

# Self-attention Multi-view Representation Learning with Diversity-promoting Complementarity

Jian-wei Liu<sup>1</sup>, Xi-hao Ding<sup>1</sup>, Run-kun Lu<sup>1</sup>, Xionglin LUO<sup>1</sup>

1. Department of automation, School of Information Science and Engineering, China University of Petroleum, Beijing, 102249  
E-mail: liujw,luo xl@cup.edu.cn, 964465871@qq.com, zslrk@gmail.com

**Abstract:** Multi-view learning attempts to generate a model with a better performance by exploiting the consensus and/or complementarity among multi-view data. However, in terms of complementarity, most existing approaches only can find representations with single complementarity rather than complementary information with diversity. In this paper, to utilize both complementarity and consistency simultaneously, give free rein to the potential of deep learning in grasping diversity-promoting complementarity for multi-view representation learning, we propose a novel supervised multi-view representation learning algorithm, called Self-Attention Multi-View network with Diversity-Promoting Complementarity (SAMVDPC), which exploits the consistency by a group of encoders, uses self-attention to find complementary information entailing diversity. Extensive experiments conducted on eight real-world datasets have demonstrated the effectiveness of our proposed method, and shows its superiority over several baseline methods, which only consider single complementary information.

**Key Words:** Multi-view Learning, Self-attention Mechanism, Complementary Information with Diversity

## 1 INTRODUCTION

Aiming to make good use of the information from multi-view data and improve the generalization performance, multi-view learning algorithms have made great progress in different tasks, such as classification, regression, and clustering, by utilizing conventional machine learning or deep learning to fully considering the relationships among multiple views [1, 2, 3, 4]. And recently, [5] analyzes these various algorithms, comes to the conclusion that there are two fundamental assumptions ensuring their success: consistency and complementarity principles. The consistency assumption suggests there is consistent information shared by all views, while the complementarity assumption states each view of multi-data may contain some knowledge that other views do not have. Based on these two assumptions, we review the literature of multi-view learning in recent years, and observe that there are still two drawbacks in many state-of-the-art multi-view learning algorithms.

First, at present, multi-view algorithms can be generally categorized into two types: the first category aims to exploit the consistency, the second one aims to leverage the complementarity among multiple views, and each category only focuses on consensus or complementarity. In detail, the first category usually tries to extract the common latent representation on which all views have minimum disagreement, such as canonical correlation analysis (CCA) class algorithms [6, 7, 8, 9, 10], which project two or more views into latent subspaces by maximizing the correlations among projected views, matrix factorization based methods [11, 12, 13], which jointly factorize multi-view data into one common centroid representation by minimizing the overall reconstruction loss of different views. And the

second category is to explicitly preserve complementary information of different views, such as co-training style algorithms [14, 15, 16, 17], which iteratively train two classifiers on two different views, and each classifier generates its complementary information to help the other classifier to train in the next iteration.

However, both consistency and complementarity of multi-view data are meaningful, the neglect of each aspect will result in the loss of valuable information. In order to address this drawback, multi-view algorithms recently began to develop the third category algorithm, which exploits the consistency and complementarity, simultaneously, such as matrix factorization based methods [18, 19], which find latent representations composed of common latent factors shared by multiple views and the specific latent factor of each view. But [18, 19] also inherit the shortcomings of matrix factorization, such as they only learn a linear map relationships, can't reflect the non-linear relationship in the multi-view dataset, and require feed all data in one time, lack the ability of dealing with large scale data.

Second, in terms of complementarity, most of existing multi-view learning algorithms only can find representations with single complementarity rather than complementary information with diversity. Srivastava and Salakhutdinov [20] propose a deep multi-modal RBM to capture the joint distribution.

et al. [21] concatenate the neural hidden coding of audio and video modalities as input, then map these inputs to a shared representation layer. Cho et al. [22] directly input multi-view sequence into RNN encoder to integrate the complementarity of multi-view data. Su et al. [23] introduce a multi-view CNN architecture that integrates complementarity among multiple 2D views of an object into a

single and compact representation by a view-specific pooling layer, which performs element-wise maximum operation across the views. In general, all the above algorithms focus on one type of complementarity among multiple views, and they don't consider mining the complementary information with diversity.

