MNIST dataset.



(b) Accuracy on CIFAR10 dataset.



(c) Accuracy on SVHN dataset.

Figure 7: A comparison of the AutoEmbedder with different label-based semi-supervised classifiers. Subfigure 7a, 7b, and 7c illustrate comparisons on MNIST, CIFAR10, and SVHN dataset, respectively. The accuracy of AutoEmbedder is reported using the ACC metric of the equation 19 and the accuracy of the label-based semi-supervised methods reported based on the percentage of positive prediction by total predictions.

Table 6: ACC, NMI, and ARI score of different unsupervised architectures tested on the REUTERS dataset.

| Model | ACC | NMI | ARI |
|---|---|---|---|
| DEC [9] | $51.4 \pm 1.9$ | $0.68 \pm 0.04$ | $0.49 \pm 0.09$ |
| AP [51] | $1.8 \pm 0.3$ | $0.27 \pm 0.07$ | $0.03 \pm 0.04$ |
| HDB [52] | $19.7 \pm 1.4$ | $0.20 \pm 0.04$ | $0.15 \pm 0.07$ |
| RCC [53] | $1.9 \pm 0.6$ | $0.32 \pm 0.03$ | $0.08 \pm 0.04$ |
| RCC-DR [53] | $1.9 \pm 0.7$ | $0.31 \pm 0.02$ | $0.08 \pm 0.05$ |
| **AutoEmbedder** | $\mathbf{72.2 \pm 1.2}$ | $\mathbf{0.80 \pm 0.03}$ | $\mathbf{0.72 \pm 0.06}$ |

We also compare the AutoEmbedder architecture with other label-based semi-supervised architectures [33, 34]. The label-based semi-supervised architectures are trained on a small portion of labeled data. Furthermore, as these DCNN methods are trained on labeled data, they contain activation functions in the last layer. Therefore the final outputs of these architectures are class labels, instead of embeddings. Although the training criteria of these methods and the AutoEmbedder is different, we perform a comparison among the architectures. Also, as the final labels of the label-based DCNN methods are not pseudo labels, the accuracy of these methods is calculated based on the ratio of the percentage of positive predictions by total predictions. Figure 7 presents comparisons based on MNIST, CIFAR-10, and SVHN datasets. In the comparisons, the AutoEmbedder is compared with label-based semi-supervised FF-CNN, BP-CNN, and ensemble architectures [33, 34]. The methods are compared based on different numbers of known labels of the dataset. The label-based semi-supervised methods are trained based on the known labels, whereas the AutoEmbedder is trained by constructing pairwise constraints from the known labels. The AutoEmbedder outperforms all the label-based DCNN architectures on CIFAR10 and SVHN datasets. However, it fails to produce the best results in the simple MNIST dataset, since MobileNet is not a suitable architecture for single-channel image datasets. The label-based DCNN architectures perform less accurately because the final layer activation functions are not optimized to generate optimal hyperplanes, while they are trained on fewer data. On the contrary, the AutoEmbedder generates embeddings based on the distance hyperparameter $\alpha$, which promises to generate low-dimensional clusterable points.

## 5    Conclusion

This paper introduces an embedding architecture, AutoEmbedder that produces meaningful clusterable embedding points. The end-to-end training process of the architecture is semi-supervised and requires a pairwise cluster linking information in the training phase. The training procedure does not include any clustering loss measures, instead, it uses Euclidean distance loss that is minimized by backpropagation. The AutoEmbedder only produces clusterable embedding points. The AutoEmbedder can be built based on any classification architecture with the required embedding dimension. From the benchmarks of this paper, it is to report that the AutoEmbedder presents better results on almost all the datasets. The embedding system constructs three-dimension embedding points from complex three-channel image datasets CIFAR10 along with SVHN and still produces better results. From the statistics of the empirical results, it may be concluded that the proposed method is beneficial to perform semi-supervised learning. We strongly believe that the overall contribution of this paper inaugurates a wider perception in the scope of embedding systems, semi-supervised learning, and image clustering research works.

## References

[1] T. N. Pappas and N. S. Jayant. An adaptive clustering algorithm for image segmentation. In *International Conference on Acoustics, Speech, and Signal Processing,*, pages 1667–1670 vol.3, May 1989.

[2] F. Hoeppner. Fuzzy shell clustering algorithms in image processing: fuzzy c-rectangular and 2-rectangular shells. *IEEE Transactions on Fuzzy Systems*, 5(4):599–613, Nov 1997.

[3] Yizhou Sun, Jiawei Han, Peixiang Zhao, Zhijun Yin, Hong Cheng, and Tianyi Wu. Rankclus: Integrating clustering with ranking for heterogeneous information network analysis. In *Proceedings of the 12th International Conference on Extending Database Technology: Advances in Database Technology*, EDBT '09, page 565–576, New York, NY, USA, 2009. Association for Computing Machinery.

