

Active Feature Acquisition with Supervised Matrix Completion

Sheng-Jun Huang
Nanjing University of Aeronautics
and Astronautics
Nanjing, China
huangsj@nuaa.edu.cn

Miao Xu
RIKEN Center for AIP
Tokyo, Japan
miao.xu@riken.jp

Ming-Kun Xie
Nanjing University of Aeronautics
and Astronautics
Nanjing, China
mkxie@nuaa.edu.cn

Masashi Sugiyama
RIKEN Center for AIP
University of Tokyo
Tokyo, Japan
sugi@k.u-tokyo.ac.jp

Gang Niu
RIKEN Center for AIP
Tokyo, Japan
gang.niu@riken.jp

Songcan Chen
Nanjing University of Aeronautics
and Astronautics
Nanjing, China
s.chen@nuaa.edu.cn

ABSTRACT

Feature missing is a serious problem in many applications, which may lead to low quality of training data and further significantly degrade the learning performance. While feature acquisition usually involves special devices or complex processes, it is expensive to acquire all feature values for the whole dataset. On the other hand, features may be correlated with each other, and some values may be recovered from the others. It is thus important to decide which features are most informative for recovering the other features as well as improving the learning performance. In this paper, we try to train an effective classification model with the least acquisition cost by jointly performing *active feature querying* and *supervised matrix completion*. When completing the feature matrix, a novel objective function is proposed to simultaneously minimize the reconstruction error on observed entries and the supervised loss on training data. When querying the feature value, the most uncertain entry is actively selected based on the variance of previous iterations. In addition, a bi-objective optimization method is presented for *cost-aware active selection* when features bear different acquisition costs. The effectiveness of the proposed approach is well validated by both theoretical analysis and experimental study.

1 INTRODUCTION

In data mining and machine learning tasks, a set of data objects is usually represented as a feature matrix, where each row is an object and each column is one dimension of the features. In many applications, the feature matrix may be partially observed with missing values due to various reasons [21]. For example, in disease diagnosis, a patient is an object, and the feature space consists of the physical examination results. Then some patients may selectively take some of the examinations, leaving the other features missing [20]. In wireless sensor network analysis, multiple sensors detect features of the environment in different aspects, among which, some expired sensors will cause missing values of the corresponding features [13].

Given that the feature values are severely missing, the performance of a classification model trained on such a dataset will be significantly degenerated. It is thus important to recover the missing values. The most reliable way is to acquire the ground-truth values for the missing features. Unfortunately, acquiring a feature value usually involves special devices or complex processes, leading

to high acquisition costs. Nevertheless, features are often correlated with each other, and redundant information is contained across features. Thus it may not be necessary to query all feature values. Instead, one can query a part of the features, and then recover the others from the observed entries.

Matrix completion would be a useful tool for recovering missing entries of the feature matrix, which has been extensively studied [6, 18, 34, 44]. However, existing approaches neglect the class labels, which may provide supervised information to guide the matrix completion to a desired solution. In practice, the observed entries may be noisy, and are not adequate to provide sufficient information to recover the missing values. Especially when the missing rate is high, there could be a large number of possible matrices that can well fit the observed values. The class labels, which strongly depend on the feature representations, are expected to narrow the choice over all possible matrices.

Furthermore, different features may have different contributions to recovering the missing values as well as improving the classification model. Some features are crucial while others may be less important. It is thus practical to actively select the most informative features to acquire their ground-truth values, and recover the missing values based on the observed features.

Traditional active learning algorithms select the most informative unlabeled instances to query their labels, and can significantly reduce the annotation cost [15, 32]. Similar ideas have been extended in order to perform active feature acquisition [1, 23, 30]. These methods typically try to estimate the expected utility of a feature value for improving the model performance, and then query the ground-truth value for the feature with maximum expected utility. However, some features with high potential utility can be recovered by matrix completion, and thus querying their values can be waste of acquisition costs.

In this paper, we jointly perform active feature querying and supervised matrix completion to minimize the acquisition cost. To exploit the label information for effective matrix completion, we propose an objective function that consists of the reconstruction error, the low rank regularizer and the empirical classification error. By minimizing this objective function, the recovered feature matrix is expected to on one hand well fit the structure in feature space, and on the other hand follow the label supervision to be discriminative. To select the most informative entry for active feature

acquisition, we propose a variation based criterion, which estimates the informativeness of a feature value on recovering the missing values as well as improving the classification model. Furthermore, we introduce a bi-objective optimization method to handle the case where the acquisition cost varies for different features.