To fight against the above mentioned serious deficiencies, in this paper, we propose a new multi-view learning paradigm based on self-attention network, called Self-Attention Multi-View network with Diversity-Promoting Complementarity (SAMVDPC). Specifically, SAMVDPC firstly encodes each view's data into a fixed-length vector representation to exploit the consistency, and then explores complementary information entailing diversity with multiple combination forms by self-attention mechanism, finally concatenates all complementary information into a vector representation, which further be used to make prediction.

We illustrate this idea using an example from a face recognition problem with two views. Given a group of people, we have collected the face information for each person to form two-view dataset. To make classification, firstly, by building a unique encoder for each view, SAMVDPC encodes each view's data into a fixed-length vector representation and outputs  $\mathbf{H} = [h_1; h_2] \in \mathcal{R}^{2 \times H}$ . Second, SAMVDPC inputs  $\mathbf{H}$  to self-attention mechanism to produce weight matrix  $\mathbf{W} = [w_1; w_2] \in \mathcal{R}^{2 \times 2}$ , then outputs two vectors:  $w_1\mathbf{H}$ , and  $w_2\mathbf{H}$ , which can utilize to combine two-view data by different ways for subsequent fusion stage. Finally, SAMVDPC incorporates the concatenated representation  $[w_1\mathbf{H}, w_2\mathbf{H}]$  as input and processes these inputs by a forward network to make prediction.

In summary, our contributions we summarize are shown as follows:

- (1) We develop a supervised multi-view deep learning algorithm, which utilizes both consistency and complementarity of multiple views, here multiple view's encoders consider the consistency, and self-attention mechanism considers the complementarity.
- (2) Compared to [18, 19], encoders in SAMVDPC can learn nonlinear and hierarchical abstract feature representation for multi-view data, which capture the non-linear relationship and real underlying properties in multi-view dataset.
- (3) SAMVDPC can find representations with complementary information possessing diversity rather than single complementarity, and sufficiently reflect the complementarity underlying multi-view data.
- (4) We have compared SAMVDPC with other state-of-the-art multi-view learning algorithms and demonstrated its effectiveness, that's more, we also build other baselines deep networks to further analyze SAMVDPC's performance, which explore single complementary by mean-pooling, max-pooling and weighted summation.

### Attention mechanism

In deep neural networks, attention mechanism [24] has been developed in the context of encoder-decoder architectures for Neural Machine Translation (NMT) [22, 25], and rapidly applied to numerous application domains and

achieved promising results on several challenging tasks, such as image captioning [26], and summarization [27]. Besides, with the development of deep learning (deep learning) of attention mechanism, self-attention mechanism,

which

utilizes a single sequence to compute content-

self-representation [5-1749] (of 5138) (the 5138) (sequence of 859) (and 5138) (has)

complementarity, abstractive summarization, and learn-

ing task-independent representations, and

arcs (ers)-287 (deceon)-287 (thirn)-286 knholledge of attention mecn-

for encoder-decoder the

and attention mechanism as follows: *Attention*

Fig. 1(a), the architecture of SAMVDPC is made up of Encoder-Block, self-attention mapping, and full connected layer, and the detailed self-attention mapping processes are shown in Fig. 1(b). We describe each of the constituents in the following subsections.

**ncoder-loc** : As shown in Fig. 1(a), Encoder-Block is composed with  $V$  same encoders to extract each view's feature. The initial model parameters of each encoder are initialized by the encoder of a corresponding auto-encoder, which will be explained more detailed in section 4.4. From these encoders,  $V$  hidden features ( $\mathbf{z}^v \in \mathbb{R}^{H \times 1}, v = 1, \dots, V$ ) can be obtained and they will be stacked horizontally and combined into a feature matrix:  $\mathbf{Z} = [\mathbf{z}^1, \dots, \mathbf{z}^V], \mathbf{z}^v \in \mathbb{R}^{H \times 1}$ , here  $H$  is the number of dimension of hidden feature vector  $\mathbf{z}^v$ .

