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

Jian-wei Liu[1], Xi-hao Ding[1], Run-kun Lu[1], Xionglin LUO[1]

1. Department of automation, School of Information Science and Engineering, China University of Petroleum, Beijing, 102249

E-mail: liujw,luo_l@cup.edu.cn, 964465871@qq.com, zs_lrk@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 show 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-views 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. Ngiam et al. [21] concatenate the final 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 b  a vie  -pooling la er,  hich performs element-  ise ma imum operation across the vie  s. In general, all the above algorithms focus on one t pe of complementarit  among multiple vie  s, and the   don't consider mining the complementar  information  ith diversit .

To   ght against the above mentioned serious de ciencies, in this paper,  e propose a ne   multi-vie   learning paradigm based on self-attention net ork, called Self-Attention Multi-Vie   net ork   ith Diversit -Promoting Complementarit  (SAMVDPC). Speciall , SAMVDPC  rst encodes each vie  's data into a    ed-length vector representation to e  ploit the consistenc , and then e  plores complementar  information entailing diversit    ith multiple combination forms b   self-attention mechanism,   nall concatenates all complementar  information into a vector representation,   hich further be used to make prediction.

 e ma   illustrate this idea using an e  ample from a face recognition problem   ith t  o vie  s.  iven a group of people,  e have collected the face information for each person to form t  o-vie   dataset. To make classi cation,  rst, b   building a unique encoder for each vie  , SAMVDPC encodes each vie  's data into a    ed-length vector representation and outputs $\mathbf{H} = [h_1; h_2] \in R^{2 \times H}$. Second, SAMVDPC inputs $\mathbf{H}$ to self-attention mechanism to produce   eight matri  $\mathbf{W} = [w_1; w_2] \in ^{2 \times 2}$, then outputs t  o vectors: $w_1\mathbf{H}$, and $w_2\mathbf{H}$,   hich can utilize to combine t  o-vie   data b   different   a s for subsequent fusion stage. Finall , SAMVDPC incorporates the concatenated representation $[w_1\mathbf{H}, w_2\mathbf{H}]$ as input and processes these inputs b   a for  ard net  ork to make prediction.

In summar , our contributions  e summarize are sho  n as follo  s:

(1)  e develop a supervised multi-vie   deep learning algorithm,   hich utilizes both consistenc  and complementarit  of multiple vie  s,   here multiple vie  s' encoders consider the consistenc , and self-attention mechanism considers the complementarit ,

(2) Compared to [18, 19], encoders in SAMVDPC can learn nonlinear and hierarchical abstract feature representation for multi-vie   data,   hich capture the non-linear relationship and real underl ing properties in multi-vie   dataset.

(3) SAMVDPC can  nd representations   ith complementar  information possessing diversit  rather than single complementarit , and suf cientl  re ect the complementarit  underl ing multi-vie   data.

(4)  e have compared SAMVDPC   ith other state-of-the-art multi-vie   learning algorithms and demonstrated its effectiveness,   hat's more,  e also build other baselines deep net  orks to further anal ze SAMVDPC's performance,   hich e  plore single complementar  b   mean-pooling, ma  -pooling and   eighted summation.

## ➚ Attention mechanism

In deep neural net  orks, attention mechanism [24] has been developed in the conte t of encoder-decoder architectures for Neural Machine Translation (NMT) [22, 25], and rapidl   applied to numerous application domains and achieved promising results on several challenging tasks, such as image captioning [26], and summarization [27]. Besides,   ith the development of deep learning9(de)-412233(la er)4372(

Fig. 1(a), the architecture of SAMVDPC is made up of Encoder-Block, self-attention mapping, and full connected la er, and the detailed self-attention mapping processes are sho n in Fig. 1(b). e describe each of the constituents in the follo ing subsections.

**ncoder- loc :** As sho n in Fig. 1(a), Encoder-Block is composed ith $V$ same encoders to e tract each vie 's feature. The initial model parameters of each encoder are initialized b the encoder of a corresponding auto-encoder, hich ill be e plained more detailed in section 4.4. From these encoders, $V$ hidden features ($\mathbf{z}^v \in R^{H \times 1}$, $v = 1, \cdots, V$) can be obtained and the ill be stacked horizontall and combined into a feature matri : $\mathbf{Z} = \left[\mathbf{z}^1, \cdots \mathbf{z}^V\right]$, $\mathbf{z}^v \in R^{H \times 1}$, here $H$ is the number of dimension of hidden feature vector $\mathbf{z}^v$.