Theoretical analysis is presented to give an upper bound on the reconstruction error of the proposed matrix completion algorithm. Further, experiments are performed on different datasets to validate the effectiveness of the proposed approach. Results demonstrate that our approach can recover the matrix accurately, and achieve effective classification with less feature acquisition cost.

The rest of the paper is organized as follows: we review related works in Section 2, and introduce the proposed approach in Section 3. Section 4 presents the settings and results of the experiments, followed by the conclusion in Section 5.

2 RELATED WORK

Active learning has been widely studied for reducing the labeling cost [15, 32]. Classical studies focused on designing a selection criterion such that selected instances can improve the model maximally. Informativeness is one of the most commonly used criteria, which estimates the ability of an instance in reducing the uncertainty of a statistical model. Typical techniques for informative sampling include statistical methods [7], SVM-based methods [37] and query-by-committee methods [11], etc.

Differently from traditional active learning that targets reducing the labeling cost, there is another branch of research employing similar ideas to reduce the feature acquisition cost [22]. These methods iteratively query the ground-truth values for the actively selected features, and are expected to improve the learning performance with least queries. Some methods tried to estimate the expected utility of each feature to improve the model, and then select the top features with maximum expected utility to query their values. For example, in [25], a criterion was proposed to estimate the expected improvement of accuracy per unit cost, and then the most cost-effective feature values were iteratively acquired. A similar approach was proposed in [39], where the learning task is clustering instead of classification, and thus the corresponding criterion estimates the expected improvement in clustering quality per unit cost. There is another category of methods called instance completion. Instead of querying one specific feature value, they selected a small batch of incomplete instances with missing features, and queried all missing values for the selected instances each time. The instances are actively selected aiming to improve the classification performance. For example, the authors of [31] proposed to estimate the expected utility of each instance for active selection, and also derived a probabilistic lower bound on the error reduction achieved with the proposed technique. The method in [9] chose the top k instances based on a derived upper bound on the expected distance between the next classifier and the final classifier.

A common limitation of these methods is that they do not consider the case where some of the missing features can be accurately recovered from the observed entries, and thus may waste the acquisition cost of unnecessary queries. There is one study that tried to query both missing features and labels, and built an imputation model for missing features [26]. However, it requires a complete

set of training examples for training a model, which may not be satisfied in real applications.

Matrix completion is a classical approach for recovering the missing entries of a partially observed matrix. It has been successfully applied to collaborative filtering [29], dimensionality reduction [40], multi-class/multi-label learning [2, 12], clustering [10, 43], etc. One main category of existing methods is statistical matrix completion based on the low-rank assumption [4, 6, 16, 17, 27, 41]. These methods usually transform the matrix completion task into an optimization problem, and try to find a low-rank matrix to fit the observed entries. There are some structural matrix completion methods which explicitly analyze the information contained in the observed entries and are capable of evaluating whether the observations are theoretically sufficient for recovering the missing values [18, 24, 33].

In some cases the observed entries are not enough to recover the others, and thus further queries are needed to acquire more ground-truth values for some missing entries. Given this background, there are some active learning approaches proposed to query the most informative entries for completion [30]. For example, a general framework was proposed in [5] for active matrix completion, where existing matrix completion methods can be enhanced with an uncertainty sampling strategy. In [35], the authors firstly estimated the posterior distribution with variational approximations or Markov chain Monte Carlo sampling, and then queried the entries for collaborative prediction. The algorithm in [30] unified active querying and matrix completion in a single framework. There are some other approaches which study active completion with specific requirements on the matrix [1, 23].

While all the above studies are not theoretically grounded, there are two works focusing on adaptive querying for matrix completion with theoretical results. One is [19] which firstly sampled several rows, and adaptively decided which columns are need to be fully observed. The other is [1] which actively completed a low-rank positive semi-definite matrix. Although these two works are theoretically sound, they do not consider any supervision information.