**Self-attention mapping**: The self-attention mechanism takes the whole hidden states matrix  $\mathbf{Z}$  as input, outputs a matrix  $\mathbf{A}$ , and each row of  $\mathbf{A}$  is a vector of weights  $\mathbf{a}_i$ :

$$\mathbf{A} = \left[ \mathbf{a}_1; \dots; \mathbf{a}_{d_s} \right] = \text{softmax}(\mathbf{W}_{s2} \tanh(\mathbf{W}_{s1} \mathbf{Z}^T)), \quad (2)$$

here,  $\mathbf{A} \in \mathbb{R}^{d_s \times V}$ ,  $\mathbf{W}_{s1} \in \mathbb{R}^{d_s \times H}$ ,  $\mathbf{W}_{s2} \in \mathbb{R}^{d_s \times d_s}$ ,  $d_s$  is a hyperparameter we can set arbitrarily, and the softmax () is operated along the second dimension of its input. Inspired by [30], Equation (2) also can be deemed as a 2-layer MLP without bias, whose hidden unit numbers are  $d_s$ , and parameters are  $\{\mathbf{W}_{s2}, \mathbf{W}_{s1}\}$ . Finally, we compute the  $d_s$  weighted sums by multiplying  $\mathbf{A}$  and  $\mathbf{Z}$ :

$$\mathbf{M} = \mathbf{AZ}^T, \mathbf{M} \in \mathbb{R}^{d_s \times H}. \quad (3)$$

It is worth noting that each row of  $\mathbf{M}$  is a unique nonlinear combination of multiple views data, and the self-attention mechanism outputs formulate  $d_s$  kinds of nonlinear combination of multiple views data. In our experiment part, the value of  $d_s$  is set to  $V$ .

**ML** : We concatenate each row of  $\mathbf{M}$  to produce a multi-view representation containing multiple combinations of multi-view data to extract complementary information entailing diversity. Then we input this representation to 2-layer MLP, and make prediction.

#### 4.2. Objective function and Regularization

The embedding matrix  $\mathbf{M}$  always suffer from redundancy problems because the self-attention mechanism often provides similar summation weights for all the  $d_s$  hops. Inspired by [1], we also add regularization to encourage the diversity of summation weights vectors across different hops of attention. Thus, in this paper, our objective function is consist of cross entropy loss and regularization, and can be formulated as follows:

$$L = \text{cross\_entropy}(y, \hat{y}) + \|\mathbf{AA}^T - \mathbf{I}\|_F^2, \quad (4)$$

here  $\lambda$  is regularization parameters, and  $\mathbf{I}$  is a unit diagonal matrix.

#### Experiment

In this section, we experimentally evaluate SAMVDPC in classification task on eight real world multi-view data sets

Table 1: Characteristics of the datasets

Data Set	Characteristics			
	Instances numbers	K		Dimension numbers
Leaves	96	3	6	64 for all
Reuters	1200	5	6	2000 for all
YaleFace	256	2	8	2016 for all
BBC	685	4	5	4659/4633/4665/4684
Cornell	195	2	5	1703/585
Texas	187	2	5	1703/561
Washington	230	2	5	1703/690
usconsin	265	2	8	1703/795

by comparing it to other baseline algorithms, and design a set of exploratory experiments to validate properties of the self-attention mechanism in SAMVDPC, finally analyse the convergence of our proposed algorithm.

#### 4.1 Datasets

In this paper, we use eight real-world multi-view data sets to verify the performance of SAMVDPC, including Leaves, Reuters, YaleFace, BBC, Cornell, Texas, Washington, and usconsin datasets. Leaves and YaleFace are two image dataset, Reuters and BBC are two text dataset, Cornell, Texas, Washington, and usconsin dataset are four subset of data sets selected from webB data sets, and webB are webpage dataset. The properties of data sets are summarized in Table 1.

#### 4.2 Comparison Algorithms and Baseline Models

We evaluate the SAMVDPC performance in classification tasks by comparing it with several state-of-the-art multi-view learning algorithms based on matrix factorization, such as MVNMF [32], multiNMF [12], MVCC [13], DICS [18], and some our designed deep neural network baseline models with three sorts of fusion strategies replacing with our self-attention mechanism, including max-pooling model, mean-pooling model, and weighted summation model. For fair comparison, in terms of matrix factorization algorithms, we choose the parameters within the range that author suggested to obtain good latent representations, and input these representations to NN( $k = 1$ ) for classification in terms of deep neural network baseline models, we instead of the self-attention mechanism with max-pooling, mean-pooling, or weighted summation fusion, maintain the remaining structure unchanged, and remove the regularization in our objective function.

