

# Learning Discrete Representations via Information Maximizing Self Augmented Training

Weihua Hu<sup>1</sup> Takeru Miyato<sup>2,3</sup> Seiya Tokui<sup>2,1</sup> Eiichi Matsumoto<sup>2,1</sup> Masashi Sugiyama<sup>4,1</sup>

## Abstract

Learning discrete representations of data is a central machine learning task because of the compactness of the representations and ease of interpretation. The task includes clustering and hash learning as special cases. Deep neural networks are promising to be used because they can model the non-linearity of data and scale to large datasets. However, their model complexity is huge, and therefore, we need to carefully regularize the networks in order to learn useful representations that exhibit intended invariance for applications of interest. To this end, we propose a method called Information Maximizing Self Augmented Training (IMSAT). In IMSAT, we use data augmentation to impose the invariance on discrete representations. More specifically, we encourage the predicted representations of augmented data points to be close to those of the original data points in an end-to-end fashion. At the same time, we maximize the information-theoretic dependency between data and their predicted discrete representations. Extensive experiments on benchmark datasets show that IMSAT produces state-of-the-art results for both clustering and unsupervised hash learning.

## 1. Introduction

The task of unsupervised discrete representation learning is to obtain a function that maps *similar* data into similar discrete representations, where the *similarity* of data is defined according to applications of interest. It is a central machine learning task because of the compactness of the representations and ease of interpretation. The task includes two important machine learning tasks as special cases: clustering and unsupervised hash learning. Clustering is widely applied to data-driven application domains (Berkhin, 2006),

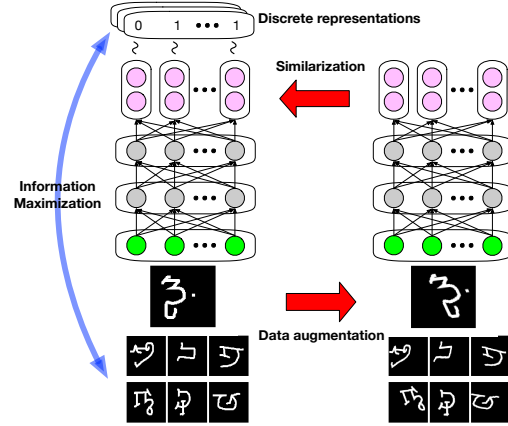


Figure 1. Basic idea of our proposed method for unsupervised discrete representation learning. We encourage the prediction of a neural network to remain unchanged under data augmentation (Red arrows), while maximizing the information-theoretic dependency between data and their representations (Blue arrow).

while hash learning is popular for an approximate nearest neighbor search for large scale information retrieval (Wang et al., 2016).

Various methods have been proposed for clustering and hash learning. The representative ones include  $K$ -means clustering and hashing (He et al., 2013), Gaussian mixture model clustering, iterative quantization (Gong et al., 2013), and minimal-loss hashing (Norouzi & Blei, 2011). However, these methods can only model linear boundaries between different representations; thus, cannot fit to non-linear structures of data. Kernel-based (Xu et al., 2004; Kulis & Darrell, 2009) and spectral (Ng et al., 2001; Weiss et al., 2009) methods can model the non-linearity of data, but they are difficult to scale to large datasets.

Deep neural networks are suitable candidates to overcome these limitations thanks to their scalability and flexibility of representing complicated, non-linear decision boundaries. However, their model complexity is huge, and therefore, regularization of the networks is crucial to learn good representations of data. Particularly, in the *unsupervised* representation learning scenario, target representations are not provided and hence are unconstrained. Therefore, we need

<sup>1</sup>University of Tokyo, Japan <sup>2</sup>Preferred Networks, Inc., Japan  
<sup>3</sup>ATR Cognitive Mechanism Laboratories, Japan <sup>4</sup>RIKEN, Japan.  
 Correspondence to: Weihua Hu <hu@ms.k.u-tokyo.ac.jp>.

to carefully regularize the networks in order to learn useful representations that exhibit intended invariance for applications of interest. Naïve regularizations to use are a weight decay and L1 regularization (Erin Liong et al., 2015). Such regularizations, however, encourage global smoothness of the function prediction; thus, may not necessarily impose the intended invariance on the predicted discrete representations.

Instead, in this paper, we use data augmentation to model the invariance of learned data representations. More specifically, we map data points into their discrete representations by a deep neural network and regularize it by encouraging its prediction to be invariant to data augmentation. The predicted discrete representations then exhibit the invariance specified by the augmentation. Our proposed regularization method is illustrated as red arrows in Figure 1. As depicted, we encourage the predicted representations of augmented data points to be close to those of the original data points in an end-to-end fashion. We term such regularization *self augmented training* (SAT). SAT is inspired by the recent success in regularization of neural networks in semi-supervised learning (Bachman et al., 2014; Miyato et al., 2016; Sajjadi et al., 2016). SAT is flexible to impose various types of invariances on the representations predicted by neural networks. For example, it is generally preferred for data representations to be locally invariant, i.e., remain unchanged under local perturbations on data points. Using SAT, we can impose the local invariance on the representations by pushing the predictions of perturbed data points to be close to those of the original data points. For image data, it may also be preferred for data representations to be invariant under affine distortion, e.g., rotation, scaling and parallel movement. We can similarly impose the invariance via SAT by using the affine distortion for the data augmentation.

We then combine the SAT with the regularized information maximization (RIM) for clustering (Gomes et al., 2010), and arrive at our Information Maximizing Self Augmented Training (IMSAT), an information-theoretic method for learning discrete representations using deep neural networks. We illustrate the basic idea of IMSAT in Figure 1. Following the RIM, we maximize information theoretic dependency between inputs and their mapped outputs, while regularizing the mapping function. IMSAT, however, differs from the original RIM in two ways. First, IMSAT deals with a more general setting of learning discrete representations; thus, is also applicable to hash learning. Second, it uses a deep neural network for the mapping function and regularizes it in an end-to-end fashion via SAT. Learning with our method can be performed by stochastic gradient descent (SGD); thus, scales well to large datasets.

In summary, our contributions are: (1) an information-theoretic method for unsupervised discrete representation

learning using deep neural networks with the end-to-end regularization, and (2) adaptations of the method to clustering and hash learning to achieve the state-of-the-art performance on several benchmark datasets.

The rest of the paper is organized as follows. Related work is summarized in Section 2, while our method, IMSAT, is presented in Section 3. Experiments on benchmark datasets are shown in Section 4 and conclusions are drawn in Section 5.

## 2. Related work

Conventional methods for clustering and hash learning can be broadly divided into two categories according to the decision boundaries they can model: linear (He et al., 2013; Gong et al., 2013; Norouzi & Blei, 2011) and non-linear (Xu et al., 2004; Kulis & Darrell, 2009; Ng et al., 2001; Weiss et al., 2009). The linear methods often scale well but cannot model non-linear boundaries between different clusters or hash codes. The non-linear methods are flexible in modeling the decision boundaries but scale poorly to large datasets.

Recently, clustering and hash learning using deep neural networks have attracted much attention due to their scalability and flexibility in representing non-linearity and hierarchical compositionality of data. In clustering, Xie et al. (2016) proposed to use deep neural networks to simultaneously learn feature representations and cluster assignments, while Dilokthanakul et al. (2016) and Zheng et al. (2016) proposed to model the data generation process by using deep generative models with Gaussian mixture models as prior distributions. Regarding hashing learning, a number of studies have used deep neural networks for supervised hash learning and achieved state-of-the-art results on image and text retrievals (Xia et al., 2014; Lai et al., 2015; Zhang et al., 2015; Xu et al., 2015; Li et al., 2015). Relatively few studies have focused on unsupervised hash learning using deep neural networks. The pioneering work is semantic hashing, which uses stacked RBM models to learn compact binary representations (Salakhutdinov & Hinton, 2009). Erin Liong et al. (2015) recently proposed to use deep neural networks for the mapping function and achieved state-of-the-art results. These unsupervised methods, however, did not explicitly impose the intended invariance on the data representations. Consequently, the predicted representations may not necessarily be intended for applications of interest.

In supervised and semi-supervised learning scenarios, data augmentation has been widely used to regularize neural networks. Leen (1995) showed that applying data augmentation to a supervised learning problem is equivalent to adding a regularization to the original cost function. Bachman et al. (2014); Miyato et al. (2016); Sajjadi et al. (2016)

showed that such regularization can be adapted to a semi-supervised learning setting to achieve state-of-the-art performance. Our study suggests that these regularization is also crucial in an *unsupervised* representation learning scenario. Relation of our work to denoising and contractive auto-encoders (Vincent et al., 2008; Rifai et al., 2011) is discussed in Appendix A.