3 THE PROPOSED APPROACH

We denote by $D = \{(\mathbf{x}_i, y_i)\}$ a dataset with n instances, where \mathbf{x} is a d -dimensional real feature vector for the i -th instance and y_i is its class label. Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be the ground-truth feature matrix of the n instances, where each column represents one dimensionality of the d -dimensional feature space. Here we consider the feature missing problem, where \mathbf{X} is only partially observed. We denote by Ω the set of indices for the observed entries of \mathbf{X} . In the rest of this section, we will firstly propose a supervised matrix completion method, and then present an active feature acquisition approach.

3.1 Supervised matrix completion

We focus on the matrix completion problem under the supervised classification setting, where the task is to learn a function f for predicting the class labels of instances. Matrix completion is a challenging problem because observed entries are usually limited, and often do not contain sufficient information for recovering missing values. Since there are an arbitrary number of possible matrices that perfectly match the observed entries, external knowledge is needed

to find the optimal one closest to the ground-truth. Low-rank is a common assumption for matrix completion, which exploits the structure information in the feature space. In this paper, we further exploit the supervised information contained in class labels to guide the matrix completion to a desirable solution. Classification function f is a mapping from the feature space to the label space, and thus can be utilized to inversely transfer the label information for feature recovering. For example, given an instance with missing features and its class label, we denote by \mathbf{x} a recovered feature vector. Assuming the classifier f is reliable, if the prediction $f(\mathbf{x})$ is faraway from the ground-truth label y , then it is more likely that feature vector \mathbf{x} is not accurately recovered. Based on this motivation, we propose to minimize the empirical classification error along with the reconstruction error and the matrix rank within one unified framework, where the feature matrix and the classification model are alternately optimized.

On one hand, we want to accurately recover the ground-truth feature matrix from the partial observation of \mathbf{X} with the low-rank assumption. On the other hand, the classification model f , which is trained with the recovered matrix $\widehat{\mathbf{X}}$, is expected to have a small empirical error. Based on this argument, we define our objective function as follows.

$$\min_{\widehat{\mathbf{X}}, f} \frac{1}{2} \|\mathcal{R}_\Omega(\widehat{\mathbf{X}} - \mathbf{X})\|_{\text{F}}^2 + \lambda_1 \|\widehat{\mathbf{X}}\|_{\text{tr}} + \lambda_2 \sum_{i=1}^n \ell(y_i, f(\widehat{\mathbf{x}}_i)), \quad (1)$$

where $\mathcal{R}_\Omega : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$,

$$[\mathcal{R}_\Omega(\mathbf{X})]_{i,j} = \begin{cases} X_{i,j} & \text{if } (i,j) \in \Omega, \\ 0 & \text{otherwise,} \end{cases}$$

$\|\cdot\|_{\text{tr}}$ is the trace norm, $\|\cdot\|_{\text{F}}$ is the Frobenius norm, and $\lambda_1, \lambda_2 \geq 0$ are regularization parameters.

We assume that the loss function ℓ can be written as a function parameterized by $\widehat{\mathbf{X}}$, and it is Lipschitz smooth with respect to $\widehat{\mathbf{X}}$. One example is the linear classifier with the squared loss, i.e., $f(\mathbf{x}_i) = \mathbf{w}^\top \mathbf{x}_i$, where $\mathbf{w} \in \mathbb{R}^d$; then we have $\sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i)) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$, where $\mathbf{y} = [y_1, y_2, \dots, y_n]$ and $\|\cdot\|$ denotes the ℓ_2 norm. In the following, we will write $\sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i))$ as $\ell(\mathbf{X}, f)$ for notational simplicity. Then the optimization problem becomes

$$\min_{\widehat{\mathbf{X}}, f} \frac{1}{2} \|\mathcal{R}_\Omega(\widehat{\mathbf{X}} - \mathbf{X})\|_{\text{F}}^2 + \lambda_1 \|\widehat{\mathbf{X}}\|_{\text{tr}} + \lambda_2 \ell(\widehat{\mathbf{X}}, f), \quad (2)$$

which can be solved by alternately optimizing $\widehat{\mathbf{X}}$ and f .

