

# Transfer Learning via Multi-View Principal Component Analysis

Yang-Sheng Ji (吉阳生), Jia-Jun Chen (陈家骏), *Member, CCF*, Gang Niu (牛 罡)  
Lin Shang (商 琳), *Member, CCF*, and Xin-Yu Dai (戴新宇), *Member, CCF*

*Department of Computer Science and Technology, Nanjing University, Nanjing 210093, China*

*National Key Laboratory for Novel Software Technology, Nanjing University, Nanjing 210093, China*

E-mail: {jyis, chenjj}@nlp.nju.edu.cn; niugang@ai.nju.edu.cn; shanglin@nju.edu.cn; dxy@nlp.nju.edu.cn

Received December 31, 2009; revised November 12, 2010.

**Abstract** Transfer learning aims at leveraging the knowledge in labeled source domains to predict the unlabeled data in a target domain, where the distributions are different in domains. Among various methods for transfer learning, one kind of algorithms focus on the correspondence between bridge features and all the other specific features from different domains, and later conduct transfer learning via the *single-view* correspondence. However, the *single-view* correspondence may prevent these algorithms from further improvement due to the problem of incorrect correlation discovery. To tackle this problem, we propose a new method for transfer learning in a *multi-view* correspondence perspective, which is called *Multi-View Principal Component Analysis* (MVPCA) approach. MVPCA discovers the correspondence between bridge features representative across all domains and specific features from different domains respectively, and conducts the transfer learning by dimensionality reduction in a *multi-view* way, which can better depict the knowledge transfer. Experiments show that MVPCA can significantly reduce the cross domain prediction error of a baseline non-transfer method. With *multi-view* correspondence information incorporated to the *single-view* transfer learning method, MVPCA can further improve the performance of one state-of-the-art *single-view* method.

**Keywords** transfer learning, multi-view principal component analysis, text mining, sentiment classification

## 1 Introduction

In machine learning, traditional classification and regression algorithms strictly rely on the i.i.d. assumption, which requires that training and testing data are independent and identically distributed. However, real world applications usually show disobedience of this assumption. Let us take the sentiment polarity classification of product reviews on Amazon as an example. For one kind of product selling on Amazon, we can use all the labeled reviews to train a sentiment polarity classifier for the reviews of this product. However, new products are emerging on Amazon every day and lots of new product reviews will be posted afterwards. How could these newly posted product reviews be automatically classified as “positive” or “negative”? Of course, we can achieve this by spending lots of human efforts and time to annotate huge amount of new reviews to train a new classifier. But it is so expensive that we cannot do this all the time. Another easier way is to predict the sentiment polarity of new reviews by adapting an existing classifier trained on another domain. However,

different types of product reviews tend to use different domain specific terms. So the term distributional difference will cause a lot of troubles for traditional classification methods. Table 1 shows some domain specific terms from four domains on Amazon: Book, DVD, Electronic, Kitchen. This gives us an intuitive impression on the distributional difference between domains. To solve this non-i.i.d. classification problem, transfer learning research has attracted a lot of interest.

Transfer learning tackles the problem of predicting

**Table 1.** Domain Specific Terms

Book	DVD	Electronic	Kitchen
two-stars	script	warranty	ease
must-read	dumb	it-worked	broken
pages	pathetic	unreliable	kitchen
repetitive	atrocious	works-great	convenient
great-book	great-dvd	stopped-working	cooking
poorly-written	predictable	memory	easy-to
excellent-book	intense	an-error	clean
required-reading	pointless	battery	sharp

Regular Paper

This work is supported by the National Natural Science Foundation of China under Grant No. 61003112, the National Basic Research 973 Program of China under Grant No. 2010CB327903, the National High-Tech Research and Development 863 Program of China under Grant No. 2006AA010109, and the Natural Science Foundation of Jiangsu Province under Grant No. 2009233.

©2011 Springer Science + Business Media, LLC & Science Press, China

testing instances drawn from a different but related distribution compared with training instances. Generally speaking, transfer learning aims at minimizing the differences between distributions of different domains to minimize the cross domain prediction error. In this paper, we assume that  $P(\mathcal{X}_s) \neq P(\mathcal{X}_t)$  and  $P(\mathcal{Y}_s) = P(\mathcal{Y}_t)$ , which are in a common setting of transfer learning.  $P(\mathcal{X}_s)$  and  $P(\mathcal{Y}_s)$  denote the distributions of instances and labels in source domain respectively.  $P(\mathcal{X}_t)$  and  $P(\mathcal{Y}_t)$  are similarly defined in the target domain.

To achieve the goal of minimizing the prediction error across domains, transfer learning methods aim at making  $P(\mathcal{X}_s)$  and  $P(\mathcal{X}_t)$  less differ. Among proposed approaches, one kind of algorithms try to learn a low dimensional semantic space, which is able to depict the common knowledge across source and target domains. Represented by the common knowledge in the semantic space, source domain and target domain become more similar. In this way, Blitzer *et al.*<sup>[1-2]</sup> proposes a *single-view* transfer learning algorithm, Structural Correspondence Learning (SCL), in which a latent semantic space is generated via a *single-view* correspondence. SCL<sup>[1-2]</sup> seeks to learn the common knowledge across domains in the semantic space by applying the SVD procedure to a correspondence matrix. This correspondence matrix represents the statistical correlations between bridge features and all other features. Briefly, bridge features are those representative words occurring frequently across all domains, with detailed specification in Subsection 3.2. It is notable that this correspondence matrix mixes up all the other features across domains except bridge features, and measures the statistical correlations between them. In this paper, we regard this matrix as a *single-view* correspondence. *Single-view* correspondence can really help knowledge transfer. If the word “well-written” is a unique term in Book domain, and the word “use-easily” is another unique term in Kitchen domain, then they may never be correlated since they only occur separately. So if a classifier  $f(\text{“well-written”})$  is trained on Book domain, it has no prediction power on the word “use-easily” in Kitchen domain. However, when you consider another bridge word “great”, which may have strong correlations with both “well-written” and “use-easily”,  $f(\text{“well-written”})$  becomes predictive on the word “use-easily”. Because when you see “well-written”, it is assumed that “great” can be seen in some extent and then the probability of seeing the word “use-easily” can be estimated.

The above example shows how SCL<sup>[1-2]</sup> works. However, SCL has some limitations by utilizing the *single-view* correspondence. Let us zoom into the limitations by digesting a negative example. In DVD domain, if

one product reviewer says “read-the-book”, he means this DVD is “boring”. So “read-the-book” has high correlation with “boring” in DVD domain. While taking a look at Book domain, if one product reviewer says “read-the-book”, he means the book is “well-written” or “great”. Surprisingly, besides the term “boring”, “read-the-book” also has high correlations with “well-written” and “great”. Obviously, it is a contradiction which is caused by the vague meanings of “read-the-book” in different domains. If we mix the instances from Book domain and DVD domain as SCL<sup>[1-2]</sup> does, and learn the *single-view* correspondence matrix on the mixed dataset, we can get a strong prediction between “boring” and “great”. In this paper, we call this phenomenon incorrect correlation discovery. As these words with vague meanings do appear in sentiment texts on the Internet, we may often encounter incorrect correlations between words if we apply the *single-view* transfer learning methods. This problem would surely decrease the performances of these *single-view* algorithms.

In order to alleviate this problem, this paper proposes a new transfer learning method by utilizing *Multi-View Principal Component Analysis* (MVPCA). In MVPCA, firstly we choose  $m$  bridge features as SCL<sup>[1-2]</sup> does. Secondly, in every single domain, we learn a correspondence between bridge features and the other features. We call these correspondences in different domains *multi-view* correspondences. In this way, we separate the *single-view* correspondence of SCL algorithm into multiple parts. In every single view of correspondence of MVPCA, those words with vague meanings can only show its domain dependent meaning and they are not correlated with those contradictory words simultaneously in the single correspondence. Thus incorrect correlations between specific words from different domains are prevented. Finally, we apply MVPCA to learn a low dimensional semantic space. In this procedure, common knowledge of cross domain correlation is preserved in the eigenvectors as SCL does. Thus, knowledge transfer can be facilitated.

Compared with previous related work, the main contributions of MVPCA algorithm include the following.

1) This paper points out that words with vague meanings may do harm to *single-view* transfer learning methods. And we propose MVPCA to tackle this problem in a *multi-view* perspective. To the best of our knowledge, no previous work performs multi-view learning before prediction function  $f(\cdot)$  is learned and no previous transfer learning method aims at this problem.

2) MVPCA opens a new scope for transfer learning. In order to learn a low dimensional semantic space,

*multi-view* correspondence learning offers a different way from *single-view* correspondence learning, which may achieve neat and clear correspondences between bridge features and different domain specific features in each domain. With the neat correspondence information incorporated into the *single-view* correspondence matrix, experimental results show that the performance of *single-view* transfer learning method can be further improved.

The rest of this paper is organized as follows. Section 2 discusses related work. Section 3 describes MVPCA algorithm in detail. Section 4 illustrates experimental results to show the effectiveness of MVPCA. Section 5 gives the conclusion and discussion on future work.

## 2 Related Work

Recently, many transfer learning algorithms have been proposed in the line of extracting common latent semantic space. Following this intuitive idea, SCL<sup>[1-2]</sup>, Co-Clustering<sup>[3]</sup>, and Transferred Component Analysis<sup>[4]</sup> are applied to learning the low dimensional semantic space. In these low dimensional semantic space learning methods, SCL<sup>[1-2]</sup> is one of the state-of-the-art transfer learning methods.

Daumé<sup>[5]</sup> suggests a new feature representation for both source domain and target domain. In this new feature space, both domain specific features and shared common features can be recognized to help improve knowledge transfer. Although [5] is easy to perform, it may suffer from “the curse of dimensionality”. Daumé et al.<sup>[6]</sup> proposes a different way of transfer via adaptation on statistical maximal entropy classifier. Dai et al.<sup>[7]</sup> carries out the co-clustering<sup>[3]</sup> procedure on data combination of both source domain and target domain to find a shared latent space across domains. Similarly, Pan et al.<sup>[4]</sup> proposes a transferred component analysis to reduce the *Maximum Mean Discrepancy* (MMD) as a criterion to measure the difference between domains. All of the methods mentioned above achieve remarkable performances on knowledge transfer via newly generated latent space.

It is notable that Pan et al.<sup>[8]</sup> also propose a *multi-view* transfer learning algorithm, which predicts physical locations via an estimation function mapping from signal space to physical location space. Resulted from time variant signal space, *multi-view* refers to one view of original estimation function and the other view of current signal space data, which is utilized to adjust the previous estimation function to predict current location. MVPCA is different from [8]. On one hand, MVPCA performs *multi-view* learning before prediction function  $f(\cdot)$  is estimated, while [8] performs *multi-view* adjustment after prediction function on

source domain has been estimated. On the other hand, before the prediction function is estimated, MVPCA can employ as much information and as neat correspondence as it could to improve the transferability of the latent semantic space. The neat correspondence information may enable MVPCA to better depict intrinsic shared structures across domains than [8] does.

Accompanying with the viewpoint of extracting latent semantic space across domains, transfer learning algorithms aiming at instance weighting are proposed as well. Dai et al.<sup>[9]</sup> assign different weights to instances in different domains, and transfer the knowledge in a similar way of boosting. Huang et al.<sup>[10]</sup>, Zadrozny<sup>[11]</sup> consider domain difference as a sample selection problem which differs the importance of instances in transferring. Jiang et al.<sup>[12]</sup> solve an application in natural language processing in the viewpoint of instance weighting. Besides above algorithms, Raina et al.<sup>[13]</sup> apply the sparse coding algorithm<sup>[14]</sup> to the problem of self-taught learning and achieves satisfactory performance. Tan et al.<sup>[15]</sup> and Dai et al.<sup>[16]</sup> propose modified versions of Naive Bayes to tackle the transfer learning in different viewpoints. Sandler et al.<sup>[17]</sup> incorporate manifold information to depict the relationship between features, and facilitate knowledge transfer via feature network constraint.

## 3 Transfer Learning via Multi-View Perspective

### 3.1 Problem Setting and Notations

Before introducing Multi-View Principal Component Analysis algorithm in detail, we should specify the problem setting and notations. In our transfer learning setting, *domain* and *task* are two basic concepts following the common transfer learning analysis.