### 3. Method

Let  $\mathcal{X}$  and  $\mathcal{Y}$  denote the domains of inputs and discrete representations, respectively. Given training samples,  $\{x_1, x_2, \dots, x_N\}$ , the task of discrete representation learning is to obtain a function,  $f : \mathcal{X} \rightarrow \mathcal{Y}$ , that maps similar inputs into similar discrete representations. The similarity of data is defined according to applications of interest.

We organize Section 3 as follows. In Section 3.1, we review the RIM for clustering (Gomes et al., 2010). In Section 3.2, we present our proposed method, IMSAT, for discrete representation learning. In Sections 3.3 and 3.4, we adapt IMSAT to the tasks of clustering and hash learning, respectively. In Section 3.5, we discuss an approximation technique for scaling up our method.

#### 3.1. Review of regularized information maximization for clustering

The RIM (Gomes et al., 2010) learns a probabilistic classifier  $p_\theta(y|x)$  such that mutual information (Cover & Thomas, 2012) between inputs and cluster assignments is maximized. At the same time, it regularizes the complexity of the classifier. Let  $X \in \mathcal{X}$  and  $Y \in \mathcal{Y} \equiv \{0, \dots, K-1\}$  denote random variables for data and cluster assignments, respectively, where  $K$  is the number of clusters. The RIM minimizes the objective:

$$\mathcal{R}(\theta) - \lambda \mathbf{I}(X; Y), \quad (1)$$

where  $\mathcal{R}(\theta)$  is the regularization penalty, and  $\mathbf{I}(X; Y)$  is mutual information between  $X$  and  $Y$ , which depends on  $\theta$  through the classifier,  $p_\theta(y|x)$ . Mutual information measures the statistical dependency between  $X$  and  $Y$ , and is 0 iff they are independent. Hyper-parameter  $\lambda \in \mathbb{R}$  trades off the two terms.

#### 3.2. Information Maximizing Self Augmented Training

Here, we present two components that make up our IMSAT. We present the information maximization part in Section 3.2.1 and the SAT part in Section 3.2.2.

##### 3.2.1. INFORMATION MAXIMIZATION FOR LEARNING DISCRETE REPRESENTATIONS

We extend the RIM and consider learning  $M$ -dimensional discrete representations of data. Let the output domain be

$\mathcal{Y} = \mathcal{Y}_1 \times \dots \times \mathcal{Y}_M$ , where  $\mathcal{Y}_m \equiv \{0, 1, \dots, V_m - 1\}$ ,  $1 \leq m \leq M$ . Let  $Y = (Y_1, \dots, Y_M) \in \mathcal{Y}$  be a random variable for the discrete representation. Our goal is to learn a multi-output probabilistic classifier  $p_\theta(y_1, \dots, y_M|x)$  that maps similar inputs into similar representations. For simplicity, we model the conditional probability  $p_\theta(y_1, \dots, y_M|x)$  by using the deep neural network depicted in Figure 1. Under the model,  $\{y_1, \dots, y_M\}$  are conditionally independent given  $x$ :

$$p_\theta(y_1, \dots, y_M|x) = \prod_{m=1}^M p_\theta(y_m|x). \quad (2)$$

Following the RIM Gomes et al. (2010), we maximize the mutual information between inputs and their discrete representations, while regularizing the multi-output probabilistic classifier. The resulting objective to *minimize* looks exactly the same as Eq. (1), except that  $Y$  is multi-dimensional in our setting.

##### 3.2.2. REGULARIZATION OF DEEP NEURAL NETWORKS VIA SELF AUGMENTED TRAINING

We present an intuitive and flexible regularization objective, termed *self augmented training* (SAT). SAT uses data augmentation to impose the intended invariance on the data representations. Essentially, SAT penalizes representation dissimilarity between augmented data points and the original data points. Let  $T : \mathcal{X} \rightarrow \mathcal{X}$  denote a pre-defined data augmentation under which the data representations should be invariant. The regularization of SAT made on data point  $x$  is

$$\begin{aligned} \mathcal{R}_{\text{SAT}}(\theta; x, T(x)) \\ = - \sum_{m=1}^M \sum_{y_m=0}^{V_m-1} p_{\hat{\theta}}(y_m|x) \log p_\theta(y_m|T(x)), \end{aligned} \quad (3)$$

where  $p_{\hat{\theta}}(y_m|x)$  is the prediction of original data point  $x$ , and  $\hat{\theta}$  is the current parameter of the network. In Eq. (3), the representations of the augmented data are pushed to be close to those of the original data. Since  $p_\theta(y|x)$  is modeled using a deep neural network, the probabilistic classifier is flexible enough to capture a wide range of invariances specified by the augmentation function  $T$ . The regularization by SAT is then the average of  $\mathcal{R}_{\text{SAT}}(\theta; T, x)$  over all the training data points:

$$\mathcal{R}_{\text{SAT}}(\theta; T) = \frac{1}{N} \sum_{n=1}^N \mathcal{R}_{\text{SAT}}(\theta; x_n, T(x_n)). \quad (4)$$

The augmentation function  $T$  can either be stochastic or deterministic. It can be designed specifically for the applications of interest. For example, for image data, affine

distortion such as rotation, scaling and parallel movement can be used for the augmentation function.

Alternatively, more general augmentation functions that do not depend on specific applications can be considered. A representative example is local perturbations. The use of local perturbations in SAT encourages the data representations to be locally invariant. The resulting decision boundaries between different representations tend to lie in low density regions of a data distribution. Such boundaries are generally preferred and follow the low-density separation principle (Grandvalet et al., 2004).

Here, we review virtual adversarial training (VAT) (Miyato et al., 2016), which is a representative regularization method based on the local perturbation. VAT was empirically shown to achieve better performance in supervised and semi-supervised learning scenarios compared to other local perturbation methods such as adversarial training (Goodfellow et al., 2015), random perturbation training (Bachman et al., 2014), and dropout (Srivastava et al., 2014).

In VAT, the augmentation function is

$$T_{\text{VAT}}(x) = x + r_{\text{adv}}, \quad (5)$$

where  $r_{\text{adv}}$  is  $\epsilon$ -perturbation to an adversarial direction:

$$r_{\text{adv}} = \arg \max_r \{ \mathcal{R}_{\text{SAT}}(\hat{\theta}; x, x + r); \|r\|^2 \leq \epsilon \}. \quad (6)$$

Hyper-parameter  $\epsilon$  controls the range of the local perturbation, and  $r_{\text{adv}}$  can be approximated efficiently by a pair of forward and backward passes. For further details, refer to Miyato et al. (2016).

### 3.3. IMSAT for clustering

In clustering, we can directly apply the RIM (Gomes et al., 2010) reviewed in Section 3.1. Unlike the original RIM, however, our method, IMSAT, uses deep neural networks for the classifiers and regularizes them via SAT. By representing mutual information as the difference between marginal entropy and conditional entropy (Cover & Thomas, 2012), we have the objective to minimize:

$$\mathcal{R}_{\text{SAT}}(\theta; T) - \lambda [H(Y) - H(Y|X)], \quad (7)$$

where  $H(\cdot)$  and  $H(\cdot|\cdot)$  are entropy and conditional entropy, respectively. Increasing the marginal entropy  $H(Y)$  encourages the cluster sizes to be uniform, while decreasing the conditional entropy  $H(Y|X)$  encourages unambiguous cluster assignments.

In practice, we can incorporate our prior knowledge on cluster sizes. Note that  $H(Y) = \log K - \text{KL}[p_\theta(y) || \mathcal{U}]$ , where  $K$  is the number of clusters,  $\text{KL}[\cdot || \cdot]$  is the Kullback-Leibler divergence, and  $\mathcal{U}$  is a uniform distribution. Hence,

maximization of  $H(Y)$  is equivalent to minimization of  $\text{KL}[p_\theta(y) || \mathcal{U}]$ , which encourages predicted cluster distribution  $p_\theta(y)$  to be close to  $\mathcal{U}$ . We can replace  $\mathcal{U}$  with any class prior probability  $q(y)$  and consider the following constrained optimization problem:

$$\begin{aligned} \min_{\theta} \quad & \mathcal{R}_{\text{SAT}}(\theta; T) + \lambda H(Y|X), \\ \text{subject to} \quad & \text{KL}[p_\theta(y) || q(y)] \leq \delta, \end{aligned} \quad (8)$$

where  $\delta > 0$  is a tolerance hyper-parameter that is set sufficiently small so that predicted cluster distribution  $p_\theta(y)$  is the same as class prior  $q(y)$  up to  $\delta$ -tolerance. Eq. (8) can be solved by using the penalty method (Bertsekas, 1999), which turns the original constrained optimization problem into a series of unconstrained optimization problems. Refer to Appendix B for the detail.