**Self-attention mapping:** The self-attention mechanism takes the hole hidden states matri $\mathbf{Z}$ as input, outputs a matri $\mathbf{A}$, and each ro of $\mathbf{A}$ is a vector of eights $\mathbf{a}_i$:

$$\mathbf{A} = \left[\mathbf{a}_1 \vdots \mathbf{a}_{d\_c}\right] = softmax\left(\mathbf{W}_{s2} \tanh\left(\mathbf{W}_{s1}\mathbf{Z}^T\right)\right), \quad (2)$$

here, $\mathbf{A} \in R^{d\_c \times V}$, $\mathbf{W}_{s1} \in R^{d\_s \times H}$, $\mathbf{W}_{s2} \in R^{d\_c \times d\_s}$, $d\_s$ is a h per parameter e can set arbitraril , and the softma () is operated along the second dimension of its input. Inspired b [30], Equation (2) also can be deemed as a 2-la er MLP ithout bias, hose hidden unit numbers are $d\_s$, and parameters are $\{\mathbf{W}_{s2}, \mathbf{W}_{s1}\}$. Finall , e compute the $d\_s$ eighted sums b multipl ing $\mathbf{A}$ and $\mathbf{Z}$:

$$\mathbf{M} = \mathbf{A}\mathbf{Z}^T, \quad \mathbf{M} \in R^{d\_c \times H}. \quad (3)$$

It is orth noting that each ro of $\mathbf{M}$ is a unique nonlinear combinations of multiple vie s data, and the self-attention mechanism outputs formulate $d\_s$ kinds of nonlinear combination of multiple vie s data. In our e periment part, the value of $d\_s$ is set to $V$.

**ML :** e concatenate each ro of to produce a multi-vie representation containing multiple combinations of multi-vie data to e tract complementar information entailing diversit . Then e input this representation to 2-la er MLP, and make prediction.

### _._ **Ob ective unction and Regulari ation**

The embedding matri $\mathbf{M}$ al a s suffer from redundanc problems because the self-attention mechanism often provides similar summation eights for all the $d\_s$ hops. Inspired b [ ], e also add regularization to encourage the diversit of summation eights vectors across different hops of attention. Thus, in this paper, our objective function is consist of cross entrop loss and regularization, and can be formulated as follo s:

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

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

### periment

In this section, e e perimentall evaluate SAMVDPC in classi cation task on eight real orld multi-vie 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 |
| XaleFace | 256 | 2 | 8 | 2016 for all |
| BBC | 685 | 4 | 5 | 4659/4633/4665/4684 |
| Cornell | 195 | 2 | 5 | 1703/585 |
| Te as | 187 | 2 | 5 | 1703/561 |
| ashington | 230 | 2 | 5 | 1703/690 |
| isconsin | 265 | 2 | 8 | 1703/795 |

b comparing it to other baseline algorithms, and design a set of e plorator e periments to validate properties of the self-attention mechanism in SAMVDPC, nall anal-se the convergence of our proposed algorithm.

### .1 Datasets

In this paper, e use eight real- orld multi-vie data sets to verif the performance of SAMVDPC, including Leaves, Reuters, XaleFace, BBC, Cornell, Te as, ashington, and isconsin datasets. Leaves and XaleFace are t o image dataset, Reuters and BBC are t o te t dataset, Cornell, Te as, ashington, and isconsin dataset are four subset of data sets selected from eb B data sets, and e-b B are ebpage dataset. The properties of data sets are summarized in Table 1.

### . Comparison Algorithms and aseline Models

e evaluate the SAMVDPC performance in classi cation tasks b comparing it ith several state-of-the-art multi-vie learning algorithms based on matri factorization, such as MVNMF [32], multiNMF [12], MVCC [13], DICS [18], and some our designed deep neural net ork baseline models ith three sorts of fusion strategies replacing ith our self-attention mechanism, including ma -pooling model, mean-pooling model, and eighted summation model. For fair comparison, in terms of matri factorization algorithms, e choose the parameters ithin the range that author suggested to obtain good latent representations, and input these representations to NN($k = 1$) for classi cation, in terms of deep neural net ork baseline models, e instead of the self-attention mechanism ith ma -pooling, mean-pooling, or eighted summation fusion, maintain the remaining structure unchanged, and remove the regularization in our objective function.