When optimizing $\widehat{\mathbf{X}}$ with fixed f , we have

$$\min_{\widehat{\mathbf{X}}} \frac{1}{2} \|\mathcal{R}_\Omega(\widehat{\mathbf{X}} - \mathbf{X})\|_{\text{F}}^2 + \lambda_1 \|\widehat{\mathbf{X}}\|_{\text{tr}} + \lambda_2 \ell(\widehat{\mathbf{X}}). \quad (3)$$

We will exploit the accelerated proximal gradient descend [38] which is a classical optimization technique in trace norm minimization to solve this problem. Let

$$g(\widehat{\mathbf{X}}) = \frac{1}{2} \|\mathcal{R}_\Omega(\widehat{\mathbf{X}} - \mathbf{X})\|_{\text{F}}^2 + \lambda_2 \ell(\widehat{\mathbf{X}}),$$

and

$$h(\widehat{\mathbf{X}}, \mathbf{Z}) = g(\mathbf{Z}) + \left\langle \nabla g(\mathbf{Z}), \widehat{\mathbf{X}} - \mathbf{Z} \right\rangle + \lambda_1 \|\widehat{\mathbf{X}}\|_{\text{tr}},$$

where

$$\nabla g(\mathbf{Z}) = \mathcal{R}_\Omega(\mathbf{Z} - \widehat{\mathbf{X}}) + \lambda_2 \frac{\partial \ell}{\partial \mathbf{Z}}.$$

We summarize the main steps here:

- Choose $\theta_0 = \theta_{-1} \in (0, 1]$, $L > 1$, $\widehat{\mathbf{X}}_0 = \widehat{\mathbf{X}}_{-1}$, $\gamma > 1$. Set $k = 0$.
- In the k -th iteration,
 - Set $\mathbf{Z}_k = \widehat{\mathbf{X}}_k + \theta_k(\theta_{k-1}^{-1} - 1)(\widehat{\mathbf{X}}_k - \widehat{\mathbf{X}}_{k-1})$.
 - Set $\widehat{\mathbf{X}}_{k+1} = \operatorname{argmin}_{\widehat{\mathbf{X}}} \left\{ h(\widehat{\mathbf{X}}, \mathbf{Z}_k) + \frac{L}{2} \|\widehat{\mathbf{X}} - \mathbf{Z}_k\|_{\text{F}}^2 \right\}$.
 - While $g(\widehat{\mathbf{X}}_{k+1}) + \lambda_1 \|\widehat{\mathbf{X}}_{k+1}\|_{\text{tr}} > h(\widehat{\mathbf{X}}_{k+1}, \mathbf{Z}_k) + \frac{L}{2} \|\widehat{\mathbf{X}}_{k+1} - \mathbf{Z}_k\|_{\text{F}}^2$:
 - * Increase $L = \gamma L$.
 - * Update $\widehat{\mathbf{X}}_{k+1} = \operatorname{argmin}_{\widehat{\mathbf{X}}} \left\{ h(\widehat{\mathbf{X}}, \mathbf{Z}_k) + \frac{L}{2} \|\widehat{\mathbf{X}} - \mathbf{Z}_k\|_{\text{F}}^2 \right\}$.
 - Set $\theta_{k+1} = \sqrt{\theta_k^4 + 4\theta_k^2} - \theta_k^2/2$.
 - Update $k = k + 1$.

The iteration continues until convergence. In the above steps, we have not specified how to obtain $\widehat{\mathbf{X}}_{k+1}$ and next we will explain this. We rewrite the problem as

$$\min_{\widehat{\mathbf{X}}} \left\langle \nabla g(\mathbf{Z}_k), \widehat{\mathbf{X}} - \mathbf{Z}_k \right\rangle + \frac{L}{2} \|\widehat{\mathbf{X}} - \mathbf{Z}_k\|_{\text{F}}^2 + \lambda_1 \|\widehat{\mathbf{X}}\|_{\text{tr}}, \quad (4)$$

which is equivalent to

$$\min_{\widehat{\mathbf{X}}} \frac{L}{2} \left\| \widehat{\mathbf{X}} - \left(\mathbf{Z}_k - \frac{1}{L} \nabla g(\mathbf{Z}_k) \right) \right\|_{\text{F}}^2 + \lambda_1 \|\widehat{\mathbf{X}}\|_{\text{tr}}. \quad (5)$$

This can be solved by Singular Value Thresholding (SVT) [3], which performs singular value decomposition on $\mathbf{Z}_k - \frac{1}{L} \nabla g(\mathbf{Z}_k) = \mathbf{U}\Sigma\mathbf{V}^\top$. Let $\widehat{\Sigma}_{ii} = \max(0, \Sigma_{ii} - \frac{\lambda_1}{L})$, the solution is given by $\mathbf{U}\widehat{\Sigma}\mathbf{V}^\top$.