### 3.4. IMSAT for hash learning

In hash learning, each data point is mapped into a  $D$ -bit-binary code. Hence, the original RIM is not directly applicable. Instead, we apply our method for discrete representation learning presented in Section 3.2.1.

The computation of mutual information  $I(Y_1, \dots, Y_D; X)$ , however, is intractable for large  $D$  because it involves a summation over an exponential number of terms, each of which corresponds to a different configuration of hash bits.

Brown (2009) showed that mutual information  $I(Y_1, \dots, Y_D; X)$  can be expanded as the sum of interaction information (McGill, 1954):

$$I(Y_1, \dots, Y_D; X) = \sum_{C \subseteq S_Y} I(C \cup X), \quad |C| \geq 1, \quad (9)$$

where  $S_Y \equiv \{Y_1, \dots, Y_D\}$ . Note that  $I$  denotes interaction information when its argument is a set of random variables. Interaction information is a generalization of mutual information and can take a negative value. When the argument is a set of two random variables, the interaction information reduces to mutual information between the two random variables. Following Brown (2009), we only retain terms involving pairs of output dimensions in Eq. (9), i.e., all terms where  $|C| \leq 2$ . This gives us

$$\sum_{d=1}^D I(Y_d; X) + \sum_{1 \leq d \neq d' \leq D} I(\{Y_d, Y_{d'}, X\}). \quad (10)$$

This approximation ignores the interactions among hash bits beyond the pairwise interactions. It is related to the orthogonality constraint that is widely used in the literature to remove redundancy among hash bits (Wang et al., 2016). In fact, the orthogonality constraint encourages the covariance between a pair of hash bits to 0. Thus, it also takes into account the pairwise interactions.



It follows from the definition of interaction information and the conditional independence in Eq. (2) that

$$\begin{aligned} I(\{Y_d, Y_{d'}, X\}) &\equiv I(Y_d; Y_{d'}|X) - I(Y_d; Y_{d'}) \\ &= -I(Y_d; Y_{d'}). \end{aligned} \quad (11)$$

In summary, our approximated objective to minimize is

$$\mathcal{R}_{\text{SAT}}(\theta; T) - \lambda \left( \sum_{d=1}^D I(X; Y_d) - \sum_{1 \leq d \neq d' \leq D} I(Y_d; Y_{d'}) \right). \quad (12)$$

The first term regularizes the neural network. The second term maximizes the mutual information between data and each hash bit, and the third term removes the redundancy among the hash bits.

### 3.5. Approximation of the marginal distribution

To scale up our method to large datasets, we would like the objective in Eq. (1) to be amenable to optimization based on mini-batch SGD. For the regularization term, we use the SAT in Eq. (4), which is the sum of per sample penalties and can be readily adapted to mini-batch computation. For the approximated mutual information in Eq. (12), we can decompose it into three parts: (i) conditional entropy  $H(Y_d|X)$ , (ii) marginal entropy  $H(Y_d)$ , and (iii) mutual information between a pair of output dimensions  $I(Y_d; Y_{d'})$ . The conditional entropy only consists of a sum over per example entropies; thus, can be adapted to mini-batch computation. However, the marginal entropy and mutual information involve the marginal distribution over a subset of target dimensions, i.e.,  $p_\theta(c) \equiv \frac{1}{N} \sum_{n=1}^N p_\theta(c|x_n)$ , where  $c \subseteq \{y_1, \dots, y_M\}$ . We see that the marginal distribution can only be calculated using the entire dataset and is not amenable to the mini-batch setting. Following [Springenberg \(2015\)](#), we approximate the marginal distributions using mini-batch data:

$$p_\theta(c) \approx \frac{1}{|\mathcal{B}|} \sum_{x \in \mathcal{B}} p_\theta(c|x) \equiv \hat{p}_\theta^{(\mathcal{B})}(c), \quad (13)$$

where  $\mathcal{B}$  is a set of data in the mini-batch.

The mini-batch approximation can be validated for the clustering scenario in Eq. (8) as follows. By the convexity of the KL divergence ([Cover & Thomas, 2012](#)) and Jensen’s inequality, we have

$$\mathbb{E}_{\mathcal{B}}[\text{KL}[\hat{p}_\theta^{(\mathcal{B})}(y)||q(y)]] \geq \text{KL}[p_\theta(y)||q(y)] \geq 0, \quad (14)$$

where the first expectation is taken with respect to the randomness of the mini-batch selection. Therefore, in the penalty method, the constraint on the exact KL divergence, i.e.,  $\text{KL}[p_\theta(y)||q(y)] \leq \delta$  can be satisfied by minimizing

its upper bound, which is the approximated KL divergence  $\mathbb{E}_{\mathcal{B}}[\text{KL}[\hat{p}_\theta^{(\mathcal{B})}(y)||q(y)]]$ . Obviously, the approximated KL divergence is amenable to the mini-batch setting; thus, can be minimized with SGD.

## 4. Experiments

In this section, we evaluate IMSAT for clustering and hash learning on benchmark datasets.

### 4.1. Implementation

In unsupervised learning, it is not straightforward to determine hyper-parameters by cross-validation. Therefore, in all the experiments with benchmark datasets, we used commonly reported parameter values for deep neural networks and avoided dataset-specific tuning as much as possible. Specifically, inspired by [Hinton et al. \(2012\)](#), we set the network dimensionality to  $d$ -1200-1200- $M$  for clustering across all the datasets, where  $d$  and  $M$  are input and output dimensionality, respectively. For hash learning, we used smaller network sizes to ensure fast computation of mapping data into hash codes. We used rectified linear units ([Jarrett et al., 2009](#); [Nair & Hinton, 2010](#); [Glorot et al., 2011](#)) for all the activations and applied batch normalization ([Ioffe & Szegedy, 2015](#)) to each layer to accelerate training. For optimization, we used Adam ([Kingma & Ba, 2015](#)) with the step size set to 0.002. Refer to Appendix C for further details. Our implementation based on Chainer ([Tokui et al., 2015](#)) is available online<sup>1</sup>.

### 4.2. Clustering

We now report the experimental results on clustering on eight benchmark datasets.

#### 4.2.1. DATASETS AND COMPARED METHODS

We evaluated our method for clustering presented in Section 3.3 on eight benchmark datasets by comparing it with other clustering methods including k-means, DEC ([Zheng et al., 2016](#)), Auto-encoder (AE)+K-means ([Zheng et al., 2016](#)), Linear RIM and Deep RIM. The latter two are variants of IMSAT with different classifiers or regularizations. Linear RIM uses a linear classifier with a weight decay while deep RIM uses a deep neural network with a weight decay. We did not include kernel classifiers because they did not scale to most of the benchmark datasets.

We conducted experiments on eight benchmark datasets. Brief dataset descriptions are given below. Refer to Appendix D for further details.

- **MNIST:** A dataset of hand-written digit classification

<sup>1</sup><https://github.com/weihua916/imsat>

Table 1. Comparison of clustering accuracy on 8 benchmark datasets (%). Results marked with † were excerpted from Xie et al. (2016).

| Method      | MNIST             | Omniglot          | STL               | CIFAR10           | CIFAR100          | SVHN              | Reuters           | 20news            |
|-------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| K-means     | 53.2              | 12.0              | 85.6              | 34.4              | 21.5              | 17.9              | 54.1              | 15.5              |
| AE+K-means  | 79.8 †            | 14.1              | 72.2              | 44.2              | 20.8              | 17.4              | 67.2              | 22.1              |
| DEC         | 84.3 †            | 5.7 (0.3)         | 78.1 (0.1)        | <b>46.9 (0.9)</b> | 14.3 (0.6)        | 11.9 (0.4)        | 67.3 (0.2)        | 30.8 (1.8)        |
| Linear RIM  | 59.6 (2.3)        | 11.1 (0.2)        | 73.5 (6.5)        | 40.3 (2.1)        | 23.7 (0.8)        | 20.2 (1.4)        | 62.8 (7.8)        | <b>50.9 (3.1)</b> |
| Deep RIM    | 58.5 (3.5)        | 5.8 (2.2)         | 92.5 (2.2)        | 40.3 (3.5)        | 13.4 (1.2)        | 26.8 (3.2)        | 62.3 (3.9)        | 25.1 (2.8)        |
| IMSAT (VAT) | <b>98.4 (0.4)</b> | <b>24.0 (0.9)</b> | <b>94.1 (0.4)</b> | 45.6 (0.8)        | <b>27.5 (0.4)</b> | <b>57.3 (3.9)</b> | <b>71.0 (4.9)</b> | 31.1 (1.9)        |

(LeCun et al., 1998).