A domain  $\mathcal{D}$  consists of a feature space  $\mathcal{X} \subset \mathbb{R}^d$  and label  $\mathcal{Y} \subset \mathbb{R}$ , distributions  $P(\mathcal{X})$  and  $P(\mathcal{Y})$ , i.e.,  $\mathcal{D} = \{\mathcal{X}, \mathcal{Y}, P(\mathcal{X}), P(\mathcal{Y})\}$ . Also  $\mathcal{D}$  can be rewritten as  $\mathcal{D} = \{(\mathbf{x}_i, y_i) | \mathbf{x}_i \in \mathcal{X}, y_i \in \mathcal{Y}\}$ , which implies  $P(\mathcal{X})$  and  $P(\mathcal{Y})$  implicitly;  $\mathbf{x}_i$  denotes an instance in domain  $\mathcal{D}$  which is a term vector of a document with Bag of Words or TFIDF representation;  $y_i$  denotes the label of instance  $\mathbf{x}_i$ ;  $[\mathbf{x}_i]_j$  denotes the  $j$ -th coordinate value of the instance and  $[\mathbf{x}_i]_j \in \mathcal{X}_j$ . Source domain is denoted by  $\mathcal{D}_s = \mathcal{D}_{sl} \cup \mathcal{D}_{su} \subseteq \mathcal{X}_s \times \mathcal{Y}_s$ ,  $\mathcal{Y}_s = \{+1, -1, 0\}$ .  $\mathcal{D}_{sl} = \{(\mathbf{x}_i^s, y_i^s) | i = 1, \dots, n_{sl}\}$  denotes labeled data;  $\mathcal{D}_{su} = \{(\mathbf{x}_i^{su}, y_i^{su} = 0) | i = 1, \dots, n_{su}\}$  denotes unlabeled data;  $\mathbf{X}_s = \{\mathbf{x}_i^s | i = 1, \dots, n_{sl} + n_{su}\}$ ;  $\mathbf{Y}_s = \{y_i^s | i = 1, \dots, n_{sl} + n_{su}\}$ ;  $\mathbf{x}_i^s$  is a positive\negative instance if  $y_i^s = +1 \setminus -1$ ,  $\mathbf{x}_i^s$  is an unlabeled instance if  $y_i^s = 0$ .  $\mathcal{D}_t$ ,  $\mathcal{D}_{tl}$ ,  $\mathcal{D}_{tu}$ ,  $\mathbf{X}_t$  and  $\mathbf{Y}_t$  are similarly defined in target domain.

A task  $\mathcal{T}$  consists of a specified domain  $\mathcal{D}$  and a prediction function  $f(\cdot)$ , i.e.,  $\mathcal{T} = \{\mathcal{D}, f(\cdot)\}$ .  $f(\cdot)$  is a prediction function learned on domain  $\mathcal{D}$ , which maximizes likelihood or minimizes a loss function, i.e.,  $f(\cdot) = \arg \min_f \sum_{\mathbf{x}_i \in \mathcal{X}} \|f(\mathbf{x}_i) - h(\mathbf{x}_i)\|_2^2$ , where  $h(\cdot)$  is a ground truth prediction function on domain  $\mathcal{D}$ .

Now we can give a formal definition of transfer learning in our setting.

**Definition 1** (Transfer Learning). *Given source domain  $\mathcal{D}_s = \mathcal{D}_{st} \cup \mathcal{D}_{su}$  and target domain  $\mathcal{D}_t = \mathcal{D}_{tu}$ , transfer learning aims at learning a task  $\mathcal{T} = \{\mathcal{D}_s \cup \mathcal{D}_t, f(\cdot)\}$  in order that  $\text{err}_{\mathcal{D}_t}(f) = \sum_{\mathbf{x}_i \in \mathcal{X}_t} \|f(\mathbf{x}_i) - h_t(\mathbf{x}_i)\|_2^2$  can be minimized. Note that,  $h_t(\cdot)$  is a ground truth prediction function on domain  $\mathcal{D}_t$ , and we often approximate  $h_t(\mathbf{x}_i^t)$  by  $y_i^t$ . This definition holds under following assumptions:*

- 1)  $\mathcal{X}_s \cap \mathcal{X}_t \neq \emptyset$ ;
- 2)  $P(\mathcal{X}_s) \neq P(\mathcal{X}_t)$ ;
- 3)  $P(\mathcal{Y}_s) = P(\mathcal{Y}_t)$ ;
- 4) *this paper only focuses on one source domain and one target domain.*

### 3.2 Bridge Features and Multi-View Correspondence Learning

Let us come to the definition of bridge feature, which is one of the key parts of MVPCA algorithm. We denote one single bridge feature as  $\mathcal{X}_{\Lambda_i}$ , and denote the bridge feature set as  $\mathcal{X}_{\Lambda}$ . One of the most important reasons that MVPCA shows good performance on domain adaptation is that our chosen bridge features contribute great discriminative ability across domains when performing multi-view transfer learning. The discriminative power of bridge features is definitely derived from high correlations with class labels. So bridge features can be regarded as pseudo labels, and these discriminative pseudo labels can bring a lot of helpful discriminative power across domains to benefit the knowledge transfer.

For short, to achieve the discriminative power on both domains, bridge features must satisfy the following assumptions:

- $\mathcal{X}_{\Lambda_i}$  is one of the representative features in both domains; all bridge features constitute bridge feature set  $\mathcal{X}_{\Lambda} = \{\mathcal{X}_{\Lambda_1}, \dots, \mathcal{X}_{\Lambda_m}\} \subseteq \mathcal{X}_s \cap \mathcal{X}_t$ ;
- each bridge feature  $\mathcal{X}_{\Lambda_i}$  is statistical correlated with class labels. In this paper, mutual information is utilized as the measurement.

To construct the bridge feature set  $\mathcal{X}_{\Lambda}$  in MVPCA, we may simply pick up  $m$  features  $\mathcal{X}_i \subseteq \mathcal{X}_s \cap \mathcal{X}_t$  with the highest mutual information values and group them into  $\mathcal{X}_{\Lambda}$ . Mutual information is computed between

feature  $\mathcal{X}_i$  and  $\mathcal{Y}_s$  on  $\mathcal{D}_{st}$  according to (1):

$$I(\mathcal{X}_i; \mathcal{Y}_s) = \sum_{y \in \mathcal{Y}_s} \sum_{x \in \mathcal{X}_i} p(x, y) \log \frac{p(x, y)}{p_1(x) p_2(y)}. \quad (1)$$

After  $\mathcal{X}_{\Lambda}$  is ready, *multi-view* correspondence can be extracted thereafter. In this paper, *multi-view* correspondence consists of two views on source and target domains respectively. One view is the correspondence between bridge features  $\mathcal{X}_{\Lambda}$  and other features  $\mathcal{X}_s - \mathcal{X}_{\Lambda}$  in  $\mathcal{D}_s$ ; the other view is the correspondence between bridge features and other features  $\mathcal{X}_t - \mathcal{X}_{\Lambda}$  in  $\mathcal{D}_t$ . Anyway, *two-view* MVPCA algorithm is a straightforward and representative version of *multi-view* MVPCA, and the *multi-view* version can be easily extended, which is specified in Theorem 2.

Then let us take source domain  $\mathcal{D}_s$  as an example to illustrate how to measure the correspondence between bridge feature  $\mathcal{X}_{\Lambda_i}$  and other text features in  $\mathcal{X}_s$  in detail. For each instance  $(\mathbf{x}_i^{su}, y_i^{su}) \in \mathcal{D}_{su}$ , check the value of bridge feature  $\mathcal{X}_{\Lambda_i}$  in instance  $(\mathbf{x}_i^{su}, y_i^{su})$ . If  $[\mathbf{x}_i^{su}]_{\Lambda_i} = 0$ , set  $y_i^{su} = -1$ ; otherwise, set  $[\mathbf{x}_i^{su}]_{\Lambda_i} = 0$ ,  $y_i^{su} = +1$ . In this way, we can construct a pseudo labeled dataset  $\mathcal{D}_s^{\Lambda_i}$ . On this pseudo labeled domain, a linear classifier  $f_{\mathcal{D}_s^{\Lambda_i}}(\cdot)$  can be learned and  $f_{\mathcal{D}_s^{\Lambda_i}}(\cdot)$  provides a classification hyper plane  $\boldsymbol{\theta}_{\mathcal{D}_s^{\Lambda_i}} \in \mathbb{R}^{d_s}$ . This hyper plane  $\boldsymbol{\theta}_{\mathcal{D}_s^{\Lambda_i}}$  can be regarded as the correspondence between bridge feature  $\mathcal{X}_{\Lambda_i}$  and other features on source domain. So we may arrange all  $\boldsymbol{\theta}_{\mathcal{D}_s^{\Lambda_i}}$  together as  $\boldsymbol{\Theta}_{\mathcal{D}_s^{\Lambda}} = [\boldsymbol{\theta}_{\mathcal{D}_s^{\Lambda_1}}, \dots, \boldsymbol{\theta}_{\mathcal{D}_s^{\Lambda_m}}] \in \mathbb{R}^{d_s \times m}$ . Similarly,  $\boldsymbol{\Theta}_{\mathcal{D}_t^{\Lambda}} = [\boldsymbol{\theta}_{\mathcal{D}_t^{\Lambda_1}}, \dots, \boldsymbol{\theta}_{\mathcal{D}_t^{\Lambda_m}}] \in \mathbb{R}^{d_t \times m}$  can be constructed in the same way. Up to now, every bridge feature  $\mathcal{X}_{\Lambda_i}$  has set up *multi-view* correspondences with features in source domain and target domain respectively:  $\boldsymbol{\theta}_{\mathcal{D}_s^{\Lambda_i}}$  and  $\boldsymbol{\theta}_{\mathcal{D}_t^{\Lambda_i}}$ . By putting them together,  $(\boldsymbol{\theta}_{\mathcal{D}_s^{\Lambda_i}}^T, \boldsymbol{\theta}_{\mathcal{D}_t^{\Lambda_i}}^T)$  represents a *multi-view* perspective on the correspondences of bridge feature  $\mathcal{X}_{\Lambda_i}$ . So we rewrite  $\boldsymbol{\Theta}_{\mathcal{D}_s^{\Lambda}}$  and  $\boldsymbol{\Theta}_{\mathcal{D}_t^{\Lambda}}$  as  $\boldsymbol{\Theta}_{\Lambda}^{(s)}$  and  $\boldsymbol{\Theta}_{\Lambda}^{(t)}$  respectively to show that they are two views of feature correspondences.

By setting up *multi-view* correspondence on both domains, strong correlations between bridge features and specific features can be detected within every domain. In this way, those words with vague meanings can only give one sentimental polarity in one view of the correspondence. And those contradictory correlations are prevented in correspondence. As a result, *multi-view* correspondences can give a cleaner and neater within domain knowledge, and preserve cross domain knowledge via bridge features as SCL does. To conclude this subsection, in *multi-view* perspective, we separate the correspondence between features from *single-view* to *multi-view* in order to better depict within domain

knowledge while preserving the same cross domain knowledge as SCL does.

### 3.3 Multi-View Principal Component Analysis

After setting up the *multi-view* correspondence, we may apply Multi-View Principal Component Analysis to these correspondence matrices to get a latent semantic space, which consists of a group of eigenvectors. These eigenvectors of MVPCA not only prevent incorrect correlation generation, but also preserve the cross domain knowledge. Thus MVPCA can further improve the performance of domain adaptation of SCL.

#### 3.3.1 Objective Function Specification and Notations

Before introducing the details of MVPCA algorithm, we specify the intuitive idea and notations first. Suppose we have *two views* of instances  $\mathbf{X}_s = \{\mathbf{x}_i^s \in \mathbb{R}^{d_s} | i = 1, \dots, m\}$  and  $\mathbf{X}_t = \{\mathbf{x}_i^t \in \mathbb{R}^{d_t} | i = 1, \dots, m\}$ . We call  $\mathbf{X}_s$  and  $\mathbf{X}_t$  source view and target view respectively. And the instances with index  $i$  in source view and target view, e.g.,  $\mathbf{x}_i^s$  and  $\mathbf{x}_i^t$ , are describing the same instance. So all the  $m$  instances in both views correspond one to one.

In order to generate a low dimensional semantic space, one common way is to compress instances from *two views* into *single view*  $\mathbf{X} = \{\mathbf{x}_i \in \mathbb{R}^{d_s+d_t} | \mathbf{x}_i = \langle \mathbf{x}_i^s, \mathbf{x}_i^t \rangle; i = 1, \dots, m\}$ .  $\mathbf{x}_i = \langle \mathbf{x}_i^s, \mathbf{x}_i^t \rangle$  denotes the spliced vector of  $\mathbf{x}_i^s$  and  $\mathbf{x}_i^t$ . Then by applying SVD or PCA on  $\mathbf{X}$ , we can get a low dimensional semantic space. However, there are some drawbacks of this method. Firstly, when applied to transfer learning, we can still see the incorrect correlation between some contradictory domain specific features. So in the low dimensional semantic space, helpful common knowledge and incorrect correlations are fused together, which can do harm to the transfer learning. As a result, we cannot expect it to show better performance than the *single-view* transfer learning method. Secondly, Vinokourov et al.<sup>[18]</sup> and Li et al.<sup>[19]</sup> point out that when performing cross-language classification, this naive method does not have promising results as Hardoon's algorithm<sup>[20]</sup> does.