MVNMF is an NMF-based algorithm b merging local geometrical structure information of each vie in a multi-vie feature e traction frame ork. The e tracted feature considered the inner-vie relatedness bet een data, and further can be used to complete various tasks. e select parameters $f$, $\mu$ to 0.01, and 10 as author suggested, respectivel .

MultiNMF is an NMF-based multi-vie algorithm, in terms of matri factorization, it requires coef cient matrices learnt from different vie s to be softl regularized to ards a common consensus matri , hich re ect the information of multi-vie data and can be used to make clas-
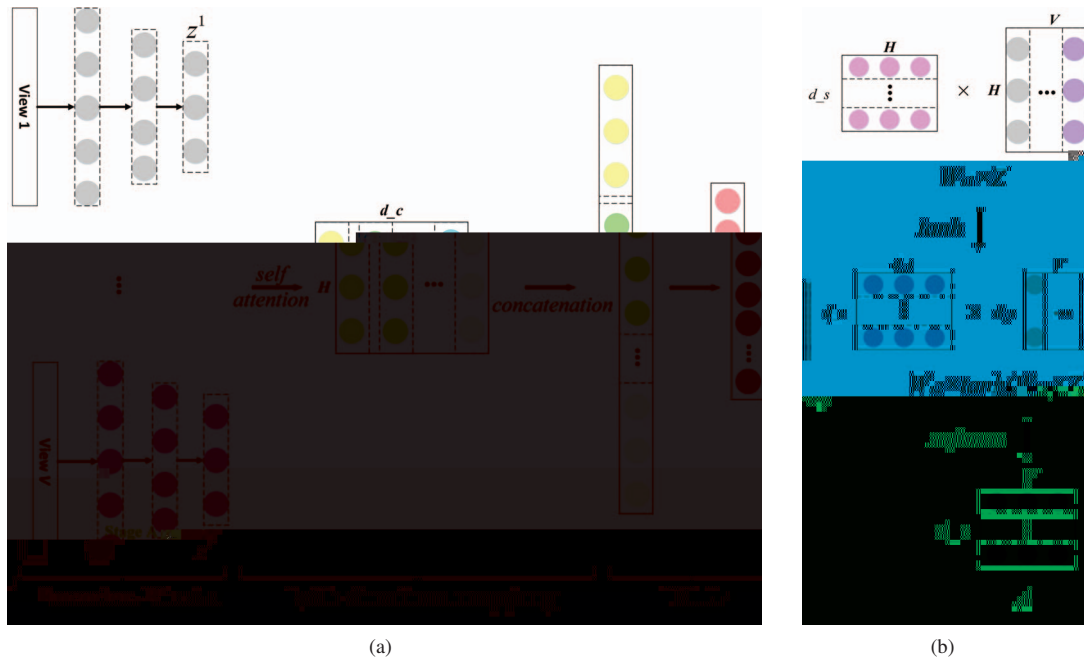
Figure 1: MVCapsNet Architecture. Fig. 1: (a) is the architecture of SAMVDPC. (b) is concrete self-attention mapping implementation processes.

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 |
| YaleFace | 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: Accurac 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 | . ± | 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 | ± .1 | 81.6±4.0 | . ± .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 | ± . | 87.2±4.5 | 72.5±13 | 76.3±4.9 | 66.9±7.8 | . ± . |
| AM D C | 100±0 | .0± | 70.0±5.2 | 93.5±2.4 | 72.2±4.9 | . ± .0 | 75.0±6.1 | 84.0±5.2 |

si cation. e select the values of regularization parameter are $10^{-3}$, $10^{-2}$, $10^{-1}$, and 1.

MVCC is a novel multi-vie method based on concept factorization ith local manifold regularization, hich also drives a common consensus representation for multiple vie s. e set parameter to 100, and both select the values of parameters and are 50, 100, 200, 500, and 1000.