Finally, the classification model f is optimized with fixed $\widehat{\mathbf{X}}$, which can be efficiently solved using existing algorithms. These two procedures are repeated until convergence.

3.2 Active feature acquisition

In this subsection, we discuss how to actively query the ground-truth values as most informative features, with the target of improving the model mostly based on the smallest number of queries. We will first present a novel criterion for estimating the informativeness of a feature, and then introduce a method to handle the case where the acquisition cost varies for different features.

3.2.1 Variance-based selection. In traditional active learning, if the model is less certain about the prediction on an instance, then the instance is considered to be more informative for improving the model, and will be more likely to be selected for label querying [15]. Inspired by this idea, we also propose an uncertainty criterion to estimate the informativeness of a feature. The challenge here is that the informativeness should reflect the usefulness of a feature both for recovering other entries and for training the classification model. Notice that the objective function defined in Eq. (1) does consider the two aspects simultaneously. At each iteration of active learning, after a small batch of feature values is acquired, the algorithm in Section 3.1 will be employed to optimize Eq. (1) for matrix completion. The output of the matrix completion may vary from iteration to iteration. If the variance of an entry over iterations is large, it implies that the entry can not be certainly decided by the algorithm, and thus may contain more useful information to recover the feature matrix and optimize the classification model. Denoting by \mathbf{X}^t the completed matrix at the t -th iteration,

Algorithm 1 The AFASMC algorithm

```

1: Input:
2:    $D$ : the data set of  $n$  instances, with  $\widehat{X}$  as the feature matrix
3:    $\Omega$ : the set of indices for observed entries
4:    $\lambda_1, \lambda_2 \geq 0$ : the parameters
5: Process:
6:   For:  $t = 1 : T$ 
7:     Repeat
8:       recover the matrix  $\widehat{X}$  with fixed classification model  $f$ 
9:       optimize the model  $f$  with fixed  $\widehat{X}$ 
10:    Until convergence
11:    For each missing entry  $\widehat{X}_{i,j}$ 
12:      calculate the average value over the  $t - 1$  iterations
13:      calculate the informativeness of  $\widehat{X}_{i,j}$  as Eq. (6)
14:    End For
15:    select a batch of entries with maximum informativeness
16:    query the ground-truth values for the selected entries
17:    set the corresponding elements of  $\Omega$  to 1
18:  End For

```

the informativeness of the j -th feature of \mathbf{x}_i is defined as:

$$I_{i,j} = \sum_{t=1}^T (X_{i,j}^t - \bar{X}_{i,j})^2, \quad (6)$$

where $\bar{X}_{i,j} = \frac{1}{T} \sum_{t=1}^T X_{i,j}^t$ is the mean value of $X_{i,j}$ over all iterations. Then a small batch of most uncertain features with largest informativeness is selected to query their ground-truth values. The pseudo code of the algorithm for active feature acquisition is summarized in Algorithm 1. We call the proposed algorithm Active Feature Acquisition Supervised Matrix Completion (AFASMC).

Note that it is not necessary to calculate the variance based on all iterations. Generally speaking, it is more important to capture the change of an entry within recent iterations. For example, if an entry has a large variance at an early stage, but becomes stable after a few queries, it implies that this entry may have been well recovered from the recently acquired features, and thus does not need to be queried any more. We will discuss this in the experiments in more detail.

3.2.2 Cost-aware selection. Finally, we discuss a more complicated case, where the cost of acquiring a feature value varies for different features. This is a common case in real applications. For example, it is much more costly to perform an fMRI scan than blood examination for diagnosing a patient. While there is typically a conflict between the informativeness and acquisition cost of a feature, we propose to balance these two factors for achieving the best cost-effectiveness. We denote the cost for acquiring the j -th dimension of the features by C_j . Note here we assume that the acquisition cost is independent of the instance. We offer two optional strategies to consider the acquisition cost. The most straightforward method is to simply divide the informativeness by the acquisition cost. So we can have the selection strategy as:

$$\operatorname{argmax}_{(i,j) \notin \Omega} \frac{I_{i,j}}{C_j}. \quad (7)$$

This strategy provides a simple solution for cost-aware selection, but may fail when one of the two factors dominates the other.