In another way, if we can learn a set of uniform orthogonal eigenvectors  $\mathbf{U} = \{\mathbf{u}_j | j = 1, \dots, p\}$  across different feature spaces of  $\mathcal{X}_s$  and  $\mathcal{X}_t$ , the set of eigenvectors  $\mathbf{U}$  may not only preserve helpful common knowledge across source and target domains, but also prevent incorrect correlation generation. While learning a set of unified eigenvectors across different feature spaces seems infeasible, we can achieve this by merging  $\mathcal{X}_s$  and  $\mathcal{X}_t$  into a unified feature space  $\hat{\mathcal{X}} = \hat{\mathcal{X}}_s = \hat{\mathcal{X}}_t \subset \mathbb{R}^d$ , which is just the same as a merged dictionary from another two. In the unified feature space,  $\mathbf{X}_s = \{\mathbf{x}_i^s \in$

$\mathbb{R}^{d_s} | i = 1, \dots, m\}$  and  $\mathbf{X}_t = \{\mathbf{x}_i^t \in \mathbb{R}^{d_t} | i = 1, \dots, m\}$  are rewritten as  $\hat{\mathbf{X}}_s = \{\hat{\mathbf{x}}_i^s \in \mathbb{R}^d | i = 1, \dots, m\}$  and  $\hat{\mathbf{X}}_t = \{\hat{\mathbf{x}}_i^t \in \mathbb{R}^d | i = 1, \dots, m\}$ . We may also denote  $\hat{\mathbf{X}}_s$  by  $[\hat{\mathbf{x}}_1^s, \dots, \hat{\mathbf{x}}_m^s]$  and denote  $\hat{\mathbf{X}}_t$  by  $[\hat{\mathbf{x}}_1^t, \dots, \hat{\mathbf{x}}_m^t]$  for convenience. Under these settings, we give the details of Multi-View Principal Component Analysis:

**Definition 2** (Objective Function of Multi-View Principal Component Analysis). *Suppose we have two views of instances  $\mathbf{X}_s = \{\mathbf{x}_i^s \in \mathbb{R}^{d_s} | i = 1, \dots, m\}$  and  $\mathbf{X}_t = \{\mathbf{x}_i^t \in \mathbb{R}^{d_t} | i = 1, \dots, m\}$ . Given a parameter  $p$ , Multi-View Principal Component Analysis tries to find a set of  $p$  uniform orthogonal vectors  $\mathbf{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_p\}$  which minimizes:*

$$\min_{\mathbf{U}} \mathcal{E}(\hat{\mathbf{X}}_s, \hat{\mathbf{X}}_t, \mathbf{U}, p) = \sum_{i=1}^m (\mathcal{E}_s^2(\hat{\mathbf{x}}_i^s, \mathbf{U}, p) + \mathcal{E}_t^2(\hat{\mathbf{x}}_i^t, \mathbf{U}, p)) \quad (2)$$

$$\begin{aligned} \text{s.t. } \mathcal{E}_s(\hat{\mathbf{x}}_i^s, \mathbf{U}, p) &= \hat{\mathbf{x}}_i^s - \sum_{j=1}^p ((\hat{\mathbf{x}}_i^s)^\top \cdot \mathbf{u}_j) \mathbf{u}_j \\ &= \hat{\mathbf{x}}_i^s - \sum_{j=1}^p \mathbf{u}_j (\hat{\mathbf{x}}_i^s)^\top \mathbf{u}_j \end{aligned} \quad (3)$$

$$\begin{aligned} \mathcal{E}_t(\hat{\mathbf{x}}_i^t, \mathbf{U}, p) &= \hat{\mathbf{x}}_i^t - \sum_{j=1}^p ((\hat{\mathbf{x}}_i^t)^\top \cdot \mathbf{u}_j) \mathbf{u}_j \\ &= \hat{\mathbf{x}}_i^t - \sum_{j=1}^p \mathbf{u}_j (\hat{\mathbf{x}}_i^t)^\top \mathbf{u}_j \end{aligned} \quad (4)$$

$$\mathbf{u}_i^\top \mathbf{u}_j = \begin{cases} 1, & \text{if } i = j, \\ 0, & \text{if } i \neq j. \end{cases} \quad (5)$$

The set of  $p$  uniform orthogonal vectors  $\mathbf{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_p\}$  build up a low dimensional semantic space. And in this semantic space, every instance can be reconstructed by the new coordinate system. The difference between original instance and reconstructed one is called residual error, denoted by  $\mathcal{E}_s^2(\hat{\mathbf{x}}_i^s, \mathbf{U}, p)$  and  $\mathcal{E}_t^2(\hat{\mathbf{x}}_i^t, \mathbf{U}, p)$ . So MVPCA aims at finding the semantic space which minimizes the total residual error of each instance under all views.

**Theorem 1.** *The first  $p$  eigenvectors with the largest eigenvalues of  $[\hat{\mathbf{X}}_s, \hat{\mathbf{X}}_t][\hat{\mathbf{X}}_s, \hat{\mathbf{X}}_t]^\top$  give the solution of the uniform orthogonal vector set to Definition 2.*

*Proof.* Firstly, let

$$\mathcal{E}(\hat{\mathbf{X}}_s, \mathbf{U}, p) = \sum_{i=1}^m \mathcal{E}_s^2(\hat{\mathbf{x}}_i^s, \mathbf{U}, p) \quad (6)$$

$$= \sum_{i=1}^m (\hat{\mathbf{x}}_i^s - \sum_{j=1}^p \mathbf{u}_j (\hat{\mathbf{x}}_i^s)^\top \mathbf{u}_j)^2 \quad (7)$$

$$= \sum_{i=1}^m ((\hat{\mathbf{x}}_i^s)^\top (\hat{\mathbf{x}}_i^s) - \sum_{j=1}^p \mathbf{u}_j^\top (\hat{\mathbf{x}}_i^s) (\hat{\mathbf{x}}_i^s)^\top \mathbf{u}_j) \quad (8)$$

$$= \text{trace}(\hat{\mathbf{X}}_s^T \hat{\mathbf{X}}_s) - \sum_{j=1}^p \mathbf{u}_j^T \hat{\mathbf{X}}_s \hat{\mathbf{X}}_s^T \mathbf{u}_j. \quad (9)$$

Similarly,

$$\mathcal{E}(\hat{\mathbf{X}}_t, \mathbf{U}, p) = \text{trace}(\hat{\mathbf{X}}_t^T \hat{\mathbf{X}}_t) - \sum_{j=1}^p \mathbf{u}_j^T \hat{\mathbf{X}}_t \hat{\mathbf{X}}_t^T \mathbf{u}_j. \quad (10)$$

Secondly,

$$\begin{aligned} & \min_{\mathbf{U}} \mathcal{E}(\hat{\mathbf{X}}_s, \hat{\mathbf{X}}_t, \mathbf{U}, p) \\ &= \min_{\mathbf{U}} (\mathcal{E}(\hat{\mathbf{X}}_s, \mathbf{U}, p) + \mathcal{E}(\hat{\mathbf{X}}_t, \mathbf{U}, p)) \end{aligned} \quad (11)$$

$$\begin{aligned} &= \min_{\mathbf{U}} \{ \text{trace}(\hat{\mathbf{X}}_s^T \hat{\mathbf{X}}_s + \hat{\mathbf{X}}_t^T \hat{\mathbf{X}}_t) - \\ & \sum_{j=1}^p \mathbf{u}_j^T (\hat{\mathbf{X}}_s \hat{\mathbf{X}}_s^T + \hat{\mathbf{X}}_t \hat{\mathbf{X}}_t^T) \mathbf{u}_j \} \end{aligned} \quad (12)$$

$$= \max_{\mathbf{U}} \sum_{j=1}^p \mathbf{u}_j^T (\hat{\mathbf{X}}_s \hat{\mathbf{X}}_s^T + \hat{\mathbf{X}}_t \hat{\mathbf{X}}_t^T) \mathbf{u}_j \quad (13)$$

$$= \max_{\mathbf{U}} \sum_{j=1}^p \mathbf{u}_j^T [\hat{\mathbf{X}}_s, \hat{\mathbf{X}}_t] [\hat{\mathbf{X}}_s, \hat{\mathbf{X}}_t]^T \mathbf{u}_j. \quad (14)$$

From the last step, we can see that  $\max_{\mathbf{U}} \sum_{j=1}^p \mathbf{u}_j^T [\hat{\mathbf{X}}_s, \hat{\mathbf{X}}_t] [\hat{\mathbf{X}}_s, \hat{\mathbf{X}}_t]^T \mathbf{u}_j$  is a standard form of Principal Component Analysis, and the solution to  $\mathbf{U}$  consists of the first  $p$  eigenvectors with largest eigenvalues of matrix  $[\hat{\mathbf{X}}_s, \hat{\mathbf{X}}_t] [\hat{\mathbf{X}}_s, \hat{\mathbf{X}}_t]^T$ .  $\square$

If we have multiple views  $\mathbf{X}_k = \{\mathbf{x}_i^k \in \mathbb{R}^{d_k} | i = 1, \dots, m\}$ ,  $k = 1, \dots, K$ , Theorem 1 can be further extended to a *multi-view* setting. After merging their feature spaces to a unified space, these multiple views are rewritten as  $\hat{\mathbf{X}}_k = \{\hat{\mathbf{x}}_i^k \in \mathbb{R}^d | i = 1, \dots, m\}$ , with  $k = 1, \dots, K$ .

**Theorem 2.** *If we have  $K$  views of instances,  $\hat{\mathbf{X}}_1 = \{\hat{\mathbf{x}}_i^1 \in \mathbb{R}^d | i = 1, \dots, m\}$  is the first view, with  $\hat{\mathbf{X}}_k$  ( $k = 2, \dots, K$ ) in the same notations. Given a parameter  $p$ , Multi-View Principal Component Analysis tries to find a set of  $p$  uniform orthogonal vectors  $\mathbf{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ , which minimizes the residual error across all the views of instances:*

$$\min_{\mathbf{U}} \mathcal{E}(\hat{\mathbf{X}}_1, \dots, \hat{\mathbf{X}}_K, \mathbf{U}, p) = \sum_{i=1}^m \sum_{k=1}^K (\mathcal{E}_k^2(\hat{\mathbf{x}}_i^k, \mathbf{U}, p)) \quad (15)$$

$$\text{s.t. } \mathcal{E}_k(\hat{\mathbf{x}}_i^k, \mathbf{U}, p) = \hat{\mathbf{x}}_i^k - \sum_{j=1}^p \mathbf{u}_j (\hat{\mathbf{x}}_i^k)^T \mathbf{u}_j, \quad k = 1, \dots, K \quad (16)$$

$$\mathbf{u}_i^T \mathbf{u}_j = \begin{cases} 1, & \text{if } i = j, \\ 0, & \text{if } i \neq j. \end{cases} \quad (17)$$

The first  $p$  eigenvectors with the largest eigenvalues of  $[\hat{\mathbf{X}}_1, \dots, \hat{\mathbf{X}}_K] [\hat{\mathbf{X}}_1, \dots, \hat{\mathbf{X}}_K]^T$  give the solution.

*Proof.* This proof is rather similar to the proof of Theorem 1. We simply omit it here.  $\square$

### 3.3.2 Multi-View Principal Component Analysis Algorithm

When we apply *Multi-View Principal Component Analysis* to transfer learning, we need to replace  $\mathbf{X}_s$  and  $\mathbf{X}_t$  in Theorem 1 by  $\boldsymbol{\theta}_{\Lambda}^{(s)}$  and  $\boldsymbol{\theta}_{\Lambda}^{(t)}$  respectively and then we get the algorithm. Now let us move to the details of *Multi-View Principal Component Analysis* algorithm, which are shown in Table 2.

**Table 2.** Multi-View Principal Component Analysis Algorithm

**Input:**  $\mathcal{D}_s = \mathcal{D}_{sl} \cup \mathcal{D}_{su}$ ,  $\mathcal{D}_t = \mathcal{D}_{tu}$ .

**Output:**  $\hat{f}_{\mathcal{D}_{sl}}$ .

**Begin**

(1) Choose  $m$  features from  $\mathcal{X}_s \cap \mathcal{X}_t$  to form bridge feature set  $\mathcal{X}_{\Lambda}$ .

(2) **For**  $\mathcal{X}_{\Lambda_i}$  in  $\mathcal{X}_{\Lambda}$

• construct pseudo labeled datasets:  $\mathcal{D}_s^{\Lambda_i}$  and  $\mathcal{D}_t^{\Lambda_i}$

• train linear classifiers on them:  $f_{\mathcal{D}_s^{\Lambda_i}}$  and  $f_{\mathcal{D}_t^{\Lambda_i}}$

• add  $\boldsymbol{\theta}_{\Lambda_i}^{(s)}$  into  $\boldsymbol{\theta}_{\Lambda}^{(s)}$ ,  $\boldsymbol{\theta}_{\Lambda_i}^{(t)}$  into  $\boldsymbol{\theta}_{\Lambda}^{(t)}$

**End**

(3) Run MVPCA:

$\mathcal{P}_{\Lambda} = \text{MVPCA}([\boldsymbol{\theta}_{\Lambda}^{(s)}, \boldsymbol{\theta}_{\Lambda}^{(t)}])$ .

(4)  $\hat{\mathcal{D}}_{sl} = \text{FeatureExpansion}(\mathcal{D}_{sl}, \mathcal{P}_{\Lambda})$

$\hat{\mathcal{D}}_t = \text{FeatureExpansion}(\mathcal{D}_t, \mathcal{P}_{\Lambda})$

(5) Train a linear classifier  $\hat{f}_{\mathcal{D}_{sl}}$  on  $\hat{\mathcal{D}}_{sl}$ .

$\hat{f}_{\mathcal{D}_{sl}}$  is the output.

Bridge features are those representative words occurring frequently across domains, which can be easily detected according to the specification in Subsection 3.2.