MVNMF is an NMF-based algorithm by merging local geometrical structure information of each view in a multi-view feature extraction framework. The extracted feature considered the inner-view relatedness between data, and further can be used to complete various tasks. We select parameters  $\lambda, \mu$  to 0.01, and 10 as author suggested, respectively.

MultiNMF is an NMF-based multi-view algorithm, in terms of matrix factorization, it requires coefficient matrices learnt from different views to be softly regularized towards a common consensus matrix, which reflect the information of multi-view data and can be used to make clas-

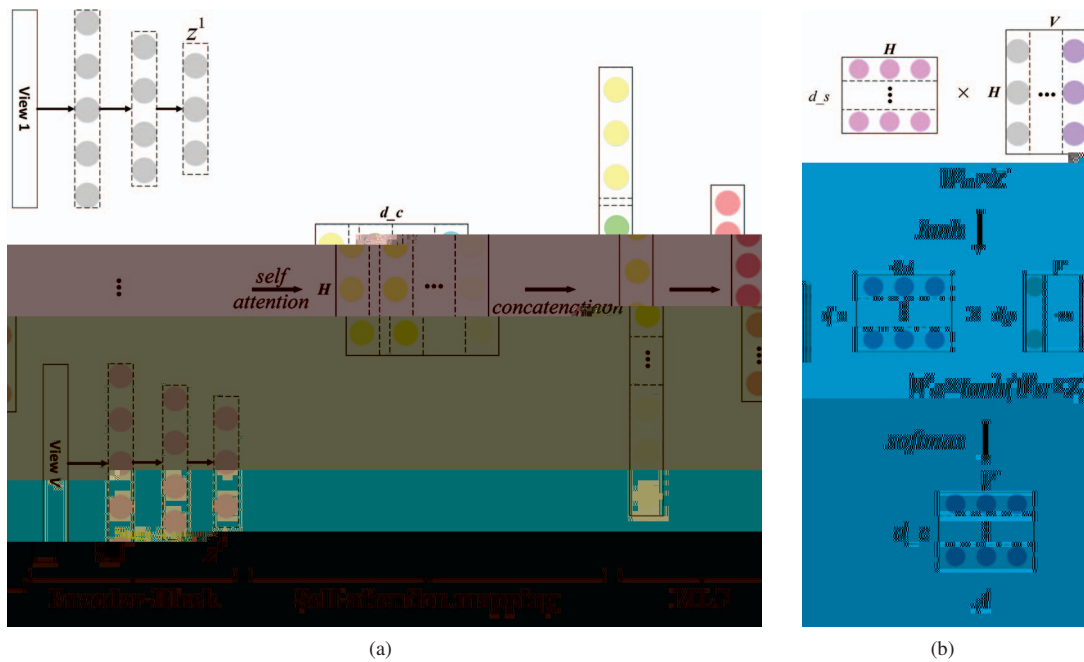


Figure 1: MVCapsNet Architecture. Fig. 1: (a) is the architecture of SAMVDPC. (b) is concrete self-attention mapping implementation processes.

Table 2: H perparameters on each data set

Data Set	Units number of each encoder layer			Diversity of complementarity	Size of mini-batch
	$l_1$	$l_2$	$l_3$	$d_c$	
Leaves	64	32	16	3	4
Reuters	2048	1024	512	2	16
aleFace	1024	512	512	2	32
BBC	1024	512	512	4	32
Cornell	1024	512	128	2	16
Te as	1024	512	128	2	16
ashington	1024	512	128	2	16
isconsin	1024	512	128	2	16