DICS is an NMF-based multi-vie learning algorithm, b e ploring the discriminative and nondiscriminating information e isting in common and vie -speci c parts among different vie s via joint non-negative matri factorization, and produce discriminative and non-discriminative feature from all subspaces. And then, discriminative and nondiscriminative features are further used to produce classi cation results. e select parameters and ithin a small range of $[0, 1]$, and set parameter to 1.

Due to no publicl available multi-vie clustering algorithm based on deep neural net ork, e generate three baseline models based on deep neural net ork. These baseline models are e plorator models to validate properties of the self-attention mechanism in SAMVDPC, the separatel use ma pooling, mean pooling, and eighted summation to fusion all multiple vie s representations produced b Encoder-Block, nd a fusion representation ith single complementarit , and then input the fusion representation to full connected la er to make prediction. And for fair comparison, e use the same settings for baseline models as hat e did in SAMVDPC.

### . Con guration and Tric s

In this subsection, e specif the con guration of SAMVDPC. In Encoder-Block, structures of all encoders are the same, each encoder has one input la er and three hidden la ers $l_1$, $l_2$, and $l_3$, the number of number in each hidden la er decrease as the la ers of encoder deepens, and the activation function of all hidden la ers is ReLU. In self-attention mapping, self-attention MLP has a hidden la er ith 300 units $d\_s$, and e al a s choose the matri embedding to have V ro s ($d\_c$). In MLP, e use a 2-la er ReLU output MLP ith 512 hidden states to output the classi cation result. For objective function, e usuall set to 0.0001. For the con guration of three baseline models, e use ma -pooling, mean-pooling, or eighted summation to take replace of the self-attention mechanism in SAMVDPC, and set in objective function to 0. The h per parameters on each data sets are summarized in Table 2.

ith regard to the initialization of SAMVDPC eights, In Encoder-Block, e pre-train $V$ auto-encoders through minimizing the reconstruction error of each vie , and then use the pre-trained parameters of auto-encoders to initialize the corresponding encoder's eight of Encoder-Block. In self-attention MLP and MLP, Xavier is used as the eight initialization method [ ].

In training process, the optimizer algorithm e used is Adam, the learning rate is al a s initialized to $10^{-3}$, $10^{-4}$, $10^{-5}$, and ill decreases graduall ith the development of training process. To avoid over tting, all la ers in Encoder-Block and MLP are regularized b dropout reg-

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

### . Result

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

All the classi cation results of eight multi-vie datasets are summarized in Table 3, and the best result on each dataset is highlighted in boldface. As e can see, the proposed SAMVDPC achieves best accurac on Te as and Yale-Face datasets, and are comparable ith other algorithms on the else datasets. The promising result ma reason from four aspects: (1) DICS, baseline models, and SAMVDPC are all algorithms e ploiting the consistenc and complementarit , simultaneousl , and compared to NMF, mult-iNMF, and MVCC, the all achieve better performance on all datasets (2) compared to matri factorization algorithms, the Encoder-Block in both baseline models and SAMVDPC can e tract features in a a of effectivel fetching consistent information and grasping the underl - ing common properties of multi-vie datasets (3) compared to baseline models, the complementar ith diversit e ploited b self-attention mechanism contains more information than ma -pooling, mean-pooling, and eighted summation.