Then let us take source domain  $\mathcal{D}_s$  as an example to illustrate how to measure the correspondence matrices  $\boldsymbol{\theta}_{\Lambda}^{(s)}$  and  $\boldsymbol{\theta}_{\Lambda}^{(t)}$ . For every bridge feature  $\mathcal{X}_{\Lambda_i}$ , we can construct a pseudo labeled dataset  $\mathcal{D}_s^{\Lambda_i}$ . Next step is to train a linear classifier  $f_{\mathcal{D}_s^{\Lambda_i}}$  on this pseudo labeled domain. In this paper, logistic regression is utilized as our linear classifier:

$$L(X, \lambda) = \max_{\boldsymbol{\theta}} \sum_{x \in X} \frac{1}{1 + e^{-y\boldsymbol{\theta}^T x}} - \lambda \|\boldsymbol{\theta}\|_2^2. \quad (18)$$

We denote the hyper plane  $\boldsymbol{\theta}$  of logistic regression as  $\boldsymbol{\theta}_{\Lambda_i}^{(s)}$  on the pseudo labeled dataset  $\mathcal{D}_s^{\Lambda_i}$ . Then we may arrange all  $\boldsymbol{\theta}_{\Lambda_i}^{(s)}$  together as  $\boldsymbol{\theta}_{\Lambda}^{(s)}$ , which is one view of feature correspondence. Similarly,  $\boldsymbol{\theta}_{\Lambda}^{(t)}$  is another view of feature correspondence. The following step is to run *Multi-View Principal Component Analysis* to get a projection matrix.

In feature generation phase, for instance  $(\mathbf{x}_i^s, y_i^s) \in \mathcal{D}_s$ ,  $\mathbf{x}_i^s$  is expanded to  $\hat{\mathbf{x}}_i^s = \langle \mathbf{x}_i^s, \mathcal{P}_\Lambda^T \mathbf{x}_i^s \rangle$ . So in this feature expansion way, we get a feature expanded labeled data  $\hat{\mathcal{D}}_{s1}$  for source domain. Similarly, we can get  $\hat{\mathcal{D}}_t$ . Finally, we train a classifier  $\hat{f}_{\mathcal{D}_{s1}}$  on source domain  $\hat{\mathcal{D}}_{s1}$ , then we apply classifier  $\hat{f}_{\mathcal{D}_{s1}}$  to predict new instances from target domain  $\hat{\mathcal{D}}_t$ .

What is more interesting, by incorporating *single-view* correspondences of SCL algorithm into MVPCA, MVPCA’s performance can be further improved via learning a semantic space from both *multi-view* and *single-view* perspectives. This variant version of MVPCA is presented in Table 3, and we denote it by Variant Multi-View Principal Component Analysis (VMVPCA), which can also be regarded as a combined algorithm of MVPCA and SCL.

On conclusion of this section, we briefly review Multi-View Principal Component Analysis method. MVPCA separates *single-view* correspondence between features into *multiple views* of feature correspondence. The separate *multiple views* prevent the generation of incorrect and contradictory correlation between domain specific features. After that, we apply Multi-View Principal Component Analysis to generate a semantic space which preserves common knowledge across domains. In this way, MVPCA algorithm overcomes the drawback of *single-view* transfer learning methods, and preserves the benefits for knowledge transfer.

**Table 3.** Variant Multi-View Principal Component Analysis Algorithm

---

**Input:**  $\mathcal{D}_s = \mathcal{D}_{s1} \cup \mathcal{D}_{su}$ ,  $\mathcal{D}_t = \mathcal{D}_{tu}$ .

**Output:**  $\hat{f}_{\mathcal{D}_{s1}}$ .

**Begin**

- (1) Choose  $m$  features from  $\mathcal{X}_s \cap \mathcal{X}_t$  to form bridge feature set  $\mathcal{X}_\Lambda$ .
- (2) **For**  $\mathcal{X}_{\Lambda_i}$  in  $\mathcal{X}_\Lambda$ 
  - construct pseudo labeled datasets:  $\mathcal{D}_s^{\Lambda_i}$  and  $\mathcal{D}_t^{\Lambda_i}$
  - train linear classifiers on them:  $f_{\mathcal{D}_s^{\Lambda_i}}$  and  $f_{\mathcal{D}_t^{\Lambda_i}}$
  - add  $\theta_{\Lambda_i}^{(s)}$  into  $\theta_\Lambda^{(s)}$ ,  $\theta_{\Lambda_i}^{(t)}$  into  $\theta_\Lambda^{(t)}$
- End**
- (3) Run SCL procedure to learn a single-view correspondence:  $\theta_\Lambda$ .
- (4) Run MVPCA procedure:  
 $\mathcal{P}_\Lambda = \text{MVPCA}([\theta_\Lambda^{(s)}, \theta_\Lambda^{(t)}, \theta_\Lambda])$ .
- (5)  $\hat{\mathcal{D}}_{s1} = \text{FeatureExpansion}(\mathcal{D}_{s1}, \mathcal{P}_\Lambda)$   
 $\hat{\mathcal{D}}_t = \text{FeatureExpansion}(\mathcal{D}_t, \mathcal{P}_\Lambda)$
- (6) Train a linear classifier  $\hat{f}_{\mathcal{D}_{s1}}$  on  $\hat{\mathcal{D}}_{s1}$ .  
 $\hat{f}_{\mathcal{D}_{s1}}$  is the output.

---

<sup>①</sup><http://www.cs.jhu.edu/~mdredze/datasets/sentiment/>.

<sup>②</sup><http://people.csail.mit.edu/jrennie/20Newsgroups/>.

## 4 Experimental Result Analysis

### 4.1 Data and Software Specification

In this paper, we perform experiments on both multi-domain sentiment dataset<sup>①</sup> and 20 Newsgroups dataset<sup>②</sup>, which show that MVPCA/VMVPCA can not only outperform baseline non-transfer method substantially, but also improve the performance of transfer learning method SCL significantly.

We choose Multi-Domain Sentiment Dataset to show that polysemous words do exist in sentiment domain. While polysemous words incorporate conflict correspondence between words into *single-view* algorithm SCL, experiments show MVPCA/VMVPCA can solve this problem and further improve the performance of transfer learning. Multi-Domain Sentiment Dataset is crawled from Amazon, where there are four types of product reviews: Book, DVD, Electronic and Kitchen. Every product review is given a rating from 1 to 5. Reviews with rating higher than 3 are regarded as positive instances, and reviews with rating lower than 3 are labeled as negative with others discarded. After this conversion, every domain contains 1000 positive instances and 1000 negative instances. Besides these labeled instances, every domain has some unlabeled instances. Book domain has 6000 unlabeled instances, DVD has 34741, Electronic has 13153, and Kitchen domain has 16785 unlabeled instances. Details are shown in Table 4.

Before applying MVPCA and VMVPCA to Multi-Domain dataset, we briefly outline the distributional shift between domains by showing different distributions of top 50 bridge features in Book, DVD, Electronic and Kitchen domains in Fig.1. The  $X$  axis denotes different bridge features in alphabetic order. The  $Y$  axis denotes the value of mutual information. From Fig.1, we can see that all of the four domains differ substantially in the mutual information values of bridge features. In this way, Fig.1 does show some intuitive differences between domains, although the differences can be detected in more ways than mutual information. These distributional differences encourage us to apply transfer learning techniques to cross domain sentiment classification.

**Table 4.** Multi-Domain Sentiment Dataset Specification

Domain	Positive	Negative	Unlabeled	Features
Book	1000	1000	6000	42835
DVD	1000	1000	34741	162682
Electronic	1000	1000	13153	47426
Kitchen	1000	1000	16785	48534

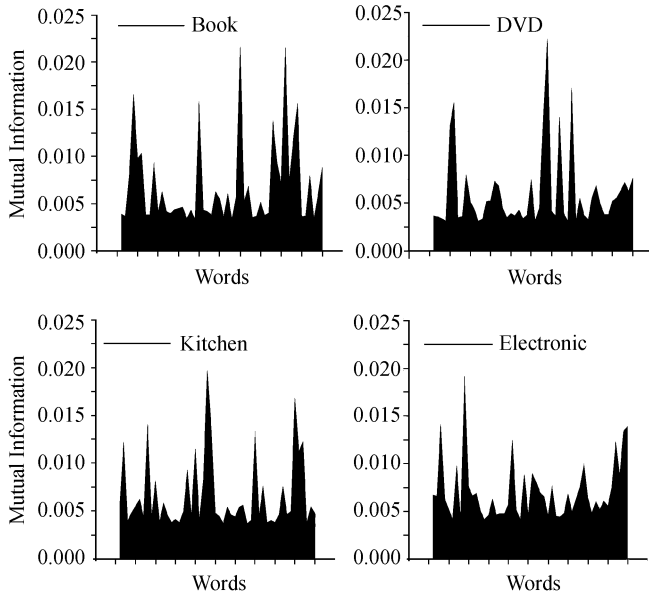


Fig.1. Distributional shift in different domains.

We choose 20 Newsgroups dataset to show that polysemous words exist not only in sentiment review

data, but also in topic domain text data. The 20 Newsgroups dataset is a set of nearly 20 000 documents collected from 20 different newsgroups. These 20 different newsgroups covers extensive topics from computer and science to forsale. These topics can be categorized into 6 top categories and 20 subcategories, with each subcategory corresponding to a newsgroup. They are shown in Table 5. In order to make the 20 Newsgroups dataset

**Table 5.** 20 Newsgroups Dataset Specification

comp.graphics	rec.autos
comp.os.ms-windows.misc	rec.motorcycles
comp.sys.ibm.pc.hardware	rec.sport.baseball
comp.sys.mac.hardware	rec.sport.hockey
comp.windows.x	
sci.crypt	talk.politics.misc
sci.electronics	talk.politics.guns
sci.med	talk.politics.mideast
sci.space	
misc.forsale	talk.religion.misc
	alt.atheism
	soc.religion.christian

**Table 6.** Cross Domain Classification Setting of the 20 Newsgroups Dataset

Dataset	Source Domain	Target Domain
comp vs. sci (C-S)	comp.graphics comp.os.ms-windows.misc sci.crypt sci.electronics	comp.sys.ibm.pc.hardware comp.sys.mac.hardware comp.windows.x sci.med sci.space
rec vs. talk (R-T)	rec.autos rec.motorcycles talk.politics.guns talk.politics.misc	rec.sport.baseball rec.sport.hockey talk.politics.mideast talk.religion.misc
rec vs. sci (R-S)	rec.autos rec.sport.baseball sci.med sci.space	rec.motorcycles rec.sport.hockey sci.crypt sci.electronics
sci vs. talk (S-T)	sci.electronics sci.med talk.politics.misc talk.religion.misc	sci.crypt sci.space talk.politics.guns talk.politics.mideast
comp vs. rec (C-R)	comp.graphics comp.sys.ibm.pc.hardware comp.sys.mac.hardware rec.motorcycles rec.sport.hockey	comp.os.ms-windows.misc comp.windows.x rec.autos rec.sport.baseball
comp vs. talk (C-T)	comp.graphics comp.sys.mac.hardware comp.windows.x talk.politics.mideast talk.religion.misc	comp.os.ms-windows.misc comp.sys.ibm.pc.hardware talk.politics.guns talk.politics.misc



suitable for our problem setting, we reorganize the 20 subcategories and put them in related but different domains. In this paper, the settings of 20 Newsgroups dataset is similar to Dai *et al.*<sup>[7]</sup> We reorganize the 20 subcategories into 6 source and target domain pairs. Within each domain pair, the texts are from only two top categories. And within each domain in the pair, positive instances consist of some subcategories in one top category while negative instances consist of some other subcategories in the other top category. Details of cross domain classification setting of the 20 Newsgroups dataset are shown in Table 6.

## 4.2 Experiments on Multi-Domain Sentiment Dataset

In order to conduct the comparisons fairly, we compare MVPCA/VMVPCA and SCL in the same experimental settings: 1) MVPCA/VMVPCA and SCL use the same parameter settings, e.g., the same number of bridge features,  $m$ , and the same dimensionality of latent semantic space,  $p$ ; 2) MVPCA/VMVPCA and SCL

utilize the same linear classifier when learning feature correspondences: logistic regression. Then we show the experimental results according to different parameter settings in Table 7 and Table 8.

In Table 7 and Table 8, “B-D” in “Domain” column means that source domain is Book domain and target domain is DVD domain. “Baseline” column denotes the accuracy of baseline method, which first trains a logistic regression model on source domain, and then evaluates the testing accuracy on target domain instances.  $m$  denotes the number of bridge features.  $p$  denotes the number of eigenvectors generated by MVPCA or SCL<sup>[1-2]</sup>, and it is also the dimensionality of cross domain semantic space. “S” denotes the accuracy of SCL<sup>[1-2]</sup>; “M” denotes the accuracy of MVPCA; “VM” denotes the accuracy of the VMVPCA algorithm described in Table 3, regarded as a mixture algorithm of SCL and MVPCA.