Table 3: Accuracy of different methods

Method	ACC(%)							
	Leaves	ale ace	Reuters	C	Cornell	Te as	Washington	Wisconsin
NM	95.0±0	50.0±2.5	40.8±1.2	38.0±1.5	41.0±1.8	57.9±1.8	69.6±2.2	52.8±1.4
MultiNM	95.0±0	64.2±4.2	52.7±0.2	73.1±0.2	49.7±7.7	68.7±3.4	59.3±2.6	50.3±3.5
M CC	100±	33.3±6.9	54.4±1.9	77.7±.7	60.8±5.0	64.7±5.5	62.8±3.8	64.3±2.7
DICS	97.9±2.5	89.1±3.2	70.3±4.0	90.2±2.4	78.7±.1	81.6±4.0	78.7±.0	85.1±4.5
MA -ooling	100±0	90.0±4.6	71.2±3.4	80.5±7.6	71.3±8.7	74.7±5.2	67.5±8.2	86.2±7.7
M AN-ooling	100±0	90.6±5.3	71.2±4.3	83.3±6.7	70.9±5.5	76.6±4.0	70.0±4.9	84.8±5.2
Weighted Sum	100±0	92.9±4.8	77.7±.7	87.2±4.5	72.5±13	76.3±4.9	66.9±7.8	84.8±5.2
AM D. C	100±0	70.0±.7	70.0±5.2	93.5±2.4	72.2±4.9	77.7±.0	75.0±6.1	84.0±5.2

sification. We select the values of regularization parameter are  $10^{-3}$ ,  $10^{-2}$ ,  $10^{-1}$ , and 1.

MVCC is a novel multi-view method based on concept factorization with local manifold regularization, which also drives a common consensus representation for multiple views. We set parameter  $\lambda$  to 100, and both select the values of parameters  $\alpha$  and  $\beta$  are 50, 100, 200, 500, and 1000. DICS is an NMF-based multi-view learning algorithm, by exploring the discriminative and nondiscriminating information existing in common and view-specific parts among different views via joint non-negative matrix factorization, and produce discriminative and non-discriminative feature from all subspaces. And then, discriminative and nondiscriminative features are further used to produce classification results. We select parameters  $\alpha$  and  $\beta$  within a small range of  $[0, 1]$ , and set parameter  $\lambda$  to 1.

Due to no publicly available multi-view clustering algorithm based on deep neural network, we generate three baseline models based on deep neural network. These baseline models are explorer models to validate properties of the self-attention mechanism in SAMVDPC, the separately use max pooling, mean pooling, and weighted summation to fusion all multiple views representations produced by Encoder-Block, and a fusion representation with single complementarity, and then input the fusion representation to fully connected layer to make prediction. And for fair comparison, we use the same settings for baseline models as that we did in SAMVDPC.

### Configuration and Tricks

In this subsection, we specify the configuration of SAMVDPC. In Encoder-Block, structures of all encoders are the same, each encoder has one input layer and three hidden layers  $l_1$ ,  $l_2$ , and  $l_3$ , the number of number in each hidden layer decrease as the layers of encoder deepens, and the activation function of all hidden layers is ReLU. In self-attention mapping, self-attention MLP has a hidden layer with 300 units  $d_s$ , and we all always choose the matrix embedding to have  $V$  rows ( $d_c$ ). In MLP, we use a 2-layer ReLU output MLP with 512 hidden states to output the classification result. For objective function, we usually set  $\lambda$  to 0.0001. For the configuration of three baseline models, we use max-pooling, mean-pooling, or weighted summation to take replace of the self-attention mechanism in SAMVDPC, and set  $\lambda$  in objective function to 0. The hyper parameters on each data sets are summarized in Table 2.

With regard to the initialization of SAMVDPC weights, In Encoder-Block, we pre-train  $V$  auto-encoders through minimizing the reconstruction error of each view, and then use the pre-trained parameters of auto-encoders to initialize the corresponding encoder's weight of Encoder-Block. In self-attention MLP and MLP, Xavier is used as the weight initialization method [15].

In training process, the optimizer algorithm we used is Adam, the learning rate is always initialized to  $10^{-3}$ ,  $10^{-4}$ ,  $10^{-5}$ , and will decrease gradually with the development of training process. To avoid overfitting, all layers in Encoder-Block and MLP are regularized by dropout reg-

ularization in training process, and dropout rate is set to 0.5.