- **Omniglot**: A dataset of hand-written character recognition (Lake et al., 2011), 100 types of characters were sampled, each containing 20 data points. Each data point was down-sampled into a 21-by-21 pixel image, and was augmented 20 times by the similar stochastic affine distortion used by Koch (2015). Refer to Appendix E for the details of the affine distortion used in our experiments.
- **STL**: A dataset of 96-by-96 color images acquired from labeled examples on ImageNet (Coates et al., 2010).
- **CIFAR10**: A dataset of 32-by-32 color images with 10 object classes, which are from the Tiny image dataset (Torralba et al., 2008).
- **CIFAR100**: A dataset 32-by-32 color images with 100 refined object classes, which are from the Tiny image dataset (Torralba et al., 2008).
- **SVHN**: A dataset with street view house numbers (Netzer et al., 2011). Training and test images were both used. Each image was represented as a 960-dimensional GIST feature (Oliva & Torralba, 2001).
- **Reuters**: A dataset with English news stories labeled with a category tree (Lewis et al., 2004). As was done by Xie et al. (2016), 10000 documents were randomly sampled and four category labels were used.
- **20news**: A dataset of newsgroup documents, partitioned nearly evenly across 20 different newsgroups<sup>2</sup>.

In the experiments, our objective is to discover clusters that correspond well with the ground-truth categories. For the STL, CIFAR10 and CIFAR100 datasets, raw pixels are not suited for such an objective because color information is dominant in the raw features. We therefore applied 50-layer pre-trained deep residual networks (He et al., 2016) to extract features and used them for clustering. Note that since the residual network is trained on ImageNet, we expect that each class of the STL dataset is well-separated in the feature space. For Reuters and 20news, we removed stop words and retained the 2000 most frequent words. We then used *tf-idf* features.

<sup>2</sup><http://qwone.com/~jason/20Newsgroups/>

#### 4.2.2. EVALUATION METRIC

Following Xie et al. (2016), we set the number of clusters to the number of ground-truth categories and evaluated clustering performance with unsupervised clustering accuracy (ACC):

$$\text{ACC} = \max_m \frac{\sum_{n=1}^N \mathbf{1}\{l_n = m(c_n)\}}{N}, \quad (15)$$

where  $l_n$  and  $c_n$  are the ground-truth label and cluster assignment produced using the algorithm for  $x_n$ , respectively. The  $m$  ranges over all possible one-to-one mappings between clusters and labels. The best mapping can be efficiently computed using the Hungarian algorithm (Kuhn, 1955).

#### 4.2.3. HYPER-PARAMETER SELECTION

We conducted experiments 12 times for each dataset and report the average and standard deviation.

In unsupervised learning, it is not straightforward to determine hyper-parameters by cross-validation. Hence, we fixed hyper-parameters across all the datasets unless there was an objective way to select them. For *K*-means, we tried 12 different initializations and reported the results with the best objectives. For AE+*K*-means and DEC (Xie et al., 2016), we used the recommended hyper-parameters for the network dimensionality and annealing speed.

For the experiments on the eight benchmark datasets, we used VAT (Miyato et al., 2016) to impose the local invariance on the discrete representations predicted by neural networks. For experiments with the Omniglot dataset, we used SAT with affine distortion as well as VAT. Inspired by the automatic kernel width selection in spectral clustering (Zelnik-Manor & Perona, 2004), the perturbation range,  $\epsilon$ , in VAT for data point  $x$  was set as

$$\epsilon(x) = \alpha \cdot \sigma_t(x), \quad (16)$$

where  $\alpha$  is a scalar and  $\sigma_t(x)$  is the Euclidian distance to the  $t$ -th neighbor of  $x$ . In our extensive experiments, we fixed  $t = 10$  and  $\alpha = 1/4$ , which performed well across the datasets for both clustering and hash learning.

Table 2. Comparison of clustering accuracy on the Omniglot dataset using IMSAT with different types of self augmented training.

| Method               | Omniglot          |
|----------------------|-------------------|
| IMSAT (VAT)          | 24.0 (0.9)        |
| IMSAT (affine)       | 45.1 (2.0)        |
| IMSAT (VAT & affine) | <b>70.0 (2.0)</b> |

For Linear RIM, Deep RIM and IMSAT, we varied one hyper-parameter and chose the best one that performed well across the datasets. More specifically, for Linear RIM and Deep RIM, we varied the decay rate over  $0.0025 \cdot 2^i, i = 0, 1, \dots, 7$ . For IMSAT, we used VAT for the regularization and varied  $\lambda$  in Eq. (18) for  $0.025 \cdot 2^i, i = 0, 1, \dots, 7$ . We set  $q$  to be the uniform distribution and let  $\delta = 0.01 \cdot H(q(y))$  in Eq. (8) for the all experiments.

Consequently, we chose 0.005 for decay rates in both Linear RIM and Deep RIM, and 0.1 for  $\lambda$  in IMSAT. We hereforth fixed these hyper-parameters throughout the experiments for both clustering and hash learning. In Appendix F, we report all the experimental results and the criteria to choose the parameters.

#### 4.2.4. EXPERIMENTAL RESULTS

In Table 1, we compare clustering performance across eight benchmark datasets. Our method, IMSAT (VAT), performed well across the datasets. The fact that IMSAT outperformed Linear RIM and Deep RIM for most datasets suggests the importance of using neural networks with an end-to-end regularization.

We further conducted experiments on the Omniglot dataset to demonstrate that clustering performance can be improved by incorporating domain-specific knowledge in the augmentation function of SAT. Specifically, we used the affine distortion in addition to VAT for the augmented function of SAT. We compared the clustering accuracy of IMSAT with three different augmentation functions: VAT, affine distortion, and the combination of VAT & affine distortion, in which we simply set the regularization to be

$$\frac{1}{2} \cdot \mathcal{R}_{\text{SAT}}(\theta; T_{\text{VAT}}) + \frac{1}{2} \cdot \mathcal{R}_{\text{SAT}}(\theta; T_{\text{affine}}), \quad (17)$$

where  $T_{\text{VAT}}$  and  $T_{\text{affine}}$  are augmentation functions of VAT and affine distortion, respectively. We used the stochastic affine distortion function,  $T_{\text{affine}}$ , defined in Appendix E.

We report the clustering accuracy of Omniglot in Table 2. We see that including affine distortion in data augmentation significantly improved clustering accuracy. Figure 2 shows ten randomly selected clusters of the Omniglot dataset found using IMSAT (VAT) and IMSAT (VAT & affine dis-

tortion). We observed that IMSAT (VAT & affine distortion) was able to discover cluster assignments that are invariant to affine distortion as we intended. These results suggest that our method successfully captured the invariance in the hand-written character recognition in an unsupervised way.

### 4.3. Hash learning

We now report the experimental results for hash learning on two benchmark datasets.

#### 4.3.1. DATASETS AND COMPARED METHODS

For the comparison with Deep Hash (Erin Liong et al., 2015), we evaluated our method for hash learning presented in Section 3.4 on 2 benchmark datasets: MNIST and CIFAR10 datasets. Each data point of CIFAR10 is represented as a 512-dimensional GIST feature (Oliva & Torralba, 2001). Our method was compared against several unsupervised hash learning methods: spectral hashing (Weiss et al., 2009), PCA-ITQ (Weiss et al., 2009), and Deep Hash (Erin Liong et al., 2015). We also compared it with the hash versions of Linear RIM and Deep RIM. We used VAT for SAT in our method.

#### 4.3.2. EVALUATION METRIC

Following Erin Liong et al. (2015), we used three evaluation metrics to measure the performance of the different methods: 1) mean average precision (mAP); 2) precision at  $N = 500$  samples; and 3) Hamming look-up result where the hamming radius is set as  $r = 2$ . We used the class labels to define the neighbors. We repeated the experiments ten times and took the average as the final result.