In what follows, we introduce another solution by bi-objective optimization. In each iteration of our algorithm, we select a small batch of missing entries of the feature matrix to acquire their ground-truth values. This is a typical subset selection problem. Generally, a subset selection problem tries to select a subset S from a large set V with an objective function \mathcal{J} and a constraint of the subset size. It can be formalized as

$$\operatorname{argmin}_{S \subseteq V} \mathcal{J}(S) \quad \text{s.t.} \quad |S| \leq b, \quad (8)$$

where $|\cdot|$ denotes the size of a set, and b is the maximum number of selected elements. Further, for convenience of presentation, the subset selection problem is reformulated as optimizing a binary vector. We introduce a binary vector $s \in \{0, 1\}^n$ to indicate the subset membership, where $s_i = 1$ if the i -th element in V is selected, and $s_i = 0$ otherwise. Following the method in [28], the subset selection problem in Eq. (8) can be written as a bi-objective minimization problem:

$$\operatorname{argmin}_{s \in \{0, 1\}^n} (\mathcal{J}_1(s), \mathcal{J}_2(s)), \quad (9)$$

$$\mathcal{J}_1(s) = \begin{cases} +\infty & \text{if } s = \{0\}^n \text{ or } |s| \geq 2b, \\ \mathcal{J}(s) & \text{otherwise.} \end{cases} \quad \mathcal{J}_2(s) = |s|,$$

where $|s|$ denotes the number of 1s in s . Obviously, the problem is for sparse selection with the target of minimizing \mathcal{J} . Here \mathcal{J}_1 is set to $+\infty$ to avoid trivial solutions or over-sized subsets. In our case, we want to maximize the informativeness in Eq. (6), and at the same time minimize the acquisition cost of the selected entries. We thus can redefine the two objective functions \mathcal{J}_1 and \mathcal{J}_2 correspondingly, and have the following bi-objective optimization problem.

$$\operatorname{argmin}_{s \in \{0, 1\}^n} (\mathcal{J}_1(s), \mathcal{J}_2(s)), \quad (10)$$

$$\mathcal{J}_1(s) = \begin{cases} +\infty & \text{if } s = \{0\}^n \text{ or } \mathcal{J}_2(s) \geq 2b, \\ -\sum_{ij} s(i,j) \cdot I_{ij} & \text{otherwise.} \end{cases}$$

$$\mathcal{J}_2(s) = \sum_{ij} s(i,j) \cdot C_j.$$

Here b is the budget for the acquisition cost in each iteration, and $s(i,j)$ is used to denote the element of s corresponding to the entry of the i -th row and j -th column in matrix X . Again, $\mathcal{J}_1(s)$ is set to $+\infty$ to exclude trivial or over-cost solutions. We employ a recently proposed Pareto optimization algorithm called Pareto Optimization for Subset Selection (POSS) [28] to solve this problem. POSS is an evolutionary style algorithm, which maintains a solution archive, and iteratively update the archive by replacing some solutions with better ones. In detail, it initializes the archive with a solution of empty subset selection. In each iteration, a solution s is selected from the current archive, and a new solution s' is generated by randomly flipping bits of s . The two objective values $\mathcal{J}_1(s')$ and $\mathcal{J}_2(s')$ are then computed to compare s' with the archived solutions. Specifically, if there exists one solution s in the archive that satisfies both the following conditions:

$$\begin{aligned} \mathcal{J}_1(s) &\leq \mathcal{J}_1(s') \text{ and } \mathcal{J}_2(s) \leq \mathcal{J}_2(s'), \\ \mathcal{J}_1(s) &< \mathcal{J}_1(s') \text{ or } \mathcal{J}_2(s) < \mathcal{J}_2(s'), \end{aligned}$$

then s' will be ignored; otherwise, s' will be added to the solution archive, and at the same time all the archived solutions s that satisfy

$$\mathcal{J}_1(s') \leq \mathcal{J}_1(s) \text{ and } \mathcal{J}_2(s') \leq \mathcal{J}_2(s)$$

will be removed from the solution archive. This process is repeated until reaching a specified number of iterations. At last, the best solution with the minimal value on J_1 and within the cost budget will be selected as the final solution.