### Result

All datasets divided into training, verification and testing data in a ratio of 0.6:0.2:0.2. For SAMVDPC comparison algorithms, and baseline models, we first run each model on each dataset to select hyper parameters that has the best accuracy and generalization performance. And then based on these hyper parameters, we run all algorithms 10 times on each dataset and report the mean values and standard deviation of accuracies.

All the classification results of eight multi-view datasets are summarized in Table 3, and the best result on each dataset is highlighted in boldface. As we can see, the proposed SAMVDPC achieves best accuracy on Teas and Yale-Face datasets, and are comparable with other algorithms on the else datasets. The promising result mainly from four aspects: (1) DICS, baseline models, and SAMVDPC are all algorithms exploiting the consistency and complementarity, simultaneously, and compared to NMF, multi-NMF, and MVCC, they all achieve better performance on all datasets (2) compared to matrix factorization algorithms, the Encoder-Block in both baseline models and SAMVDPC can extract features in a way of effectively fetching consistent information and grasping the underlying common properties of multi-view datasets (3) compared to baseline models, the complementarity with diversity exploited by self-attention mechanism contains more information than max-pooling, mean-pooling, and weighted summation.

### Convergence Analysis of Training Process

In order to empirically investigate the convergence property of SAMVDPC, we plot the iterative curves of objective function and the corresponding classification accuracies on three typical data sets, Leaves, BBC, and Teas in Fig. 2. From Fig. 2, we can observe that: (1) the objective function values drop sharply and meanwhile the classification accuracies increase rapidly within the previous rounds of iterative process, and then the objective function and the accuracy curves begin to decrease/grow mildly, finally converge to a value or fluctuate around a constant (2) with respect to convergence speed, the objective function values of SAMVDPC converge in the least iterations, in contrast, max-pooling corresponds to the most iterations, because max-pooling operation is loss compression process and the backpropagation process doesn't make full use of information from multiple views data (3) in respect of convergence result, the objective function of SAMVDPC can finally converge to a fixed value on every dataset, but the objective function of baseline models always fluctuate around a constant, that's more, compared to baseline models, we can find that the classification accuracy curves of SAMVDPC often fluctuate within a narrow range. In conclusion, compared to baseline models, SAMVDPC get a better performance on the iterative curves of objective function and the corresponding classification accuracy.