#### 4.3.3. EXPERIMENTAL RESULTS

The MNIST and CIFAR10 datasets both have 10 classes, and contain 70000 and 60000 data points, respectively. Following Erin Liong et al. (2015), we randomly sampled 1000 samples, 100 per class, as the query data and used the remaining data as the gallery set.

We tested performance for 16 and 32-bit hash codes. In practice, fast computation of hash codes is crucial for fast information retrieval. Hence, small networks are preferable. We therefore tested our method on three different network sizes: the same ones as Deep Hash (Erin Liong et al., 2015),  $d$ -200-200- $M$ , and  $d$ -400-400- $M$ . Deep Hash used  $d$ -60-30- $M$  for 16-bit hash codes and  $d$ -80-50- $M$  for 32-bit hash codes.

Table 3 lists the results for 16-bit hash. Due to the space constraint, we report the results for 32-bit hash codes in Appendix G. Note that the results showed a similar tendency as that of 16-bit hash codes. From Table 3, we

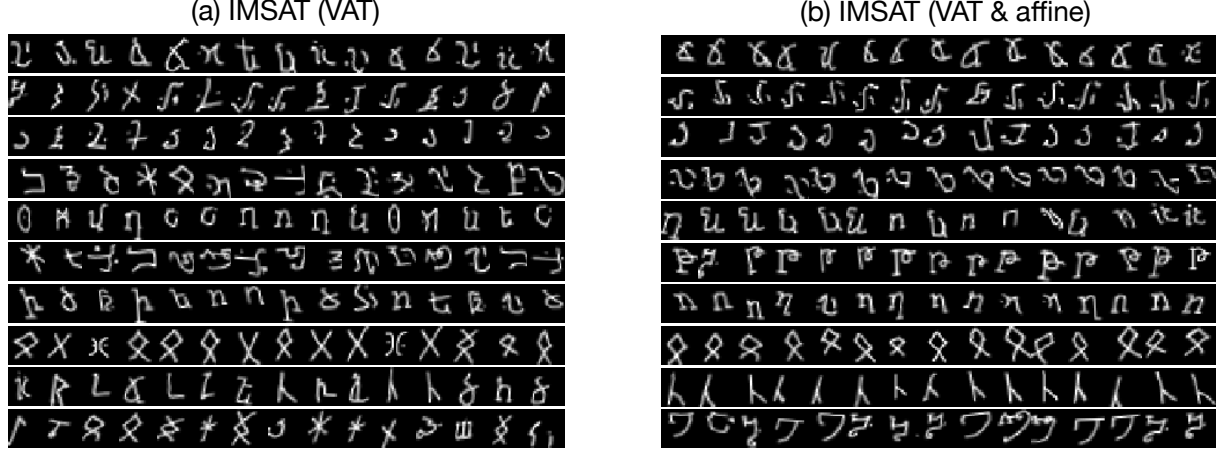


Figure 2. Randomly sampled clusters of Omniglot discovered using (a) IMSAT (VAT) and (b) IMSAT (VAT & affine). Each row contains randomly sampled data points in same cluster.

Table 3. Comparison of hash performance for 16-bit hash codes (%). Experimental results of Deep Hash and the previous methods are excerpted from [Erin Liong et al. \(2015\)](#).

| Method<br>(Dimensions of hidden layers)              | Hamming ranking (mAP) |                   | precision @ sample = 500 |                   | precision @ r = 2 |                   |
|--|-----------------------|-------------------|--------------------------|-------------------|-------------------|-------------------|
|  | MNIST                 | CIFAR10           | MNIST                    | CIFAR10           | MNIST             | CIFAR10           |
| Spectral hash ( <a href="#">Weiss et al., 2009</a> ) | 26.6                  | 12.6              | 56.3                     | 18.8              | 57.5              | 18.5              |
| PCA-ITQ ( <a href="#">Gong et al., 2013</a> )        | 41.2                  | 15.7              | 66.4                     | 22.5              | 65.7              | 22.6              |
| Deep Hash (60-30 or 80-50)                           | 43.1                  | 16.2              | 67.9                     | 23.8              | 66.1              | 23.3              |
| Linear RIM   | 35.9 (0.6)            | <b>24.0 (3.5)</b> | 68.9 (1.1)               | 15.9 (0.5)        | 71.3 (0.9)        | 14.2 (0.3)        |
| Deep RIM (60-30)                                     | 42.7 (2.8)            | 15.2 (0.5)        | 67.9 (2.7)               | 21.8 (0.9)        | 65.9 (2.7)        | 21.2 (0.9)        |
| Deep RIM (200-200)                                   | 43.7 (3.7)            | 15.6 (0.6)        | 68.7 (4.9)               | 21.6 (1.2)        | 67.0 (4.9)        | 21.1 (1.1)        |
| Deep RIM (400-400)                                   | 43.9 (2.7)            | 15.4 (0.2)        | 69.0 (3.2)               | 21.5 (0.4)        | 66.7 (3.2)        | 20.9 (0.3)        |
| <b>IMSAT (VAT)</b> (60-30)                           | 61.2 (2.5)            | 19.8 (1.2)        | 78.6 (2.1)               | 21.0 (1.8)        | 76.5 (2.3)        | 19.3 (1.6)        |
| <b>IMSAT (VAT)</b> (200-200)                         | 80.7 (2.2)            | 21.2 (0.8)        | 95.8 (1.0)               | <b>27.3 (1.3)</b> | 94.6 (1.4)        | 26.1 (1.3)        |
| <b>IMSAT (VAT)</b> (400-400)                         | <b>83.9 (2.3)</b>     | 21.4 (0.5)        | <b>97.0 (0.8)</b>        | <b>27.3 (1.1)</b> | <b>96.2 (1.1)</b> | <b>26.4 (1.0)</b> |

see that IMSAT with the largest network sizes (400-400) achieved competitive performance in both datasets. The performance of IMSAT improved significantly when we used slightly bigger networks (200-200), while the performance of Deep RIM did not improve much with the larger networks. We deduce that this is because we can better model the local invariance with VAT by using more flexible networks. Deep RIM, on the other hand, did not significantly benefit from the larger networks, because the additional flexibility of the networks was not used by the global function regularization via weight-decay<sup>3</sup>. In Appendix H, our deduction is supported using a toy dataset.

<sup>3</sup>Hence, we deduce that Deep Hash, which is only regularized by weight-decay, would not benefit much by using larger networks.

## 5. Conclusion & Future work

In this paper, we presented IMSAT, an information-theoretic method for unsupervised discrete representation learning using deep neural networks. Through extensive experiments, we showed that intended discrete representations can be obtained by directly imposing the invariance to data augmentation on the prediction of neural networks in an end-to-end fashion. For future work, it is interesting to apply our method to structured data, i.e., graph data or sequential data, by considering appropriate data augmentation.

## Acknowledgements

We thank Vogel Brian for helpful discussions and insightful reviews on the paper. This paper is based on results obtained from an internship at Preferred Networks, Inc.