From Table 7 and Table 8, it is clear that MVPCA and VMVPCA win most of the time over SCL algorithm and baseline method. The values of accuracy in bold mean they are the winners, and the values in bold

**Table 7.** Classification Accuracy Comparison on Multi-Domain Sentiment Dataset (%)

Domain	Baseline	$m = 100, p = 10$			$m = 100, p = 30$			$m = 100, p = 50$			$m = 100, p = 70$		
		S	M	VM	S	M	VM	S	M	VM	S	M	VM
B-D	79.30	78.45	77.05	77.85	79.15	78.90	79.05	79.80	80.60	81.05	80.65	80.70	<b>80.90</b>
B-E	77.20	76.25	76.55	76.50	79.30	78.80	77.95	79.55	78.90	79.95	79.05	78.60	<b>80.00</b>
B-K	77.15	77.40	76.10	75.90	79.20	78.75	79.25	80.00	80.05	79.85	<b>80.70</b>	79.25	<b>80.70</b>
D-B	75.10	76.45	74.95	74.60	78.50	77.75	77.90	78.25	<b>78.80</b>	78.30	78.60	78.55	78.35
D-E	76.50	77.60	76.20	76.40	76.90	78.25	78.15	77.20	78.25	<b>79.90</b>	77.60	78.70	79.15
D-K	76.75	76.95	77.50	77.05	80.25	78.50	78.75	80.30	78.50	79.65	<b>80.85</b>	79.15	79.80
E-B	70.30	72.75	70.45	71.90	72.15	72.45	72.80	71.80	71.75	<b>73.30</b>	71.75	71.85	72.70
E-D	71.55	71.20	69.45	71.45	73.55	72.80	71.85	72.75	<b>73.80</b>	73.30	72.75	75.50	73.10
E-K	84.75	85.35	85.45	85.55	85.75	86.95	86.85	86.55	87.65	<b>87.95</b>	86.85	87.60	87.50
K-B	71.70	72.85	72.50	73.05	73.85	73.30	72.05	<b>74.70</b>	72.75	73.30	74.35	72.55	74.30
K-D	72.05	73.30	73.30	71.80	75.85	74.40	75.25	74.10	<b>76.45</b>	75.85	75.00	75.80	76.20
K-E	84.25	84.50	85.10	84.30	85.45	85.60	85.70	86.15	86.15	<b>86.55</b>	85.55	85.50	86.05

**Table 8.** Classification Accuracy Comparison on Multi-Domain Sentiment Dataset (%)

Domain	Baseline	$m = 200, p = 10$			$m = 200, p = 30$			$m = 200, p = 50$			$m = 200, p = 70$		
		S	M	VM	S	M	VM	S	M	VM	S	M	VM
B-D	79.30	79.95	79.05	78.75	79.20	79.55	79.65	80.05	81.00	80.30	<b>81.80</b>	81.65	81.75
B-E	77.20	76.40	76.95	76.70	77.70	77.30	77.30	79.35	78.60	79.05	78.95	79.35	<b>80.25</b>
B-K	77.15	76.75	76.80	76.60	78.05	78.50	78.25	81.25	80.00	80.10	81.40	80.15	<b>81.45</b>
D-B	75.10	77.05	75.95	74.90	77.30	78.60	78.40	76.40	77.15	76.65	<b>77.90</b>	77.25	<b>77.90</b>
D-E	76.50	77.50	76.20	76.20	73.30	77.65	79.10	75.20	77.65	78.70	77.15	78.30	<b>79.55</b>
D-K	76.75	77.60	76.95	77.50	80.20	78.45	78.35	79.35	77.60	77.65	<b>80.95</b>	78.35	80.35
E-B	70.30	72.45	70.45	70.85	72.35	<b>72.55</b>	72.30	72.10	71.30	71.95	71.30	70.25	72.50
E-D	71.55	71.30	69.20	69.25	73.55	72.55	72.20	71.60	74.20	73.15	73.20	<b>76.65</b>	75.35
E-K	84.75	85.70	85.70	85.45	85.35	86.80	86.25	85.95	87.10	86.55	87.00	<b>87.60</b>	87.35
K-B	71.70	71.95	72.45	73.05	73.05	71.90	72.05	73.70	72.20	73.15	<b>74.70</b>	72.05	73.20
K-D	72.05	72.60	72.15	71.25	71.20	74.10	75.70	73.90	74.75	75.20	73.30	75.25	<b>75.20</b>
K-E	84.25	84.60	83.05	83.30	83.65	85.60	85.55	85.25	85.40	85.60	85.70	<b>85.85</b>	85.65

with underline show the winners are MVPPCA or VMVPCA. To see how they perform in different settings, we need to conduct further analysis.

Firstly, we put some efforts on the win/draw/loss statistics of these methods under different settings of  $p$  shown in Table 9, where “M vs. Baseline” denotes the comparison between MVPPCA and baseline method. Similarly for VM vs. Baseline. “M-best vs. S-best” denotes the comparison between the best results of MVPPCA/VMVPCA and the best results of SCL, under different settings of  $p$ . “M/VM vs. S” compares the better results of MVPPCA and VMVPCA with SCL. “M vs. S” denotes the comparison between MVPPCA and SCL. Similarly for “VM vs. S”, “VM vs. M”. Obviously, both MVPPCA and VMVPCA outperform baseline method in almost all time of transfer learning when  $p \geq 30$ . When the parameter  $p = 10$  or 30, MVPPCA and VMVPCA do not outperform SCL or perform even worse when  $p = 10$ . When  $p > 30$ , performances of MVPPCA and VMVPCA are enhanced and they outperform SCL significantly. And they achieve the most significant improvement over SCL when  $p = 50$ . When  $m = 100$ ,  $p \geq 50$ , VMVPCA performs better than SCL and baseline at the 5% significance level, while MVPPCA performs comparable with SCL and significant better than baseline; when  $m = 200$ ,  $p \geq 30$ , VMVPCA performs better than SCL and baseline at the 5% significance level, while MVPPCA performs comparable with SCL and significant better than baseline. However, it is notable that when  $p$  reaches 70, MVPPCA degenerates its performance to be comparable with SCL. In general, when  $p$  reaches 50, MVPPCA and VMVPCA

can obviously improve SCL’s performance. The reason why they need a certain number of eigenvectors to get better knowledge transfer performance is that they performs a multi-view component analysis which deals with more complex structural information across domains than SCL does. So compared with SCL, MVPPCA and VMVPCA need more eigenvectors to depict the knowledge transfer clearly.

**Table 9.** Win/Draw/Loss Statistics Analysis on Multi-Domain Sentiment Dataset

Parameter Settings	$p = 10$	$p = 30$	$p = 50$	$p = 70$
M vs. Baseline	12/0/12	23/0/1	24/0/0	23/0/1
VM vs. Baseline	9/0/15	23/0/1	24/0/0	24/0/0
M-best vs. S-best	17/2/5	17/2/5	17/2/5	17/2/5
M/VM vs. S	9/0/15	13/0/11	17/0/7	18/0/6
M vs. S	7/2/15	11/0/13	14/1/9	12/0/12
VM vs. S	7/0/17	12/0/12	16/0/8	17/0/7
VM vs. M	11/1/12	11/1/12	16/0/8	16/0/8

Fig.2 and Fig.3 demonstrate the similar statistics of the comparison between SCL and MVPPCA/VMVPCA. X axis denotes the settings of source and target domain pairs. Y axis denotes the classification accuracy. When  $m = 100$ ,  $p = 10$  and  $m = 200$ ,  $p = 10$ , it is shown that MVPPCA/VMVPCA only achieve comparable performance with baseline method. And at the same time, SCL performs better than our algorithms. When  $p = 30$  and  $m = 200$ , SCL degenerates its performance, and both MVPPCA and VMVPCA achieve obvious improvement over SCL and baseline method.

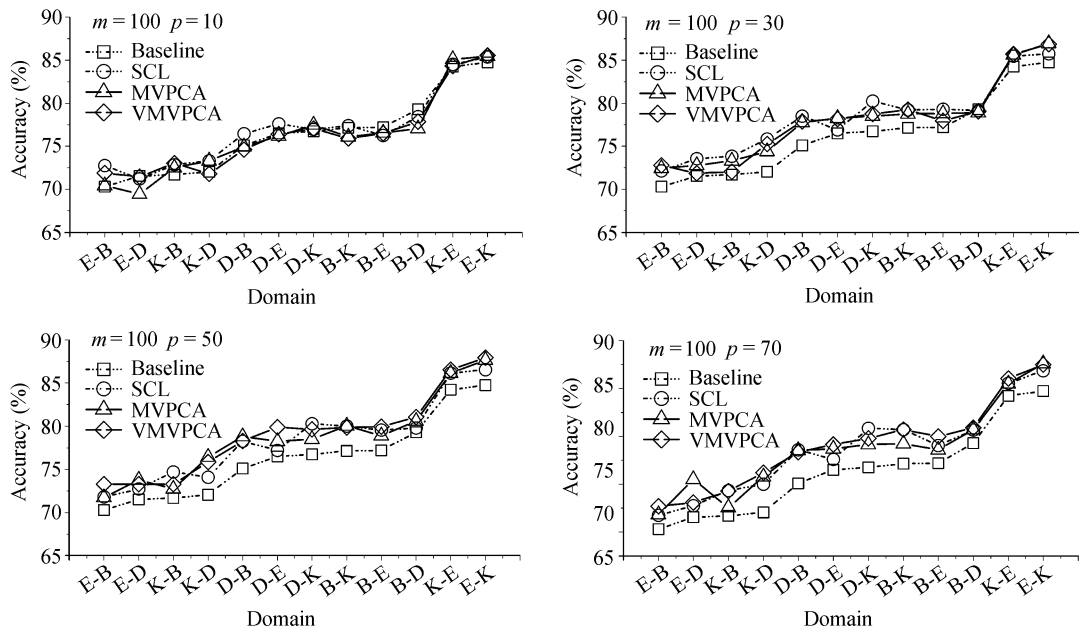


Fig.2. Transfer learning performance comparison on multi-domain sentiment dataset when  $m = 100$ .

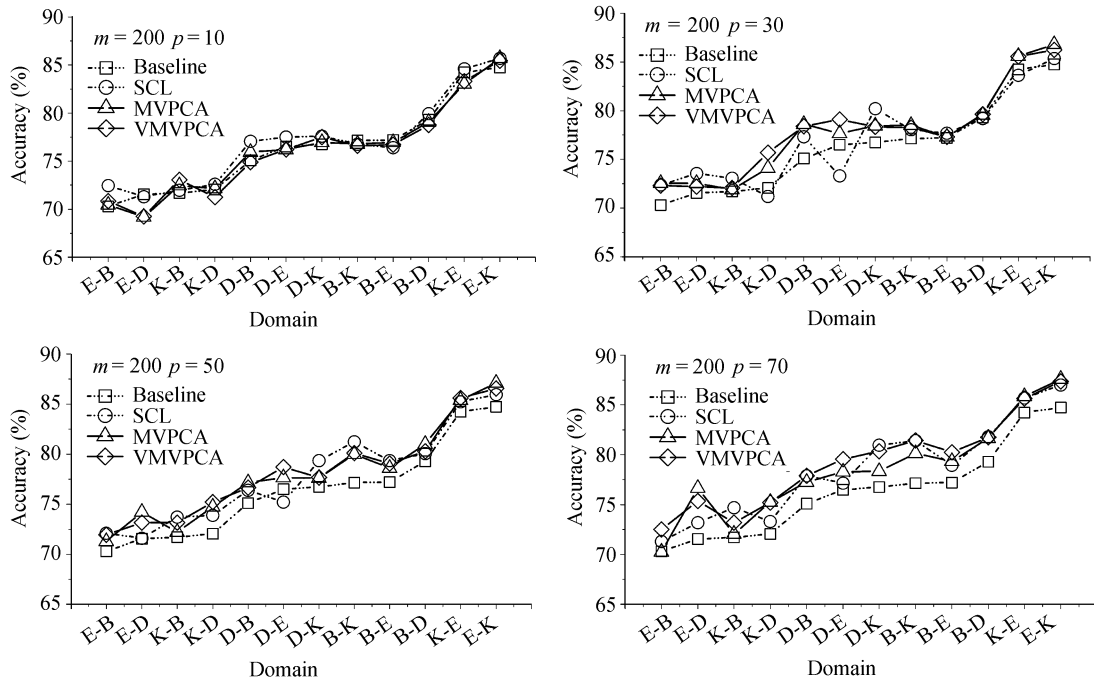


Fig.3. Transfer learning performance comparison on Multi-Domain Sentiment Dataset when  $m = 200$ .

When  $p = 30$  and  $m = 100$ , both SCL and MVPCA/VMVPCA get enhanced performance, with SCL achieving the best performance. When  $p = 50$  and  $p = 70$ , the performance curves of MVPCA/VMVPCA stay over the curve of SCL most of the time. And MVPCA/VMVPCA win 35 in 48 times of comparison when  $p = 50$  and  $p = 70$ . In general, MVPCA and VMVPCA achieve significant improvement over SCL when  $p$  reaches 50 or 70.

Then we study how the performances of MVPCA/VMVPCA are affected by parameters  $p$  and  $m$ . Table 10 shows the average classification accuracy of the algorithms in comparison under different settings of  $p$ . If we focus on one line, when  $p$  increases, the cross domain classification accuracy of every algorithm will increase. This trend makes sense, because if one classifier owns more information to benefit knowledge transfer, we can expect the continuous improvement. From Fig.4, we can also see that at first when  $p$  is small, SCL performs the best. As soon as MVPCA collects enough eigenvectors to facilitate transfer learning, MVPCA outperforms SCL. One more remark about

the parameter  $m$  is that as  $m$  increases, we cannot always expect to improve its performance. Because choose more bridge features cannot guarantee that the discriminative power increases at the same time, they will always increase the information. Bridge features within one domain have so similar explanations that 100 bridge features can give sufficient helpful information for transfer learning. This also leaves an open question for future work on choosing bridge features.