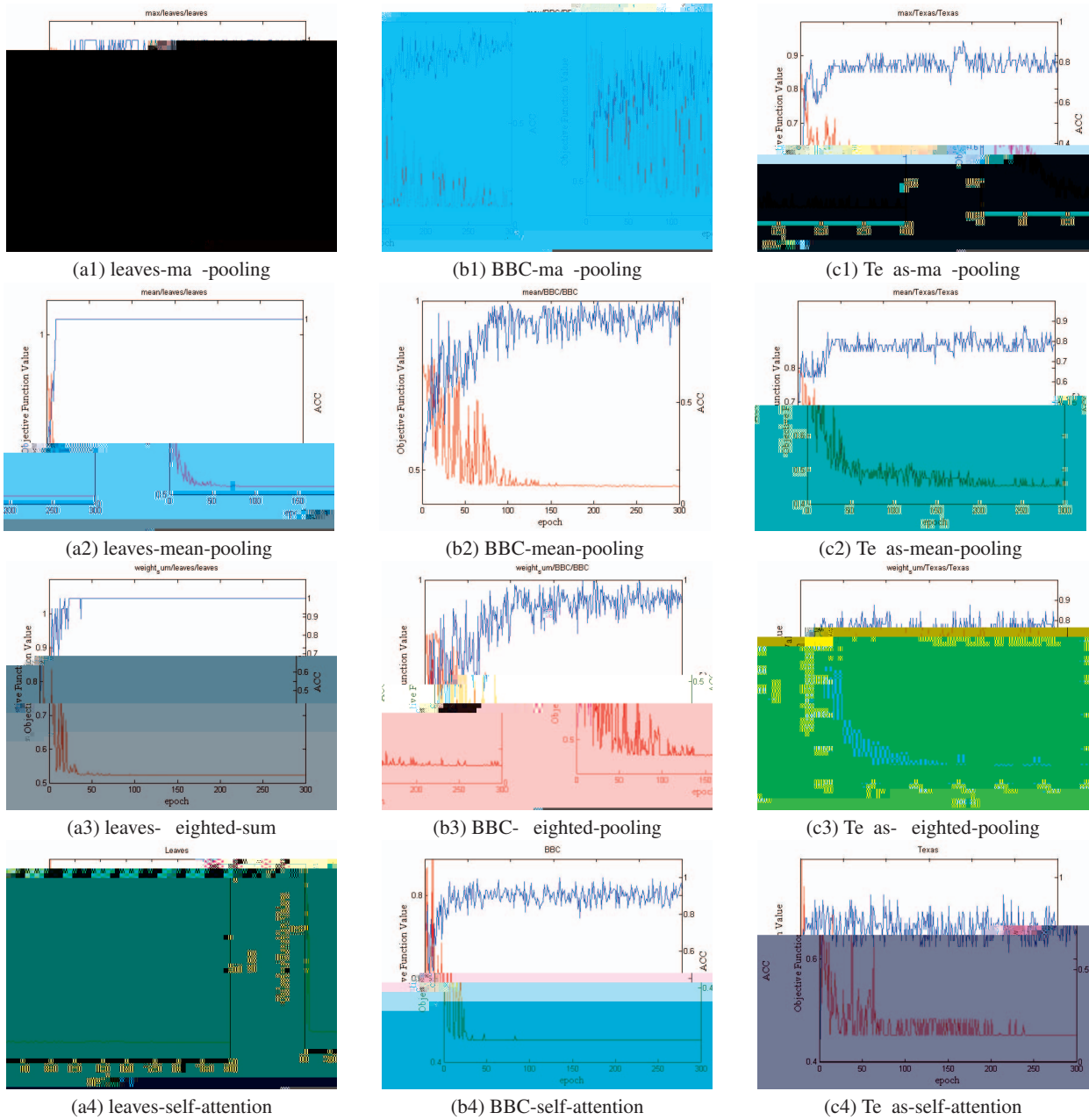


Figure 2: The convergence property of SAMVDP

## Conclusion and Future Work

In this paper, we propose a novel multi-view network,

- [2] Y. Li, B. Feng, J. Hua, D. Tao, L. Yang, and C. Xu. 2011. Difficult guided image retrieval using linear multiview embedding. In 19th ACM Multimedia Proceedings, 1169-1172.
- [3] C. Gan, R. Pan, and J. Li. 2011. Bi-directional domain adaptation for cross-language text classification. In 22nd IJCAI Proceedings, 1535-1540.
- [4] B. Xie, Y. Mu, D. Tao, and S. Sridharan. 2011. Huang, m-sne: Multiview stochastic neighbor embedding. *IEEE Trans. Systems, Man, and Cybernetics*, 41(4):1088-1096.
- [5] Ch. Xu, D. Tao, and C. Xu. 2013. A Survey on Multi-view Learning. arXiv preprint arXiv, 1304.5634.
- [6] S. Chaudhuri, S. Mukherjee, S. Livescu, and S. Sridharan. 2009. Multi-view clustering via canonical correlation analysis. In 26th ICML Proceedings, 129-136.
- [7] J. D. R. Farquhar, D. R. Hardoon, H. Meng, J. Shave-Talor, and S. Szepesvári. 2005. Transfer learning: SVM-2, Theory and Practice. In 18th NIPS Proceedings, 355-362.
- [8] D. R. Hardoon, S. Szepesvári, and J. Shave-Talor. 2004. Canonical correlation analysis: an overview with application to learning methods. *Neural Computation*, 16(12):2639-2664.
- [9] M. Fan, S. Shan, H. Liang, S. Lao, and X. Chen. 2016. Multi-view discriminant analysis. *IEEE Trans. Pattern Anal. Mach. Intell.*, 38(1):188-194.
- [10] A. Sharma, A. Kumar, H. Daum<sup>3</sup>, D. and S. Jacobs. 2012. Generalized multiview analysis: a discriminative latent space. In 20th CVPR Proceedings, 2160-2167.
- [11] Y. Yuan, L. Huang, J. Peng, and J. Fan. 2015. Multi-view concept learning for data representation. *IEEE Trans. on Data Eng.*, 27(11): 3016-3028.
- [12] J. Yao, J. Han, J. Liu, and C. Yang. 2013. Multi-view clustering via joint nonnegative matrix factorization. In 13th SDM Proceedings, 252-260.
- [13] H. Yang, Y. Yang, and T. Li. 2016. Multi-view clustering via concept factorization with local manifold regularization. In 16th ICDM Proceedings, 1245-1250.
- [14] A. Blum, T. M. Mitchell. 1998. Combining labeled and unlabeled data with co-training. In 16th COLT Proceedings, 92-100.
- [15] A. Kumar, and H. Daum<sup>3</sup>. 2011. A co-training approach for multi-view spectral clustering. In 28th ICML Proceedings, 393-400.
- [16] C. Yang, B. Hou. 2010. A new analysis of co-training. In 27th ICML Proceedings, 1135-1142.
- [17] M. Yang, B. Hou. 2011. CoTrade: consistent co-training with data editing. *IEEE Trans. Systems, Man, and Cybernetics*, 41(6):1612-1626.