## References

- Bachman, Philip, Alsharif, Ouais, and Precup, Doina. Learning with pseudo-ensembles. In *NIPS*, 2014.
- Berkhin, Pavel. A survey of clustering data mining techniques. In *Grouping multidimensional data*, pp. 25–71. Springer, 2006.
- Bertsekas, Dimitri P. *Nonlinear programming*. Athena scientific Belmont, 1999.
- Brown, Gavin. A new perspective for information theoretic feature selection. In *AISTATS*, 2009.
- Coates, Adam, Lee, Honglak, and Ng, Andrew Y. An analysis of single-layer networks in unsupervised feature learning. *Ann Arbor*, 1001(48109):2, 2010.
- Cover, Thomas M and Thomas, Joy A. *Elements of information theory*. John Wiley & Sons, 2012.
- Dilokthanakul, Nat, Mediano, Pedro AM, Garnelo, Marta, Lee, Matthew CH, Salimbeni, Hugh, Arulkumaran, Kai, and Shanahan, Murray. Deep unsupervised clustering with gaussian mixture variational autoencoders. *arXiv preprint arXiv:1611.02648*, 2016.
- Erin Liong, Venice, Lu, Jiwen, Wang, Gang, Moulin, Pierre, and Zhou, Jie. Deep hashing for compact binary codes learning. In *CVPR*, 2015.
- Glorot, Xavier, Bordes, Antoine, and Bengio, Yoshua. Deep sparse rectifier neural networks. In *AISTATS*, 2011.
- Gomes, Ryan, Krause, Andreas, and Perona, Pietro. Discriminative clustering by regularized information maximization. In *NIPS*, 2010.
- Gong, Yunchao, Lazebnik, Svetlana, Gordo, Albert, and Perronnin, Florent. Iterative quantization: A procrustean approach to learning binary codes for large-scale image retrieval. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(12):2916–2929, 2013.
- Goodfellow, Ian J, Shlens, Jonathon, and Szegedy, Christian. Explaining and harnessing adversarial examples. In *ICLR*, 2015.
- Grandvalet, Yves, Bengio, Yoshua, et al. Semi-supervised learning by entropy minimization. In *NIPS*, 2004.
- He, Kaiming, Wen, Fang, and Sun, Jian. K-means hashing: An affinity-preserving quantization method for learning binary compact codes. In *CVPR*, 2013.
- He, Kaiming, Zhang, Xiangyu, Ren, Shaoqing, and Sun, Jian. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *CVPR*, 2015.
- He, Kaiming, Zhang, Xiangyu, Ren, Shaoqing, and Sun, Jian. Deep residual learning for image recognition. In *CVPR*, 2016.
- Hinton, Geoffrey E, Srivastava, Nitish, Krizhevsky, Alex, Sutskever, Ilya, and Salakhutdinov, Ruslan R. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012.
- Ioffe, Sergey and Szegedy, Christian. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- Jarrett, Kevin, Kavukcuoglu, Koray, Ranzato, Marc’Aurelio, and LeCun, Yann. What is the best multi-stage architecture for object recognition? In *ICCV*, 2009.
- Kingma, Diederik and Ba, Jimmy. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- Koch, Gregory. *Siamese neural networks for one-shot image recognition*. PhD thesis, University of Toronto, 2015.
- Kuhn, Harold W. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.
- Kulis, Brian and Darrell, Trevor. Learning to hash with binary reconstructive embeddings. In *NIPS*, 2009.
- Lai, Hanjiang, Pan, Yan, Liu, Ye, and Yan, Shuicheng. Simultaneous feature learning and hash coding with deep neural networks. In *CVPR*, 2015.
- Lake, Brenden M, Salakhutdinov, Ruslan, Gross, Jason, and Tenenbaum, Joshua B. One shot learning of simple visual concepts. In *CogSci*, 2011.
- LeCun, Yann, Bottou, Léon, Bengio, Yoshua, and Haffner, Patrick. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Leen, Todd K. From data distributions to regularization in invariant learning. *Neural Computation*, 7(5):974–981, 1995.
- Lewis, David D, Yang, Yiming, Rose, Tony G, and Li, Fan. Rcv1: A new benchmark collection for text categorization research. *Journal of machine learning research*, 5 (Apr):361–397, 2004.
- Li, Wu-Jun, Wang, Sheng, and Kang, Wang-Cheng. Feature learning based deep supervised hashing with pairwise labels. In *IJCAI*, 2015.
- McGill, William J. Multivariate information transmission. *Psychometrika*, 19(2):97–116, 1954.

- Miyato, Takeru, Maeda, Shin-ichi, Koyama, Masanori, Nakae, Ken, and Ishii, Shin. Distributional smoothing with virtual adversarial training. In *ICLR*, 2016.
- Nair, Vinod and Hinton, Geoffrey E. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010.
- Netzer, Yuval, Wang, Tao, Coates, Adam, Bissacco, Alessandro, Wu, Bo, and Ng, Andrew Y. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, 2011.
- Ng, Andrew Y, Jordan, Michael I, Weiss, Yair, et al. On spectral clustering: Analysis and an algorithm. In *NIPS*, 2001.
- Norouzi, Mohammad and Blei, David M. Minimal loss hashing for compact binary codes. In *ICML*, 2011.
- Oliva, Aude and Torralba, Antonio. Modeling the shape of the scene: A holistic representation of the spatial envelope. *International journal of computer vision*, 42(3): 145–175, 2001.
- Rifai, Salah, Vincent, Pascal, Muller, Xavier, Glorot, Xavier, and Bengio, Yoshua. Contractive auto-encoders: Explicit invariance during feature extraction. In *ICML*, 2011.
- Sajjadi, Mehdi, Javanmardi, Mehran, and Tasdizen, Tolga. Regularization with stochastic transformations and perturbations for deep semi-supervised learning. In *NIPS*, 2016.
- Salakhutdinov, Ruslan and Hinton, Geoffrey. Semantic hashing. *International Journal of Approximate Reasoning*, 50(7):969–978, 2009.
- Springenberg, Jost Tobias. Unsupervised and semi-supervised learning with categorical generative adversarial networks. In *ICLR*, 2015.
- Srivastava, Nitish, Hinton, Geoffrey E, Krizhevsky, Alex, Sutskever, Ilya, and Salakhutdinov, Ruslan. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1): 1929–1958, 2014.
- Tokui, Seiya, Oono, Kenta, Hido, Shohei, and Clayton, Justin. Chainer: a next-generation open source framework for deep learning. In *Proceedings of workshop on machine learning systems (LearningSys) in the twenty-ninth annual conference on neural information processing systems (NIPS)*, 2015.
- Torralba, Antonio, Fergus, Rob, and Freeman, William T. 80 million tiny images: A large data set for nonparametric object and scene recognition. *IEEE transactions on pattern analysis and machine intelligence*, 30(11):1958–1970, 2008.
- Vincent, Pascal, Larochelle, Hugo, Bengio, Yoshua, and Manzagol, Pierre-Antoine. Extracting and composing robust features with denoising autoencoders. In *ICML*, 2008.
- Wang, Jun, Liu, Wei, Kumar, Sanjiv, and Chang, Shih-Fu. Learning to hash for indexing big data—a survey. *Proceedings of the IEEE*, 104(1):34–57, 2016.
- Weiss, Yair, Torralba, Antonio, and Fergus, Rob. Spectral hashing. In *NIPS*, 2009.
- Xia, Rongkai, Pan, Yan, Lai, Hanjiang, Liu, Cong, and Yan, Shuicheng. Supervised hashing for image retrieval via image representation learning. In *AAAI*, 2014.
- Xie, Junyuan, Girshick, Ross, and Farhadi, Ali. Unsupervised deep embedding for clustering analysis. In *ICML*, 2016.
- Xu, Jiaming, Wang, Peng, Tian, Guanhua, Xu, Bo, Zhao, Jun, Wang, Fangyuan, and Hao, Hongwei. Convolutional neural networks for text hashing. In *IJCAI*, 2015.
- Xu, Linli, Neufeld, James, Larson, Bryce, and Schuurmans, Dale. Maximum margin clustering. In *NIPS*, 2004.
- Zelnik-Manor, Lihi and Perona, Pietro. Self-tuning spectral clustering. In *NIPS*, 2004.
- Zhang, Ruimao, Lin, Liang, Zhang, Rui, Zuo, Wangmeng, and Zhang, Lei. Bit-scalable deep hashing with regularized similarity learning for image retrieval and person re-identification. *IEEE Transactions on Image Processing*, 24(12):4766–4779, 2015.
- Zheng, Yin, Tan, Huachun, Tang, Bangsheng, Zhou, Hanning, et al. Variational deep embedding: A generative approach to clustering. *arXiv preprint arXiv:1611.05148*, 2016.

## A. Relation to denoising and contractive auto-encoders

Our method is related to denoising auto-encoders (Vincent et al., 2008). Auto-encoders maximize a lower bound of mutual information (Cover & Thomas, 2012) between inputs and their hidden representations (Vincent et al., 2008), while the denoising mechanism regularizes the auto-encoders to be locally invariant. However, such a regularization does not necessarily impose the invariance on the *hidden representations* because the decoder network also has the flexibility to model the invariance to data perturbations. SAT is more direct in imposing the intended invariance on hidden representations predicted by the encoder network.

Contractive auto-encoders (Rifai et al., 2011) directly impose the local invariance on the encoder network by minimizing the Frobenius norm of the Jacobian with respect to the weight matrices. However, it is empirically shown that such regularization attained lower generalization performance in supervised and semi-supervised settings than VAT, which regularizes neural networks in an end-to-end fashion (Miyato et al., 2016). Hence, we adopted the end-to-end regularization in our unsupervised learning. In addition, our regularization, SAT, has the flexibility of modeling other types invariance such as invariance to affine distortion, which cannot be modeled with the contractive regularization. Finally, compared with the auto-encoders approaches, our method does not require learning the decoder network; thus, is computationally more efficient.