3.3 Theoretical analysis

In this subsection, we will present a theoretical bound on the reconstruction error of the supervised matrix completion method introduced in Section 3.1. For the loss between $\mathbf{X}\mathbf{w}$ and \mathbf{y} , i.e., the term $\sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i))$ in Eq. (1), here we discuss a more strict case by enforcing $\mathbf{X}\mathbf{w}$ and \mathbf{y} to be equal. It is reasonable to relax this strict constraint as in Eq. (1) to cope with possible noises. The relaxation is also benefited by more flexible choice of the loss function, for example $\|\mathbf{X}\mathbf{w} - \mathbf{y}\|$, as well as ease of optimization. For convenience of presentation, we rewrite the noiseless counterpart of Eq. (1) as:

$$\min_{\widehat{\mathbf{X}}} \sum_{(i,j) \in \Omega} (\widehat{X}_{i,j} - X_{i,j})^2 \quad \text{s.t. } \|\widehat{\mathbf{X}}\|_{\text{tr}}^2 \leq \beta\sqrt{r}nd, \quad f(\widehat{\mathbf{X}}) = \mathbf{y}, \quad (11)$$

where β and r are constants. We assume $\widehat{\mathbf{X}}^*$ is the optimal solution for Eq. (11), and try to analyze the difference between the solution of our algorithm and the optimal solution. Before discussing the property of the solution, we first define the *coherence* of a matrix, which will be used later.

DEFINITION 1. For a rank- r matrix $\mathbf{M} \in \mathbb{R}^{n \times m}$ whose SVD is $\mathbf{M} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$, we use the following value as the coherence,

$$\mu(\mathbf{M}) = \max\left\{ \max_{1 \leq i \leq n} \|\mathbf{U}_{i,*}\|, \max_{1 \leq j \leq m} \|\mathbf{V}_{j,*}\| \right\},$$

where $\mathbf{U}_{i,*}$ ($\mathbf{V}_{j,*}$) denotes the i th (j th) row of \mathbf{U} (\mathbf{V}).

Note that compared to [4, 42], for ease of use, in this paper we do not normalize the coherence by the size of the matrix. Coherence measures how the values of the entries are distributed in a matrix. The lower the coherence is, the more average the values of the entries are distributed. Apparently, if there is no entry that has a ‘‘peak’’ value in a matrix, the matrix is easier to be completed with partial observations. Based on this definition, we give our theoretical results in Theorem 1.

THEOREM 1. Suppose that $\|\mathbf{X}\|_{\text{tr}}^2 \leq \beta\sqrt{r}nd$, $f(\mathbf{X}) = \mathbf{y}$ and Ω is chosen independently at random following a binomial model with probability $|\Omega|/(nd)$. Let $\widehat{\mathbf{X}}^*$ be the solution to the optimization problem Eq. (11) and $\mu = \max_{\widehat{\mathbf{X}} \in G} \mu(\widehat{\mathbf{X}})$, where $G \subset \mathbb{R}^{n \times d}$ is $G = \{\widehat{\mathbf{X}} \in \mathbb{R}^{n \times d} \mid \|\widehat{\mathbf{X}}\|_{\text{tr}}^2 \leq \beta\sqrt{r}nd, f(\widehat{\mathbf{X}}) = \mathbf{y}\}$ for some $r \leq \min\{n, d\}$, and $\beta \geq 0$. Then with probability at least $1 - C/(n+d)$, we have

$$\frac{1}{nd} \|\widehat{\mathbf{X}}^* - \mathbf{X}\|_{\text{F}}^2 \leq 2 \left(C_0 \mu^2 \beta \sqrt{\frac{r(n+d)}{|\Omega|}} \sqrt{1 + \frac{(n+d) \log(n+d)}{|\Omega|}} \right).$$

Theorem 1 provides an upper bound of the reconstruction error for the proposed supervised matrix completion algorithm. Moreover, it is obvious that a smaller upper bound can be expected by increasing $|\Omega|$. This also motivates us to iteratively acquire more

feature values. Below we present a sketch of proof of Theorem 1. A detailed proof is available in a longer version on arXiv [14].