- [18] C. Yang, Y. Yin, P. Li, Y. Yang, and J. Shao. 2018. Multi-view discriminative learning via joint non-negative matrix factorization. In 23rd DASFAA Proceedings, 542-557.
- [19] A. P. Singh, S. J. Gordon. 2008. Relational learning via collective matrix factorization. In 14th UDD Proceedings, 650-658.
- [20] N. Srivastava, R. Salakhutdinov. 2014. Multimodal learning with deep Boltzmann machines. *Journal of Machine Learning Research*, 15(1):2949-2980.
- [21] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Ng. 2011. Multimodal Deep Learning. In 14th UDD Proceedings, 689-696.
- [22] K. Cho, B. V. Merriënboer, D. Sussan, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. 2014. Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. In 2014 EMNLP Proceedings, 1724-1734.
- [23] H. Su, S. Maji, E. Sifert, and E. Learned-Miller. 2015. Multi-view Convolutional Neural Networks for 3D Shape Recognition. In 2015 ICCV Proceedings, 945-953.
- [24] D. Bahdanau, K. Cho, and Y. Bengio. 2015. Neural Machine Translation by Jointly Learning to Align and Translate. In 3rd ICLR Proceedings.
- [25] I. Sutskever, O. Vinyals, and P. V. Le. 2014. Sequence to Sequence Learning with Neural Networks. In 28th NIPS Proceedings, 3104-3112.
- [26] K. Xu, J. Ba, R. Kiros, K. Cho, A. C. Courville, R. Salakhutdinov, R. S. Smeed, and Y. Bengio. 2015. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. In 32nd ICML Proceedings, 2048-2057.
- [27] A. M. Rush, S. Chopra, and J. Weston. 2015. A Neural Attention Model for Abstractive Sentence Summarization. In 2015 EMNLP Proceedings, 379-389.
- [28] J. Cheng, L. Dong, and M. Lapata. 2015. Long Short-Term Memory-Net for Machine Reading. In 2016 EMNLP Proceedings, 551-561.
- [29] A. P. Parikh, O. Tackstam, D. Das, and J. Uszkoreit. 2016. A Decomposable Attention Model for Natural Language Inference. In 2016 EMNLP Proceedings, 2249-2255.
- [30] J. Lin, M. Feng, C. N. d. Santos, M. Yu, B. Xiang, B. Hou, and Y. Bengio. 2017. A Structured Self-Attentive Sentence Embedding. In 5th ICLR Proceedings.
- [31] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. 2017. Attention is All you Need. In 31st NIPS Proceedings, 6000-6010.
- [32] C. Yang, X. Gong, H. Fu, M. Li, and Y. Huang. 2015. Feature extraction via multi-view non-negative matrix factorization with local graph regularization. In 2015 ICIP Proceedings, 3500-3504.
- [33] X. Glorot, Y. Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In 13th AISTATS Proceedings, 249-256.
- [34] C. Finn, P. Abbeel, and S. Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In 34th ICML Proceedings, 1126-1135.