[4] Shudong Huang, Zenglin Xu, and Jiancheng Lv. Adaptive local structure learning for document co-clustering. *Knowledge-Based Systems*, 148:74–84, 2018.

[5] Han Kyul Kim, Hyun joong Kim, and Sungzoon Cho. Bag-of-concepts: Comprehending document representation through clustering words in distributed representation. *Neurocomputing*, 266:336–352, 2017.

[6] Xiaofei He, Deng Cai, and Partha Niyogi. Laplacian score for feature selection. In *Proceedings of the 18th International Conference on Neural Information Processing Systems*, NIPS'05, page 507–514, Cambridge, MA, USA, 2005. MIT Press.

[7] Yazhou Ren, Guoji Zhang, Guoxian Yu, and Xuan Li. Local and global structure preserving based feature selection. *Neurocomputing*, 89:147 – 157, 2012.

[8] Fernando De la Torre and Takeo Kanade. Discriminative cluster analysis. In *Proceedings of the 23rd International Conference on Machine Learning*, ICML '06, page 241–248, New York, NY, USA, 2006. Association for Computing Machinery.

[9] Junyuan Xie, Ross Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering analysis, 2015.

[10] Uri Shaham, Kelly Stanton, Henry Li, Boaz Nadler, Ronen Basri, and Yuval Kluger. Spectralnet: Spectral clustering using deep neural networks, 2018.

[11] Ioannis Athanasiadis, Panagiotis Mousouliotis, and Loukas Petrou. A framework of transfer learning in object detection for embedded systems, 2018.

[12] Yazhou Ren, Kangrong Hu, Xinyi Dai, Lili Pan, Steven C.H. Hoi, and Zenglin Xu. Semi-supervised deep embedded clustering. *Neurocomputing*, 325:121 – 130, 2019.

[13] Chris Ding and Xiaofeng He. K-means clustering via principal component analysis. In *Proceedings of the Twenty-First International Conference on Machine Learning*, ICML '04, page 29, New York, NY, USA, 2004. Association for Computing Machinery.

[14] E. Min, X. Guo, Q. Liu, G. Zhang, J. Cui, and J. Long. A survey of clustering with deep learning: From the perspective of network architecture. *IEEE Access*, 6:39501–39514, 2018.

[15] Bo Yang, Xiao Fu, Nicholas D. Sidiropoulos, and Mingyi Hong. Towards k-means-friendly spaces: Simultaneous deep learning and clustering, 2016.

[16] P. Huang, Y. Huang, W. Wang, and L. Wang. Deep embedding network for clustering. In *2014 22nd International Conference on Pattern Recognition*, pages 1532–1537, Aug 2014.

[17] Sohil Atul Shah and Vladlen Koltun. Deep continuous clustering, 2018.

[18] Kamran Ghasedi Dizaji, Amirhossein Herandi, Cheng Deng, Weidong Cai, and Heng Huang. Deep clustering via joint convolutional autoencoder embedding and relative entropy minimization, 2017.

[19] Dongdong Chen, Jiancheng Lv, and Yi Zhang. Unsupervised multi-manifold clustering by learning deep representation, 2017.

[20] Pan Ji, Tong Zhang, Hongdong Li, Mathieu Salzmann, and Ian Reid. Deep subspace clustering networks. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 24–33. Curran Associates, Inc., 2017.

[21] Joris Guérin, Olivier Gibaru, Stéphane Thiery, and Eric Nyiri. Cnn features are also great at unsupervised classification, 2017.

[22] Jianwei Yang, Devi Parikh, and Dhruv Batra. Joint unsupervised learning of deep representations and image clusters. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.

[23] Jianlong Chang, Lingfeng Wang, Gaofeng Meng, Shiming Xiang, and Chunhong Pan. Deep adaptive image clustering. *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 5880–5888, 2017.

[24] Nat Dilokthanakul, Pedro A. M. Mediano, Marta Garnelo, Matthew C. H. Lee, Hugh Salimbeni, Kai Arulkumaran, and Murray Shanahan. Deep unsupervised clustering with gaussian mixture variational autoencoders, 2016.

[25] Zhuxi Jiang, Yin Zheng, Huachun Tan, Bangsheng Tang, and Hanning Zhou. Variational deep embedding: An unsupervised and generative approach to clustering, 2016.

[26] Eric Bair. Semi-supervised clustering methods. *Wiley Interdisciplinary Reviews: Computational Statistics*, 5(5):349–361, Jul 2013.

[27] Mikhail Bilenko, Sugato Basu, and Raymond J. Mooney. Integrating constraints and metric learning in semi-supervised clustering. In *Proceedings of the Twenty-First International Conference on Machine Learning*, ICML '04, page 11, New York, NY, USA, 2004. Association for Computing Machinery.