## B. Penalty method and its implementation

Our goal is to optimize the constrained objective of Eq. (8):

$$\begin{aligned} \min_{\theta} \quad & \mathcal{R}_{\text{SAT}}(\theta; T) + \lambda H(Y|X), \\ \text{subject to} \quad & \text{KL}[p_{\theta}(y) || q(y)] \leq \delta. \end{aligned}$$

We use the penalty method (Bertsekas, 1999) to solve the optimization. We introduce a scalar parameter  $\mu$  and consider minimizing the following unconstrained objective:

$$\mathcal{R}_{\text{SAT}}(\theta; T) + \lambda H(Y|X) + \mu \max\{\text{KL}[p_{\theta}(y) || q(y)] - \delta, 0\}. \quad (18)$$

We increase  $\mu$  and solve the optimization of Eq. (18) for a fixed  $\mu$ . Let  $\mu^*$  be the smallest value for which the solution of Eq. (18) satisfies the constraint of Eq. (8). The penalty method ensures that the solution obtained by solving Eq. (18) with  $\mu = \mu^*$  is the same as that of the constrained optimization of Eq. (8).

In experiments in Section 4.2, we increased  $\mu$  in the order of  $\lambda, 2\lambda, 4\lambda, 6\lambda, \dots$  until the solution of Eq. (18) satisfied the constraint of Eq. (8).

## C. Implementation detail

We set the size of mini-batch to 250, and ran 50 epochs for each dataset. We initialized weights following He et al. (2015): each element of the weight is initialized by the value drawn independently from Gaussian distribution whose mean is 0, and standard deviation is  $scale \times \sqrt{2/fan_{in}}$ , where  $fan_{in}$  is the number of input units. We set the  $scale$  to be 0.1-0.1-0.0001 for weight matrices from the input to the output. The bias terms were all initialized with 0.

## D. Datasets description

We now describe the datasets used in the experiments in Section 4.2. The dataset statistics are summarized in Table 4.

- **MNIST**: A dataset of hand-written digit classification (LeCun et al., 1998). The value of each pixel was transformed linearly into an interval  $[-1, 1]$ .
- **Omniglot**: A dataset of hand-written character recognition (Lake et al., 2011), containing examples from 50 alphabets ranging from well-established international languages. We sampled 100 types of characters from four alphabets, Magi, Anglo-Saxon Futhorc, Arcadian, and Armenian. Each character contains 20 data points. Since the original data have high resolution (105-by-105 pixels), each data point was down-sampled to 21-by-21 pixels. We also augmented each data point 20 times by the stochastic affine distortion explained in Appendix E.

Table 4. Summary of dataset statistics.

| Dataset  | #Points | #classes | Dimension | %Largest class |
|----------|---------|----------|-----------|----------------|
| MNIST    | 70000   | 10       | 784       | 11%            |
| Omniglot | 40000   | 100      | 441       | 1%             |
| STL      | 13000   | 10       | 2048      | 10%            |
| CIFAR10  | 60000   | 10       | 2048      | 10%            |
| CIFAR100 | 60000   | 100      | 2048      | 1%             |
| SVHN     | 99289   | 10       | 960       | 19%            |
| Reuters  | 10000   | 4        | 2000      | 43%            |
| 20news   | 18040   | 20       | 2000      | 5%             |

- **STL**: A dataset of 96-by-96 color images acquired from labeled examples on ImageNet (Coates et al., 2010). Features were extracted using 50-layer pre-trained deep residual networks (He et al., 2016) available online as a caffe model. Note that since the residual network is also trained on ImageNet, we expect that each class is separated well in the feature space.
- **CIFAR10**: A dataset of 32-by-32 color images with ten object classes, which are from the Tiny image dataset (Torralba et al., 2008). Features were extracted using the 50-layer pre-trained deep residual networks (He et al., 2016).
- **CIFAR100**: A dataset 32-by-32 color images with 100 refined object classes, which are from the Tiny image dataset (Torralba et al., 2008). Features were extracted using the 50-layer pre-trained deep residual networks (He et al., 2016).
- **SVHN**: A dataset with street view house numbers (Netzer et al., 2011). Training and test images were both used. Each image was represented as a 960-dimensional GIST feature (Oliva & Torralba, 2001).
- **Reuters**: A dataset with English news stories labeled with a category tree (Lewis et al., 2004). Following DEC (Xie et al., 2016), we used four categories: corporate/industrial, government/social, markets, and economics as labels. The preprocessing was the same as that used by Xie et al. (2016), except that we removed stop words. As Xie et al. (2016) did, 10000 documents were randomly sampled, and *tf-idf* features were used.
- **20news**: A dataset of newsgroup documents, partitioned nearly evenly across 20 different newsgroups<sup>4</sup>. As Reuters dataset, stop words were removed, and the 2000 most frequent words were retained. Documents with less than ten words were then removed, and *tf-idf* features were used.

For the STL, CIFAR10 and CIFAR100 datasets, each image was first resized into a 224-by-224 image before its feature was extracted using the deep residual network.

## E. Affine distortion for the Omniglot dataset

We applied stochastic affine distortion to data points in Omniglot. The affine distortion is similar to the one used by Koch (2015), except that we applied the affine distortion on down-sampled images in our experiments. The followings are the stochastic components of the affine distortion used in our experiments. Our implementation of the affine distortion is based on scikit-image<sup>5</sup>.

- Random scaling along  $x$  and  $y$ -axis by a factor of  $(s_x, s_y)$ , where  $s_x$  and  $s_y$  are drawn uniformly from interval  $[0.8, 1.2]$ .
- Random translation along  $x$  and  $y$ -axis by  $(t_x, t_y)$ , where  $t_x$  and  $t_y$  are drawn uniformly from interval  $[-0.4, 0.4]$ .
- Random rotation by  $\theta$ , where  $\theta$  is drawn uniformly from interval  $[-10^\circ, 10^\circ]$ .
- Random shearing along  $x$  and  $y$ -axis by  $(\rho_x, \rho_y)$ , where  $\theta$  is drawn uniformly from interval  $[-0.3, 0.3]$ .

Figure. 3 shows examples of the random affine distortion.

<sup>4</sup><http://qwone.com/~jason/20Newsgroups/>

<sup>5</sup><http://scikit-image.org/>



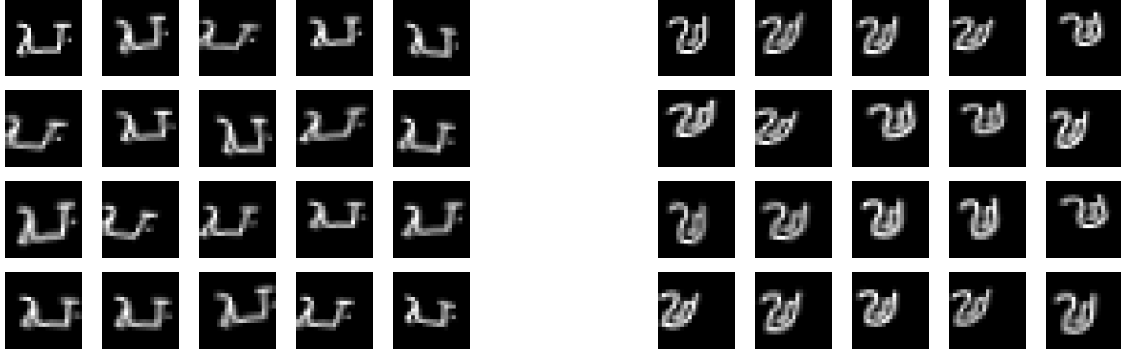


Figure 3. Examples of the random affine distortion used in our experiments. Images in the top left side are stochastically transformed using the affine distortion.

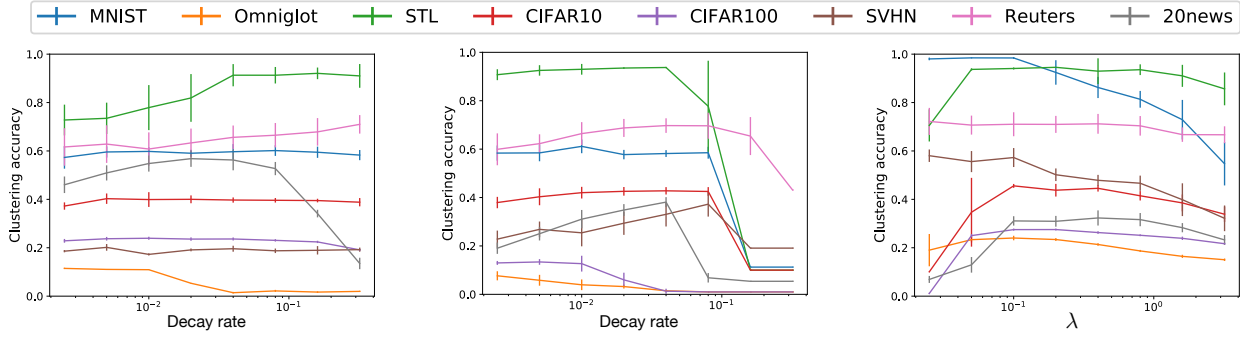


Figure 4. Relationship between hyper-parameters and clustering accuracy for 8 benchmark datasets with different methods: (a) Linear RIM, (b) Deep RIM, and (c) IMSAT (VAT).