From above discussions, we get the point that MVPCA/VMVPCA need a certain number of eigenvectors to facilitate cross domain knowledge transfer, and when  $p$  reaches 50, MVPCA/VMVPCA achieve significant improvement over SCL and baseline method. So let us take a detailed look at the improvement of our algorithms. Fig.5 shows error reduction ratios of MVPCA/VMVPCA over baseline method and SCL algorithm.  $X$  axis denotes source and target domain pairs;  $Y$  axis denotes error reduction ratio, and the curves of MVPCA/VMVPCA are computed with the minimum error of MVPCA and VMVPCA. In this paragraph, we simply denote MVPCA and VMVPCA

**Table 10.** Average Performance Comparison on Multi-Domain Sentiment Dataset (%)

	$m = 100$				$m = 200$			
	$p = 10$	$p = 30$	$p = 50$	$p = 70$	$p = 10$	$p = 30$	$p = 50$	$p = 70$
Baseline	76.38							
SCL	<b>76.92</b>	<b>78.33</b>	78.43	78.64	<b>76.99</b>	77.08	77.84	78.61
MVPCA	76.22	78.04	78.64	78.65	76.24	77.80	78.08	78.56
VMVPCA	76.36	77.96	<b>79.08</b>	<b>79.06</b>	76.15	<b>77.93</b>	<b>78.17</b>	<b>79.21</b>

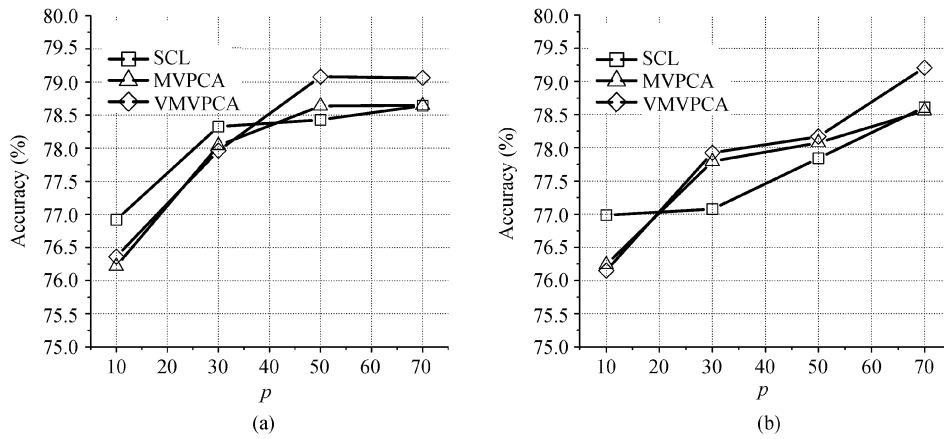


Fig.4. Average transfer learning performance comparison on multi-domain sentiment dataset. (a)  $m = 100$ . (b)  $m = 200$ .

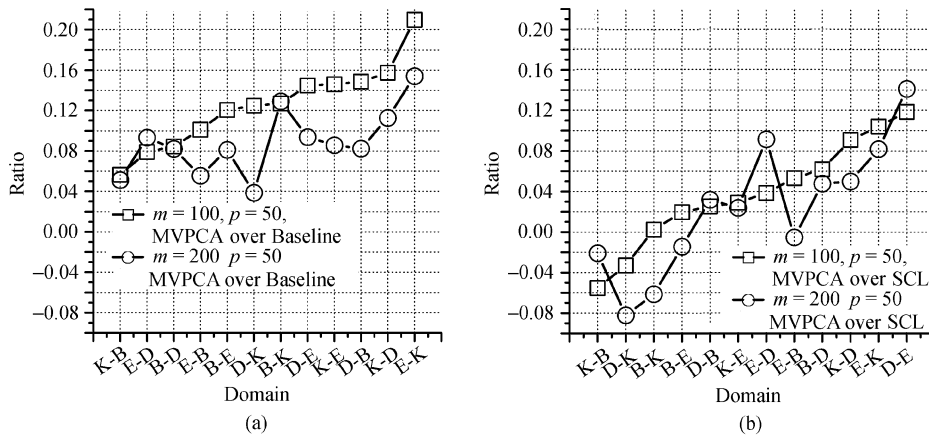


Fig.5. Error reduction ratio of MVPCA on multi-domain sentiment dataset.

together by MVPCA. Compared with baseline method, when  $p = 50$ , MVPCA always gets a positive error reduction, which means MVPCA always outperforms baseline method under this setting. MVPCA reduces the error of baseline method transferring from Electronic to Kitchen up to 21% when  $m = 100$  and  $p = 50$ , and achieves an average of 11.43% error reduction over baseline when  $m = 100$ ,  $p = 50$ . Compared with SCL algorithm, it seems that transferring from Kitchen domain to Book domain and transferring from DVD domain to Kitchen domain are two difficult tasks for MVPCA/VMVPCA, and they cannot outperform SCL algorithm in many different parameter settings. However, MVPCA can reduce SCL's error of transferring from DVD to Electronic up to 14% when  $m = 200$ ,  $p = 50$ , and MVPCA achieves an average of 2.7% error reduction over SCL when  $m = 100$ ,  $p = 50$ .

Finally, we can give the following remarks of our algorithm. When  $p$  increases to a certain extent, MVPCA collects enough information from the complex structure of *multi-view* feature correspondences, and these information keeps MVPCA away from incorrect

correlations. Thus, MVPCA can outperform the *single-view* transfer learning algorithm SCL. It is more interesting that when we fuse the *single-view* correspondences with the *multi-view* correspondences, we can learn a more powerful semantic space to benefit our knowledge transfer. Thus VMVPCA achieves an average error reduction of 11.43% over baseline method and 2.7% over SCL algorithm.

### 4.3 Experiments on 20 Newsgroups Dataset

Similar to the experiments on multi-domain sentiment dataset, we perform MVPCA/VMVPCA and SCL under the setting of the same parameters and the same base classifier: logistic regression. The experimental results of comparisons between MVPCA/VMVPCA and SCL are listed in Table 11, where "Domain" tells the setting of source and target domains, and "C-S" in this column denotes that source and target domains are from dataset "comp vs. sci" in Table 5. "Baseline", "S", "M" and "VM" are defined similarly with Table 7. On 20 Newsgroups dataset, we

**Table 11.** Classification Accuracy Comparison on 20 Newsgroups Dataset (%)

Domain	Baseline	$m = 100, p = 10$			$m = 100, p = 30$			$m = 100, p = 50$			$m = 100, p = 70$		
		S	M	VM	S	M	VM	S	M	VM	S	M	VM
C-S	68.51	68.74	68.90	72.03	66.00	70.99	72.09	65.86	70.48	<b>71.81</b>	66.20	<b>71.78</b>	70.23
C-R	81.95	88.40	86.24	87.84	<b>89.78</b>	85.47	88.73	88.55	86.92	88.17	88.12	86.11	87.33
C-T	87.03	90.95	90.59	90.65	91.50	93.32	<b>93.85</b>	91.20	93.32	93.68	88.61	91.72	93.16
R-S	72.65	75.30	76.48	78.75	76.23	77.39	78.68	76.31	78.98	<b>79.38</b>	76.81	78.15	77.82
R-T	74.29	66.58	59.80	72.77	67.23	59.66	75.41	62.28	62.87	<b>76.01</b>	66.47	62.64	75.78
S-T	75.41	73.89	70.00	73.63	74.26	68.24	74.81	75.10	69.68	75.49	<b>75.80</b>	70.81	74.99

only perform experiments when  $m = 100$  to show the superior of MVPCA/VMVPCA over SCL.

We need to conduct further analysis to see how they perform in different settings. Table 12 gives the win/draw/loss statistics of MVPCA/VMVPCA compared with other methods. When we conduct t-test, VMVPCA is significantly better than baseline method and SCL, while MVPCA is significantly better than baseline method and comparable with SCL. Fig.6 shows the comparisons between MVPCA/VMVPCA and other methods intuitively. X axis denotes different source and target domain pairs. Y axis denotes the prediction accuracy. Fig.6 shows that VMVPCA almost always stays above SCL and baseline. It also shows that R-T and S-T are two hard domains for transfer learning. While SCL performs much worse than baseline on domain R-T, which is called negative transfer, VMVPCA performs better than both SCL and baseline when  $p > 10$ .

Table 13 and Fig.7 give the overview of how these methods perform on 20 Newsgroups dataset. The numbers in bold with underline in Table 13 denote the best performances from MVPCA/VMVPCA. In Fig.7,

it can be seen that the performance curves are different from those on multi-domain sentiment dataset.

**Table 12.** Win/Draw/Loss Statistics Analysis on 20 Newsgroups Dataset

Parameter Settings	$p = 10$	$p = 30$	$p = 50$	$p = 70$
M vs. Baseline	4/0/2	4/0/2	4/0/2	4/0/2
VM vs. Baseline	4/0/2	5/0/1	6/0/0	5/0/1
(M/VM)-best vs. S-best	4/0/2	4/0/2	4/0/2	4/0/2
M/VM vs. S	3/0/3	5/0/1	5/0/1	4/0/2
M vs. S	2/0/4	3/0/3	4/0/2	3/0/3
VM vs. S	3/0/3	5/0/1	5/0/1	5/0/1
VM vs. M	6/0/0	6/0/0	6/0/0	4/0/2

**Table 13.** Average Performance Comparison on 20 Newsgroups Dataset (%)

$m = 100$	$p = 10$	$p = 30$	$p = 50$	$p = 70$
Baseline	76.64	76.64	76.64	76.64
SCL	77.31	77.50	76.55	77.00
MVPCA	75.34	75.85	77.04	76.87
VMVPCA	<b>79.28</b>	<b>80.60</b>	<b>80.76</b>	<b>79.89</b>

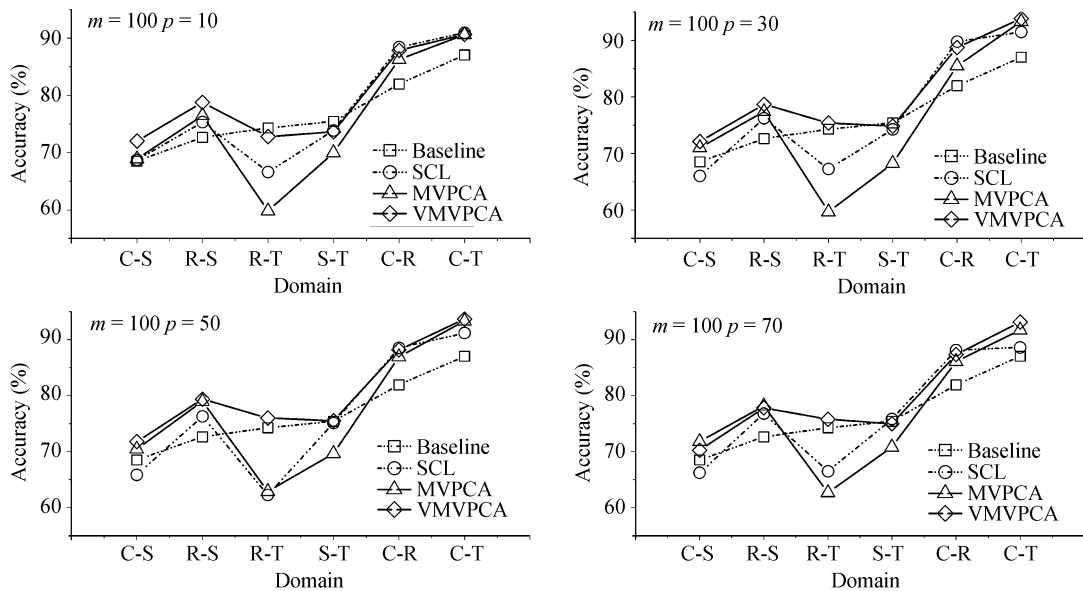


Fig.6. Transfer learning performance comparison between MVPCA/VMVPCA and other methods on 20 Newsgroups dataset.

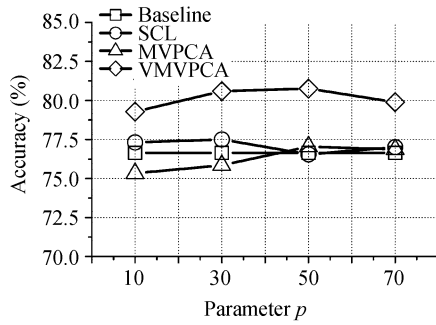
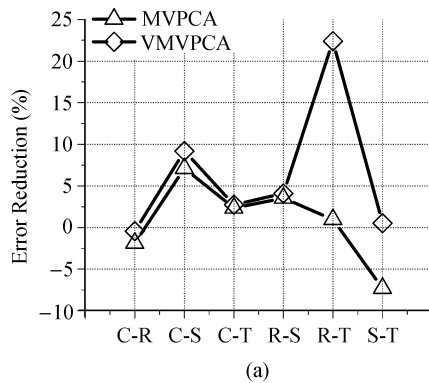


Fig.7. Average transfer learning performance comparison between baseline method, SCL algorithm, and MVPCA/VMVPCA on 20 Newsgroups dataset.