[28] Kiri Wagstaff, Claire Cardie, Seth Rogers, and Stefan Schroedl. Constrained k-means clustering with background knowledge. In *In ICML*, pages 577–584. Morgan Kaufmann, 2001.

[29] Dan Klein, Sepandar D. Kamvar, and Christopher D. Manning. From instance-level constraints to space-level constraints: Making the most of prior knowledge in data clustering. Technical Report 2002-10, Stanford InfoLab, February 2002.

[30] Yazhou Ren, Xiaohui Hu, Ke Shi, Guoxian Yu, Dezhong Yao, and Zenglin Xu. Semi-supervised denpeak clustering with pairwise constraints. In Xin Geng and Byeong-Ho Kang, editors, *PRICAI 2018: Trends in Artificial Intelligence*, pages 837–850, Cham, 2018. Springer International Publishing.

[31] Gang Chen. Deep transductive semi-supervised maximum margin clustering, 2015.

[32] Chuanqi Tan, Fuchun Sun, Tao Kong, Wenchang Zhang, Chao Yang, and Chunfang Liu. A survey on deep transfer learning, 2018.

[33] Yueru Chen, Yijing Yang, Min Zhang, and C. C. Jay Kuo. Semi-supervised learning via feedforward-designed convolutional neural networks, 2019.

[34] Quanshi Zhang, Ying Nian Wu, and Song-Chun Zhu. Interpretable convolutional neural networks, 2017.

[35] Cha Zhang and Yunqian Ma, editors. *Ensemble Machine Learning*. Springer US, 2012.

[36] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 1735–1742. IEEE, 2006.

[37] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2014.

[38] Jorge M Santos and Mark Embrechts. On the use of the adjusted rand index as a metric for evaluating supervised classification. In *International conference on artificial neural networks*, pages 175–184. Springer, 2009.

[39] François Chollet et al. Keras. https://keras.io, 2015.

[40] Stefan Van Der Walt, S. Chris Colbert, and Gaël Varoquaux. The numpy array: a structure for efficient numerical computation, 2011.

[41] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

[42] Philip Häusser, Johannes Plapp, Vladimir Golkov, Elie Aljalbout, and Daniel Cremers. Associative deep clustering: Training a classification network with no labels. In *GCPR*, 2017.

[43] Jost Tobias Springenberg. Unsupervised and semi-supervised learning with categorical generative adversarial networks, 2015.

[44] Sugato Basu, Arindam Banerjee, and Raymond J Mooney. Active semi-supervision for pairwise constrained clustering. In *Proceedings of the 2004 SIAM international conference on data mining*, pages 333–344. SIAM, 2004.

[45] Sudipto Mukherjee, Himanshu Asnani, Eugene Lin, and Sreeram Kannan. ClusterGAN: Latent space clustering in generative adversarial networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33:4610–4617, July 2019.

[46] Xu Yang, Cheng Deng, Feng Zheng, Junchi Yan, and Wei Liu. Deep spectral clustering using dual autoencoder network, 2019.

[47] Fengfu Li, Hong Qiao, Bo Zhang, and Xuanyang Xi. Discriminatively boosted image clustering with fully convolutional auto-encoders, 2017.

[48] Xifeng Guo, En Zhu, Xinwang Liu, and Jianping Yin. Deep embedded clustering with data augmentation. In Jun Zhu and Ichiro Takeuchi, editors, *Proceedings of The 10th Asian Conference on Machine Learning*, volume 95 of *Proceedings of Machine Learning Research*, pages 550–565. PMLR, 14–16 Nov 2018.

[49] Oliver Nina, Jamison Moody, and Clarissa Milligan. A decoder-free approach for unsupervised clustering and manifold learning with random triplet mining. In *The IEEE International Conference on Computer Vision (ICCV) Workshops*, Oct 2019.

[50] Zhangyang Wang, Shiyu Chang, Jiayu Zhou, Meng Wang, and Thomas S Huang. Learning a task-specific deep architecture for clustering. In Sanjay Chawla Venkatasubramanian and Wagner Meira, editors, *16th SIAM International Conference on Data Mining 2016, SDM 2016*, 16th SIAM International Conference on Data Mining 2016, SDM 2016, pages 369–377, United States, 1 2016. Society for Industrial and Applied Mathematics Publications.

[51] B. J. Frey and D. Dueck. Clustering by passing messages between data points. *Science*, 315(5814):972–976, February 2007.

[52] Ricardo J. G. B. Campello, Davoud Moulavi, and Joerg Sander. Density-based clustering based on hierarchical density estimates. In *Advances in Knowledge Discovery and Data Mining*, pages 160–172. Springer Berlin Heidelberg, 2013.

[53] Sohil Atul Shah and Vladlen Koltun. Robust continuous clustering. *Proceedings of the National Academy of Sciences*, 114(37):9814–9819, August 2017.