## F. Hyper-parameter selection

In Figure 4 we report the experimental results for different hyper-parameter settings. We used Eq. (19) as a criterion to select hyper-parameter,  $\beta^*$ , which performed well across the datasets.

$$\beta^* = \arg \max_{\beta} \sum_{\text{dataset}} \frac{\text{ACC}(\beta, \text{dataset})}{\text{ACC}(\beta_{\text{dataset}}^*, \text{dataset})}, \quad (19)$$

where  $\beta_{\text{dataset}}^*$  is the best hyper-parameter for the dataset, and  $\text{ACC}(\beta, \text{dataset})$  is the clustering accuracy when hyper-parameter  $\beta$  is used for the dataset. According to the criterion, we set the decay rates to 0.005 for both Linear RIM and Deep RIM, and  $\lambda$  to 0.1 for IMSAT.

## G. Experimental results on hash learning with 32-bit hash codes

Table 5 lists the results on hash learning when 32-bit hash codes were used. For the experiments involving 16-bit hash codes, we observe that IMSAT with the largest network sizes (400-400) exhibited competitive performance in both datasets. The performance of IMSAT improved significantly when we used slightly larger networks (200-200), while the performance of Deep RIM did not improve much with the larger networks.

## H. Comparisons of hash learning with different regularizations and network sizes using toy dataset

We used a toy dataset to illustrate that IMSAT can benefit from larger networks sizes by better modeling the local invariance. We also illustrate that weight-decay does not benefit much from the increased flexibility of neural networks.

Table 5. Comparison of hash performance for 32-bit hash codes (%). Experimental results of Deep Hash and the previous methods are excerpted from [Erin Liong et al. \(2015\)](#).

| Method<br>(Network dimensionality)                   | Hamming ranking (mAP) |                   | precision @ sample = 500 |                   | precision @ r = 2 |                   |
|--|-----------------------|-------------------|--------------------------|-------------------|-------------------|-------------------|
|  | MNIST                 | CIFAR10           | MNIST                    | CIFAR10           | MNIST             | CIFAR10           |
| Spectral hash ( <a href="#">Weiss et al., 2009</a> ) | 25.7                  | 12.4              | 61.3                     | 19.7              | 65.3              | 20.6              |
| PCA-ITQ ( <a href="#">Gong et al., 2013</a> )        | 43.8                  | 16.2              | 74.0                     | 25.3              | 73.1              | 15.0              |
| Deep Hash (60-30 or 80-50)                           | 45.0                  | 16.6              | 74.7                     | 26.0              | 73.3              | 15.8              |
| Linear RIM   | 29.7 (0.4)            | <b>21.2 (3.0)</b> | 68.9 (0.9)               | 16.7 (0.8)        | 60.9 (2.2)        | 15.2 (0.9)        |
| Deep RIM (80-50)                                     | 34.8 (0.7)            | 14.2 (0.3)        | 72.7 (2.2)               | 24.0 (0.9)        | 72.6 (2.1)        | 23.5 (1.0)        |
| Deep RIM (200-200)                                   | 36.5 (0.8)            | 14.1 (0.2)        | 76.2 (1.7)               | 23.7 (0.7)        | 75.9 (1.6)        | 23.3 (0.7)        |
| Deep RIM (400-400)                                   | 37.0 (1.2)            | 14.2 (0.4)        | 76.1 (2.2)               | 23.9 (1.3)        | 75.7 (2.3)        | 23.7 (1.2)        |
| <b>IMSAT (VAT)</b> (80-50)                           | 55.4 (1.4)            | 20.0 (5.5)        | 87.6 (1.3)               | 23.5 (3.4)        | 88.8 (1.3)        | 22.4 (3.2)        |
| <b>IMSAT (VAT)</b> (200-200)                         | 62.9 (1.1)            | 18.9 (0.7)        | 96.1 (0.6)               | 29.8 (1.6)        | 95.8 (0.4)        | <b>29.1 (1.4)</b> |
| <b>IMSAT (VAT)</b> (400-400)                         | <b>64.8 (0.8)</b>     | 18.9 (0.5)        | <b>97.3 (0.4)</b>        | <b>30.8 (1.2)</b> | <b>96.7 (0.6)</b> | <b>29.2 (1.2)</b> |

For the experiments, we generated a spiral-shaped dataset, each arc containing 300 data points. For IMSAT, we used VAT regularization and set  $\epsilon = 0.3$  for all the data points. We compared IMSAT with Deep RIM, which also uses neural networks but with weight-decay regularization. We set the decay rate to 0.0005. We varied three settings for the network dimensionality of the hidden layers: 5-5, 10-10, and 20-20.

Figure 5 shows the experimental results. We see that IMSAT (VAT) can model the complicated decision boundaries by using the increased network dimensionality. On the contrary, the decision boundaries of Deep RIM do not adapt to the non-linearity of data even when the network dimensionality is increased. This observation may suggest why IMSAT (VAT) benefited from the large networks in the benchmark datasets, while Deep RIM did not.

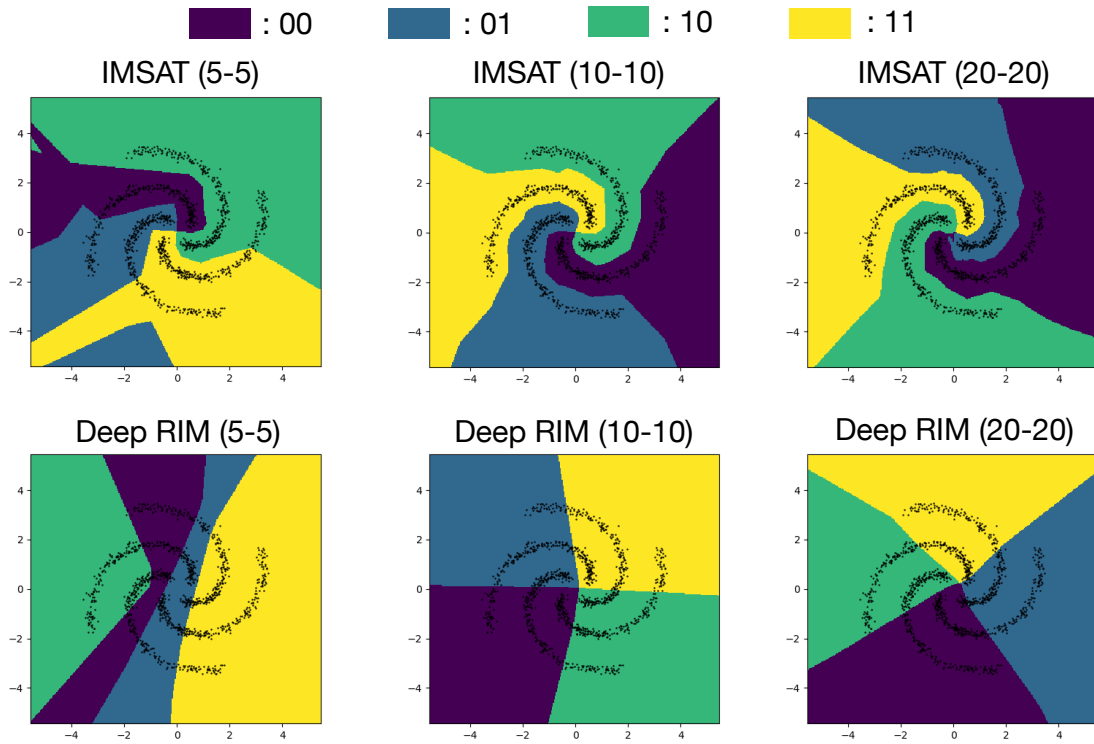


Figure 5. Comparisons of hash learning with the different regularizations and network sizes using toy datasets.