The curves in Fig.7 are rather stable compared to those in Fig.4, with VMVPCA staying higher above the curves of baseline and SCL methods. However, MVPCA needs more eigenvectors to achieve its best performance than SCL method. And this phenomenon is the same as that MVPCA shown on Multi-Domain Sentiment Dataset. Fig.8 shows the error reduction of MVPCA/VMVPCA over SCL algorithm. When  $p = 50$ , VMVPCA reduces the prediction error of SCL up to 22% on domain R-T. On 20 Newsgroups dataset, we can conclude that VMVPCA performs significantly better than SCL with MVPCA achieving comparable performance with SCL, while all of MVPCA/VMVPCA perform substantially better than baseline method. The improvement of MVPCA/VMVPCA over *single-view* transfer learning algorithm SCL and baseline non-transfer method shows that *multi-view* transfer learning technique works well in topic domain text data.

#### 4.4 Experiments on the Performance of MVPCA/VMVPCA under Different Quantities of Polysemous Words

To show MVPCA/VMVPCA outperforms SCL actually due to the problem of polysemous words, we



further study the performance of MVPCA/VMVPCA and SCL under different quantities of vague words in the dataset. We specify a polysemous word by its correlation with class labels in different domains, and to change the number of polysemous words in a domain, we artificially create some polysemous words to incorporate into the dataset. Thus the perturbed datasets provide different settings of quantities of vague words to study how the performance is affected by the number of polysemous words. Experiments show that the more polysemous words in the dataset, MVPCA/VMVPCA the better choice MVPCA/VMVPCA is when we need transfer learning techniques.

To further study the performance of MVPCA/VMVPCA under different quantities of polysemous words, we randomly choose “B-E”, “E-K” from multi-domain sentiment dataset, and “C-S”, “R-T” from 20 Newsgroups dataset. To specify whether a word  $w$  is a polysemous word, we simply compute

$$p(w = 1, y = +1) = \frac{\text{count}(w = 1, y = +1)}{\sum_{w,y} \text{count}(w, y)}$$

and

$$p(w = 1, y = -1) = \frac{\text{count}(w = 1, y = -1)}{\sum_{w,y} \text{count}(w, y)}.$$

If  $\frac{p(w=1,y=+1)}{p(w=1,y=-1)} > \delta$  with  $\delta \geq 1$ , it tells that  $w$  occurs more often with positive label. If  $\frac{p(w=1,y=+1)}{p(w=1,y=-1)} < \frac{1}{\delta}$  with  $\delta \geq 1$ , it tells that  $w$  occurs more often with negative label. In this paper, we set  $\delta = 1$ . So if  $\frac{p(w=1,y=+1)}{p(w=1,y=-1)} > 1$  occurs in source domain and meanwhile  $\frac{p(w=1,y=+1)}{p(w=1,y=-1)} < 1$  occurs in target domain, or vice versa, we regard  $w$  as a polysemous word with different meanings in different domains. Based on the definition of polysemous words, we can artificially create some new polysemous words. For example, if a word “cars”

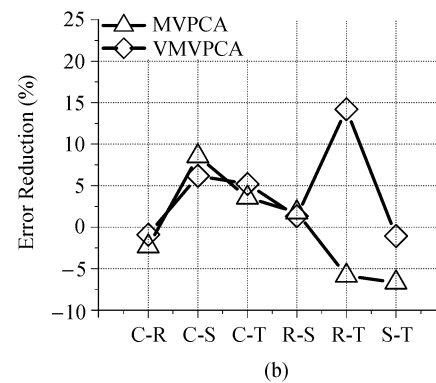


Fig.8. Error reduction of MVPCA/VMVPCA over SCL on 20 Newsgroups dataset. X axis denotes different domain pairs, Y axis denotes the percentage of prediction error reduction of MVPCA/VMVPCA over SCL. (a)  $m = 100, p = 50$ . (b)  $m = 100, p = 70$ .

has strong correlation with positive label in source domain, while “israeli” strongly correlates with negative label in target domain, we generate “cars\_israeli” as a polysemous word. Similarly, if “fbi” correlates more with negative label in source domain while “games” highly correlates with positive label in target domain, “fbi\_games” is artificially generated as a polysemous word. Table 14 illustrates some of the artificial polysemous words in different domains.

After the generation of new words, we replace those original words with these artificially generated

**Table 14.** Artificially Generated Polysemous Words in Different Datasets

B-E	excellent_not work, a must_not, wonderful_not buy, great_return, highly recommend_waste, i highly_terrible, an excellent_waste your, easy_work, a great_refund, a wonderful_after, easy_poor,
E-K	great_not, excellent_waste, price_not buy, highly_disappointed, a great_was, easy to_returned, the price_return, the best_after, perfect_your money,
C-S	windows_space, graphics_earth, dos_orbit, file_years, files_nasa, thanks_writes, microsoft_medical, image_shuttle, gif_article, win_disease
R-T	bike_jews, car_government, cars_israeli, engine_arab, ride_turkish, motorcycle_war, honda_by, bmw_people, riding_children, bikes_such

polysemous words so that we can perturb the dataset. In the experiments, we perturb the datasets with artificially generating 40, 80, 120, 160 and 200 new vague words respectively. After the perturbation, the occurrence ratio of polysemous words in the dataset varies depending on the number of artificial polysemous words. Table 15 shows the different polysemous words ratios before and after perturbation.

When we perform MVPCA/VMVPCA and other algorithms on the perturbed datasets, we get the results in Table 16 and Fig.9. In Table 16, “B-E”, “E-K”, “C-S”, “R-T” denote different sub-datasets. “AW” denotes different settings of artificially generated polysemous words. “Baseline”, “S”, “M”, “VM” denote the same results as defined in Table 7.  $m$  and  $p$  are algorithm parameters. Under different settings of  $m$  and  $p$ , we compare the cross domain prediction accuracy between

**Table 15.** Polysemous Words Ratios with Perturbation in Different Domains

No. Artificial Words	B-E	E-K	C-S	R-T
0	0.26	0.19	0.45	0.27
40	0.29	0.22	0.49	0.31
80	0.31	0.24	0.52	0.35
120	0.33	0.27	0.56	0.37
160	0.35	0.28	0.58	0.40
200	0.37	0.30	0.61	0.44

**Table 16.** Experiments on the Performance of MVPCA/VMVPCA under Different Quantities of Polysemous Words

AW	Base-line	$m = 100, p = 10$			$m = 100, p = 30$			$m = 100, p = 50$			$m = 100, p = 70$			
		S	M	VM	S	M	VM	S	M	VM	S	M	VM	
B-E	40	71.25	69.15	64.50	<b>69.20</b>	<b>72.25</b>	67.50	70.05	71.65	70.35	<b>71.70</b>	72.00	70.10	<b>72.35</b>
	80	67.75	<b>67.65</b>	63.45	66.70	<b>68.60</b>	65.50	67.75	<b>68.50</b>	65.85	68.15	<b>68.40</b>	66.90	67.95
	120	66.20	<b>66.20</b>	60.80	65.75	<b>66.95</b>	63.60	66.40	<b>67.45</b>	63.70	66.85	<b>67.30</b>	64.10	66.80
	160	65.35	<b>66.40</b>	61.45	64.75	65.85	61.20	<b>65.95</b>	<b>66.25</b>	62.75	65.35	64.95	63.15	<b>65.50</b>
	200	63.30	<b>63.35</b>	60.35	63.10	<b>64.05</b>	59.85	63.90	62.25	59.10	<b>63.45</b>	<b>61.65</b>	59.10	61.60
E-K	40	81.90	82.60	82.15	<b>82.95</b>	82.70	81.75	<b>82.90</b>	82.75	82.30	<b>82.90</b>	82.55	81.95	<b>82.90</b>
	80	82.10	81.95	81.75	<b>82.05</b>	81.15	81.20	<b>82.15</b>	<b>82.10</b>	80.70	82.00	81.60	81.40	<b>82.15</b>
	120	80.80	<b>80.45</b>	<b>80.45</b>	<b>80.45</b>	80.75	80.20	<b>81.55</b>	81.20	80.40	<b>81.55</b>	80.15	<b>81.20</b>	<b>81.20</b>
	160	80.25	<b>80.25</b>	79.60	80.00	79.05	79.80	<b>80.90</b>	79.75	<b>80.35</b>	<b>80.35</b>	80.00	79.75	<b>80.55</b>
	200	80.05	79.35	79.45	<b>79.60</b>	79.50	79.15	<b>80.45</b>	79.65	79.15	<b>80.10</b>	<b>79.90</b>	79.45	79.75
C-S	40	61.21	61.87	61.89	<b>63.16</b>	58.39	<b>63.14</b>	62.99	58.66	<b>62.69</b>	62.63	59.01	<b>61.99</b>	61.64
	80	59.80	59.60	59.27	<b>61.03</b>	55.76	<b>60.15</b>	59.37	57.06	58.68	<b>59.21</b>	<b>59.13</b>	56.76	58.31
	120	58.60	59.62	59.97	<b>60.66</b>	57.19	<b>59.87</b>	58.84	56.04	58.66	<b>58.78</b>	55.82	58.19	<b>58.82</b>
	160	58.23	<b>60.42</b>	59.27	59.87	56.27	57.43	<b>59.48</b>	55.57	57.98	<b>59.19</b>	55.67	58.37	<b>59.82</b>
	200	59.25	<b>60.95</b>	59.37	60.48	57.31	56.78	<b>59.15</b>	57.64	57.35	<b>61.52</b>	56.51	58.80	<b>61.50</b>
R-T	40	48.80	46.05	42.42	<b>48.64</b>	37.38	41.55	<b>46.61</b>	39.63	42.67	<b>45.34</b>	40.93	43.57	<b>46.67</b>
	80	44.30	<b>45.65</b>	41.55	45.51	35.95	39.72	<b>44.44</b>	36.51	40.00	<b>46.24</b>	32.60	40.23	<b>47.06</b>
	120	40.14	39.78	40.73	<b>42.90</b>	35.78	37.58	<b>38.54</b>	34.18	38.34	<b>40.51</b>	34.21	37.78	<b>42.78</b>
	160	41.24	44.02	40.90	<b>46.41</b>	36.06	39.63	<b>41.60</b>	36.60	38.48	<b>44.36</b>	35.67	40.00	<b>42.48</b>
	200	39.41	40.25	41.21	<b>43.94</b>	37.72	40.73	<b>43.07</b>	35.95	39.83	<b>42.64</b>	36.34	39.41	<b>42.03</b>

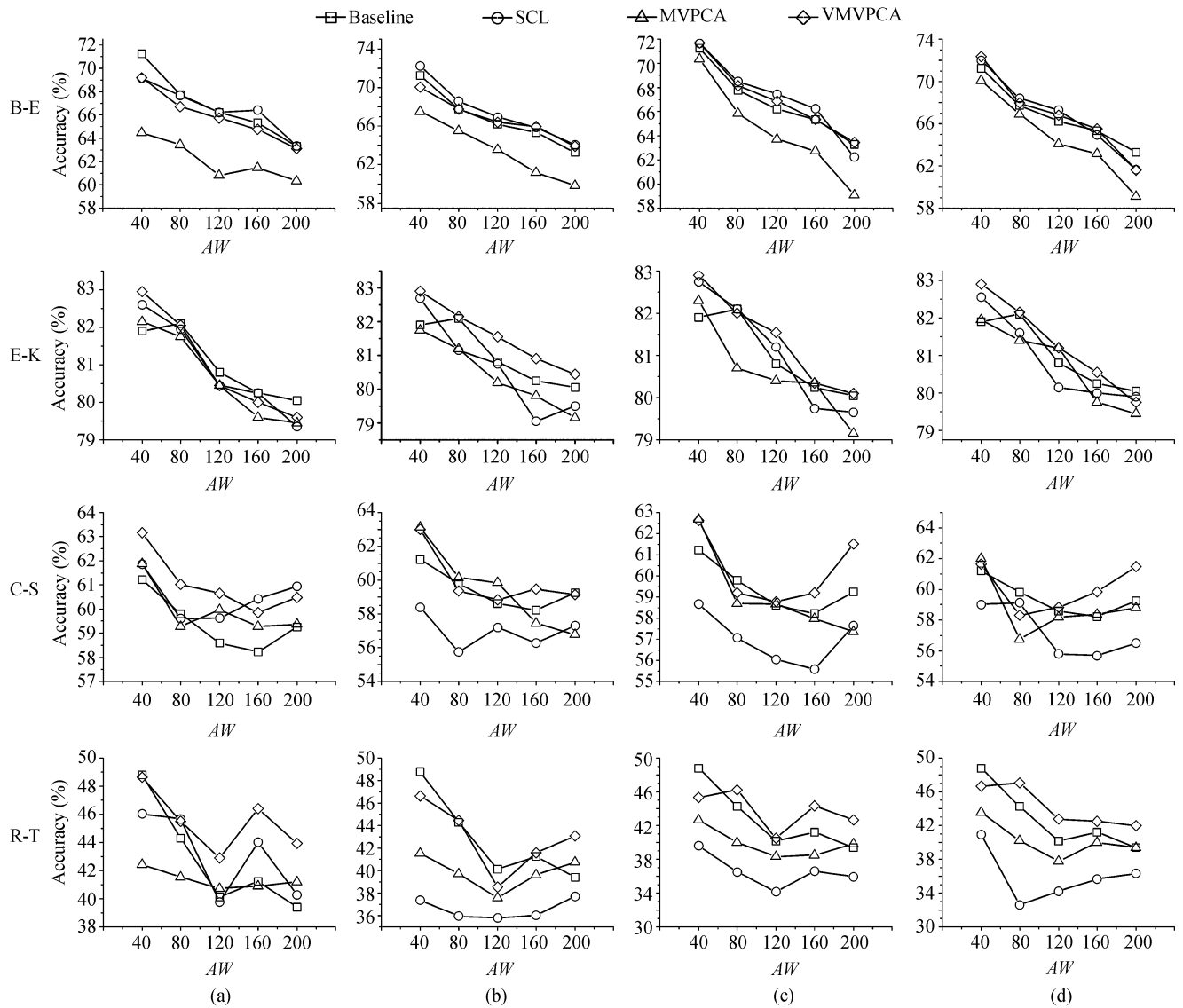


Fig.9. Transfer learning performance vs. quantity of artificial polysemous words. (Every row corresponds to a sub-dataset, which is labeled to the left of it. In each row, every subgraph illustrates the performances of all the algorithms under a fixed parameter setting. In every subgraph, X axis denotes the number of artificial generated vague words, Y axis denotes the transfer prediction accuracy on the perturbed dataset.) (a)  $m = 100$ ,  $p = 10$ . (b)  $m = 100$ ,  $p = 30$ . (c)  $m = 100$ ,  $p = 50$ . (d)  $m = 100$ ,  $p = 70$ .