### . Convergence Analysis of Training rocess

In order to empiricall investigate the convergence propert of SAMVDPC, e plot the iterative curves of objective function and the corresponding classi cation accuracies on three t pical data sets, Leaves, BBC, and Te as in Fig. 2. From Fig. 2, e can observe that: (1) the objective function values drop sharpl and mean hile the classi cation accuracies increase rapidl ithin the previous rounds of iterative process, and then the objective function and the accurac curves begin to decrease/gro mildl , nall converge to a value or uctuate around a constant (2) ith respect to convergence speed, the objective function values of SAMVDPC converge in the least iterations, in contrast, ma -pooling corresponds to the most iterations, because ma -pooling operation is loss compression process and the backpropagation process doesn't make full use of information from multiple vie s data (3) in respect of convergence result, the objective function of SAMVDPC can nall converge to a ed value on ever dataset, but the objective function of baseline models al a s nall uctuate around a constant, hat's more, compared to baseline models, e can nd that the classi cation accurac curves of SAMVDPC often uctuate ithin a narro range. In conclusion, compared to baseline models, SAMVDPC get a better performance on the iterative curves of objective function and the corresponding classi cation accurac .
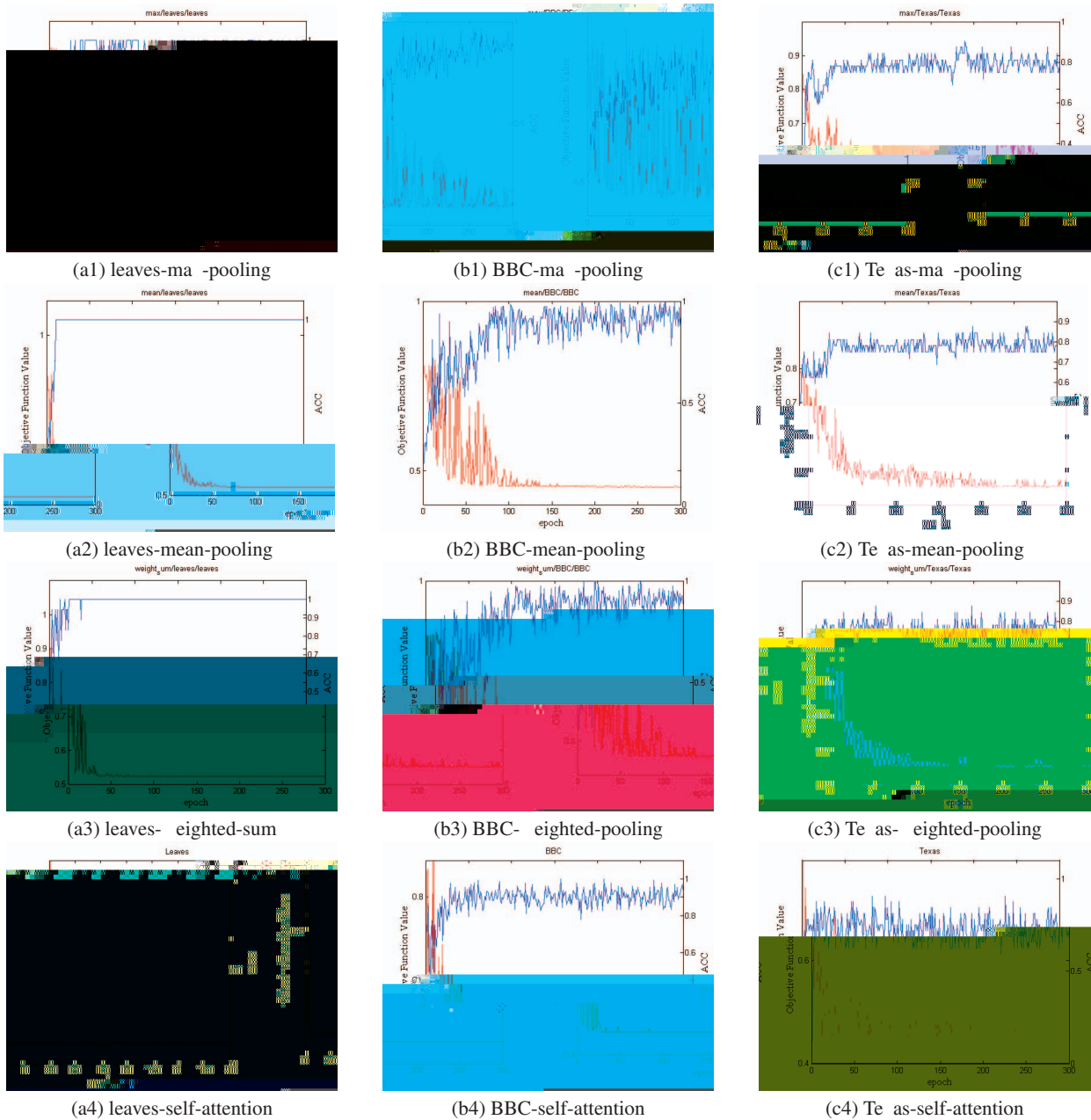
(a1) leaves-ma -pooling