SCL and MVPCA/VMVPCA. Numbers in bold with underline show that MVPCA/VMVPCA are winners; numbers in bold show that SCL wins. Table 16 illustrates that VMVPCA achieves comparable performance with SCL in the sub-dataset of B-E, while MVPCA is hurt by the artificially generated polysemous words performing worse than baseline and SCL. In the sub-dataset of E-K, while MVPCA performs comparable with SCL, VMVPCA performs better than SCL at the 5% significance level. In the sub-datasets of C-S and R-T, MVPCA/VMVPCA perform much better than SCL under different settings of  $m$ ,  $p$  and quantity of artificial polysemous words. When  $m = 100$ ,

$p = 70$  and  $AW = 80$ , VMVPCA improves the prediction accuracy of SCL up to 44.56% on R-T sub-dataset.

One important advantage of VMVPCA over SCL is that VMVPCA rarely perform worse than baseline method when  $p \geq 30$  shown in Fig.9. Another advantage of VMVPCA is that the performance of SCL decreases much faster than VMVPCA, when the number of polysemous words increases. To understand why VMVPCA does not improve SCL significantly on sub-dataset B-E due to the artificial polysemous words, we may need to refer to Blitzer *et al.*<sup>[2]</sup>, who proposes a distance measure between domains, which is called  $\mathcal{A}$ -distance. Measured by  $\mathcal{A}$ -distance, Book and Electronic



domains have so large a distance that incorporating 200 new polysemous words into B-E sub-dataset does not change it much. Thus, the advantage of VMVPCA over SCL is not significant. Another point needs to be noted is that the improvement of VMVPCA over SCL on C-S and R-T sub-datasets is much higher than on B-E and E-K sub-datasets. This difference is caused by the different polysemous word ratios of each sub-dataset. In Table 15, C-S and R-T sub-datasets have much higher polysemous words ratios than B-E and E-K sub-datasets. With higher polysemous words ratio, we can get more improvement of VMVPCA over SCL. Based on above discussions, we recommend the *multi-view* transfer learning algorithm VMVPCA compared with SCL, especially when polysemous words occur a lot in source and target domains.

## 5 Conclusion and Future Work

This paper proposes a new method called *Multi-View Principal Component Analysis* for transfer learning. MVPCA/VMVPCA separates the *single-view* of correspondences into several parts, which can be regarded as *multi-view* correspondences. In these separated correspondences, contradict correlations are eliminated, and then the cross domain correlation could be extracted by *multi-view* PCA. In this way, MVPCA/VMVPCA learns a neater and cleaner semantic space to help knowledge transfer.

However, in order to learn the semantic space of *multi-view* correspondences, Canonical Correlation Analysis<sup>[20-21]</sup> is a natural way to solve the problem. In our future work, we will study Canonical Correlation Analysis in depth, and apply it to transfer learning. What is more, we can embed the discriminative information to learn a cross domain semantic space<sup>[22-24]</sup>. Besides, considering text summarization<sup>[25]</sup>, Canonical Correlation Analysis could also be helpful. And there are some open questions, for example, how to choose a set of bridge features that are of high quality.

**Acknowledgements** We thank Yue Lu, Yu-Feng Li, Xiang-Nan Kong, Jun-Ming Xu, Yin Zhu, and Cheng-Xiang Zhai for helpful comments on this draft. We also thank anonymous reviewers for great comments on improving the work and presentation of this paper, and editors for their hard work.

## References

- [1] Blitzer J, McDonald R, Pereira F. Domain adaptation with structural correspondence learning. In *Proc. Conference on Empirical Methods in Natural Language*, Sydney, Australia, Jul. 22-23, 2006, pp.120-128.
- [2] Blitzer J, Dredze M, Pereira F. Biographies, Bollywood, Boomboxes and Blenders: Domain adaptation for sentiment classification. In *Proc. the 45th Annual Meeting of the Association of Computational Linguistics*, Prague, Czech, Jun. 25-27, 2007, pp.432-439.
- [3] Dhillon I S, Mallela S, Modha D S. Information-theoretic co-clustering. In *Proc. the 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, Washington DC, USA, Aug. 24-27, 2003, pp.89-98.
- [4] Pan S J, Tsang I W, Kwok J T, Yang Q. Domain adaptation via transfer component analysis. In *Proc. the 21st International Joint Conference on Artificial Intelligence*, Pasadena, USA, Jul. 11-17, 2009, pp.1187-1192.
- [5] Daumé H. Frustratingly easy domain adaptation. In *Proc. the 45th Annual Meeting of the Association for Computational Linguistics*, Prague, Czech, Jun. 25-27, 2007, pp.256-263.
- [6] Daumé H, Marcu D. Domain adaptation for statistical classifiers. *Journal of Artificial Intelligence Research*, 2006, 26: 101-126.
- [7] Dai W Y, Xue G R, Yang Q, Yu Y. Co-clustering based classification for out-of-domain documents. In *Proc. the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, San Jose, USA, Aug. 12-15, 2007, pp.210-219.
- [8] Pan S J, Kwok J T, Yang Q, Pan J J. Adaptive localization in a dynamic WiFi environment through multi-view learning. In *Proc. the 22nd AAAI Conference on Artificial Intelligence*, Vancouver, Canada, Jul. 22-26, 2007, pp.1108-1113.
- [9] Dai W Y, Yang Q, Xue G R, Yu Y. Boosting for transfer learning. In *Proc. the 24th International Conference on Machine Learning*, Corvallis, USA, Jun. 20-24, 2007, pp.193-200.
- [10] Huang J, Smola A, Gretton A, Borgwardt K M, Schölkopf B. Correcting sample selection bias by unlabeled data. In *Proc. the 19th Annual Conference on Neural Information Processing Systems*, Vancouver, Canada, Dec. 4-6, 2006, pp.601-608.
- [11] Zadrozny B. Learning and evaluating classifiers under sample selection bias. In *Proc. the 21st International Conference on Machine Learning*, Banff, Canada, Jul. 4-8, 2004, pp.903-910.
- [12] Jiang J, Zhai C X. Instance weighting for domain adaptation in NLP. In *Proc. the 45th Annual Meeting of the Association of Computational Linguistics*, Prague, Czech, Jun. 25-27, 2007, pp.264-271.
- [13] Raina R, Battle A, Lee H, Packer B, Ng A Y. Self-taught learning: Transfer learning from unlabeled data. In *Proc. the 24th International Conference on Machine Learning*, Corvallis, USA, Jun. 20-24, 2007, pp.759-766.
- [14] Lee H, Battle A, Raina R, Ng A Y. Efficient sparse coding algorithms. In *Proc. the 19th Annual Conference on Neural Information Processing Systems*, Vancouver, Canada, Dec. 4-6, 2006, pp.801-808.
- [15] Tan S, Cheng X, Wang Y, Xu H. Adapting naive Bayes to domain adaptation for sentiment analysis. In *Proc. the 31st European Conference on IR Research*, Toulouse, France, Apr. 7-9, 2009, pp.337-349.
- [16] Dai W Y, Xue G R, Yang Q, Yu Y. Transferring Naive Bayes classifiers for text classification. In *Proc. the 22nd AAAI Conference on Artificial Intelligence (AAAI 2007)*, Vancouver, Canada, Jul. 22-26, 2007, pp.540-545.
- [17] Sandler T, Blitzer J, Talukdar P, Pereira F. Regularized learning with networks of features. In *Proc. the 22nd Annual Conference on Neural Information Processing Systems*, Vancouver, Canada, Dec. 8-11, 2008, pp.1401-1408.
- [18] Vinokourov A, Shawe-Taylor J, Cristianini N. Inferring a semantic representation of text via cross-language correlation analysis. In *Proc. the 15th Advances of Neural Information Processing Systems*, Vancouver, Canada, Dec. 9-12, 2002, pp.1473-1480.
- [19] Li Y, Shawe-Taylor J. Using KCCA for Japanese-English

cross-language information retrieval and classification. *Journal of Intelligent Information Systems*, 2006, 27(2): 117-133.

- [20] Hardoon D R, Szedmak S, Shawe-Taylor J. Canonical correlation analysis: An overview with application to learning methods. *Neural Computation*, 2004, 16(12): 2639-2664.
- [21] Hotelling H. Relations between two sets of variates. *Biometrika*, 1936, 28(3/4): 321-377.
- [22] Diethe T, Hardoon D R, Shawe-Taylor J. Multiview Fisher discriminant analysis. In *Proc. Learning from Multiple Sources Workshop (NIPS 2008)*, Whistler, Canada, Dec. 13, 2008.
- [23] Zhou Z H, Zhan D C, Yang Q. Semi-supervised learning with very few labeled training examples. In *Proc. the 22nd AAAI Conference on Artificial Intelligence (AAAI 2007)*, Vancouver, Canada, Jul. 22-26, 2007, pp.675-680.
- [24] Zhang D Q, Zhou Z H, Chen S C. Semi-supervised dimensionality reduction. In *Proc. the 7th SIAM International Conference on Data Mining (SDM 2007)*, Minneapolis, USA, Apr. 26-28, 2007, pp.629-634.
- [25] Long C, Huang M L, Zhu X Y, Ming Li. A new approach for multi-document update summarization. *Journal of Computer Science and Technology*, 2010, 25(4): 739-749.



**Yang-Sheng Ji** received his Bachelor's degree in computer science from Nanjing University, China, in 2006 and he is now pursuing a Ph.D. degree in Department of Computer Science and Technology, Nanjing University. He visited Machine Learning Group in Microsoft Research Asia as an intern in 2008.

From November 2009 to November 2010, he visited University of Illinois at Urbana-Champaign, and worked with Prof. Cheng-Xiang Zhai. His research interests include sentiment analysis, review summarization, and transfer learning.



**Jia-Jun Chen** is a professor of the Department of Computer Science & Technology at Nanjing University, China. He received his Ph.D. degree in computer science from Nanjing University, China, in 1998. Currently, Prof. Chen is the director of the Natural Language Processing Lab at Nanjing University. He is a member of China Computer Federation (CCF), and a member of the Machine Translation Specialization Committee of Chinese Information Processing Society of China. His research interests include machine translation, text categorization and information extraction. His research has been supported by the National 863 High-Tech Projects of China, the National Natural Science Foundation of China (NSFC), the Social Science Foundation of China and the Natural Science Foundation of Jiangsu Province of China.



**Gang Niu** received his Master's degree from Department of Computer Science and Technology in Nanjing University, China. He is currently a Ph.D. candidate in Department of Computer Science in Tokyo Institute of Technology, Japan. His research interests include support vector machines, spectral clustering and manifold learning algorithms.



**Lin Shang** is an associate professor in the Department of Computer Science and Technology at the Nanjing University. She received her B.S. and M.S. degrees from Lanzhou University, China, in 1995 and 1998, and Ph.D. degree in computer science from the Nanjing University in 2004. She has visited Regina University, Canada, in May, 2007. Her

research interests mainly include artificial intelligence, machine learning, data mining, and rough sets. In these areas she has published over 50 papers in the past five years. She is the program committee member of RSKT'07, RSKT'08, CRSSC 2006, 2007, 2008, 2009, and organization committee chair/co-chair of agent2008, CRSSC 2010, China. She is a member of CCF.



**Xin-Yu Dai** received his Ph.D. degree in computer science from Nanjing University in 2005, China. He is currently an associate professor of the Department of Computer Science and Technology at Nanjing University. His research interests focus on machine translation and information retrieval. He is a member of CCF.