(b1) BBC-ma -pooling



(c1) Te as-ma -pooling



(a2) leaves-mean-pooling



(b2) BBC-mean-pooling



(c2) Te as-mean-pooling



(a3) leaves- eighted-sum



(b3) BBC- eighted-pooling



(c3) Te as- eighted-pooling



(a4) leaves-self-attention



(b4) BBC-self-attention



(c4) Te as-self-attention

Figure 2: The convergence propert  of SAMVDP

## Conclusion and  uture Wor

In this paper,  e propose a novel multi-vie  net ork,

[2] X. Li, B. Geng, J. Zha, D. Tao, L. Xang, and C. Xu. 2011. Difficult guided image retrieval using linear multiview embedding. In 19th ACM Multimedia Proceedings, 1169-1172.

[3] C. Wan, R. Pan, and J. Li. 2011. Bi-weighting domain adaptation for cross-language text classification. In 22nd IJCAI Proceedings, 1535-1540.

[4] B. Xie, X. Mu, D. Tao, and . 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] . Chaudhuri, S. M. Kakade, . Livescu, and . 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. Shawe-Taylor, and S. Szedm'k. 2005. Two view learning: SVM-2 , Theory and Practice. In 18th NIPS Proceedings, 355-362.

[8] D. R. Hardoon, S. Szedm'k, and J. Shawe-Taylor. 2004. Canonical correlation analysis: an overview with application to learning methods. Neural Computation, 16(12):2639-2664.

[9] M. Wan, S. Shan, H. Zhang, S. Lao, and X Chen. 2016. Multiview discriminant analysis. IEEE Trans. Pattern Anal. Mach. Intell, 38(1):188-194.

[10] A. Sharma, A. Kumar, H. Daum', D. and . Jacobs. 2012. Generalized multiview analysis: a discriminative latent space. In 20th CVPR Proceedings, 2160-2167.

[11] . Guan, L. Zhang, J. Peng, and J Fan. 2015. Multi-view concept learning for data representation. IEEE Trans. Knowl. Data Eng., 27(11): 3016-3028.

[12] J. Gao, J. Han, J. Liu, and C. Wang. 2013. Multi-view clustering via joint nonnegative matrix factorization. In 13th SDM Proceedings, 252-260.

[13] H. Wang, X. 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'. 2011. A co-training approach for multi-view spectral clustering. In 28th ICML Proceedings, 393-400.

[16] . Wang, . Zhou. 2010. A new analysis of co-training. In 27th ICML Proceedings, 1135-1142.

[17] M. Zhang, . Zhou. 2011. CoTrade: confident co-training with data editing. IEEE Trans. Systems, Man, and Cybernetics, 41(6):1612-1626.

[18] . Zhang, . Jin, P. Li, . Xang, and J. Shao. 2018. Multiview discriminative learning via joint non-negative matrix factorization. In 23rd DASFAA Proceedings, 542-557.

[19] A. P. Singh, . J. Gordon. 2008. Relational learning via collective matrix factorization. In 14th KDD 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. Y. Ng: 2011. Multimodal Deep Learning. In 1rth KDD Proceedings, 689-696.

[22] . Cho, B. V. Merrienboer, . '?ehre, 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. Kalogerakis, and E. . Learned-Miller. 2015. Multi-view Convolutional Neural Networks for 3D Shape Recognition. In 2015 ICCV Proceedings, 945-953.

[24] D. Bahdanau, . 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 . V. Le. 2014. Sequence to Sequence Learning with Neural Networks. In 28th NIPS Proceedings, 3104-3112.

[26] . Xu, J. Ba, R. Kiros, . Cho, A. C. Courville, R. Salakhutdinov, R. S. Zemel, 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-Networks for Machine Reading. In 2016 EMNLP Proceedings, 551-561.

[29] A. P. Parikh, O. T?ckstr?m, D. Das, and J. Uszkoreit. 2016. A Decomposable Attention Model for Natural Language Inference. In 2016 EMNLP Proceedings, 2249-2255.

[30] . Lin, M. Feng, C. N. d. Santos, M. Xu, B. Xiang, B. Zhou, 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 31t h NIPS Proceedings, 6000-6010.

[32] . Wang, X. Kong, H. Fu, M. Li, and Y. Zhang. 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.