To prove Theorem 1, we firstly define $\mathcal{L}_\Omega(\widehat{\mathbf{X}}) = -\sum_{(i,j) \in \Omega} (\widehat{X}_{i,j} - X_{i,j})^2$ and $\bar{\mathcal{L}}_\Omega(\widehat{\mathbf{X}}) = \mathcal{L}_\Omega(\widehat{\mathbf{X}}) - \mathcal{L}_\Omega(\mathbf{X})$. Note that because \mathbf{X} is a constant matrix, subtracting $\mathcal{L}_\Omega(\mathbf{X})$ will not affect optimization of the objective function, i.e., they will both have the same $\widehat{\mathbf{X}}^*$ leading to the optimum. Then we will use the following three lemmas to prove Theorem 1:

LEMMA 1. Assume that $\|\mathbf{X}\|_{\text{tr}}^2 \leq \beta\sqrt{r}nd$, and $\mu = \max_{\widehat{\mathbf{X}} \in G} \mu(\widehat{\mathbf{X}})$, then with probability at least $1 - C/(n+d)$ we have

$$\sup_{\widehat{\mathbf{X}} \in G} \left\| \bar{\mathcal{L}}_\Omega(\widehat{\mathbf{X}}) - \mathbb{E}[\bar{\mathcal{L}}_\Omega(\widehat{\mathbf{X}})] \right\| \leq \left(C_0 \mu^2 \beta \sqrt{r} \sqrt{|\Omega|(n+d) + nd \log(n+d)} \right),$$

where the expectation is over the choice of Ω .

We can also easily derive the following result from [8]:

LEMMA 2. If $E \in \mathbb{R}^{d_1 \times d_2}$, and each entry $E_{i,j}$ is a Radamacher random variable; $\Delta \in \{0, 1\}^{d_1 \times d_2}$, and each entry $\Delta_{i,j}$ is independently sampled when $\Delta_{i,j} = 1$ with probability $n/(d_1 d_2)$ and 0 with probability $1 - n/(d_1 d_2)$. Then we have

$$\mathbb{E} \left[\|E \circ \Delta\|^h \right] \leq C 2^h (1 + \sqrt{6})^h \left(\frac{n(d_1 + d_2) + d_1 d_2 \log(d_1 + d_2)}{d_1 d_2} \right)^{h/2},$$

provided that $h \geq 1$ and C is a constant.

Further, the trace norm of the Hadamard product of two matrices is bounded as follows:

LEMMA 3. Assume that there are two matrices \mathbf{A} and \mathbf{B} that have the same shape, then we have $\|\mathbf{A} \circ \mathbf{B}\|_{\text{tr}} \leq \mu(\mathbf{A})^2 \|\mathbf{A}\|_{\text{tr}} \|\mathbf{B}\|_{\text{tr}}$, where \circ is the Hadamard product.

The proof of the lemmas is available in a longer version at arXiv [14]. Next we show how from Lemma 1 we can derive our Theorem 1. Note that for any choice of matrix $\mathbf{A} \in G$, we have

$$\begin{aligned} \mathbb{E} \left[\bar{\mathcal{L}}_\Omega(\mathbf{A}) - \bar{\mathcal{L}}_\Omega(\mathbf{X}) \right] &= \mathbb{E} \left[\mathcal{L}_\Omega(\mathbf{A}) - \mathcal{L}_\Omega(\mathbf{X}) \right] \\ &= \mathbb{E} \left[- \sum_{(i,j) \in \Omega} (A_{i,j} - X_{i,j})^2 \right] = - \frac{|\Omega|}{nd} \left(\sum_{(i,j)} (A_{i,j} - X_{i,j})^2 \right), \end{aligned}$$

where the expectation is over Ω . We can also have

$$\begin{aligned} \bar{\mathcal{L}}_\Omega(\mathbf{A}) - \bar{\mathcal{L}}_\Omega(\mathbf{X}) &= \bar{\mathcal{L}}_\Omega(\mathbf{A}) - \mathbb{E} \left[\bar{\mathcal{L}}_\Omega(\mathbf{A}) \right] - \bar{\mathcal{L}}_\Omega(\mathbf{X}) \\ &+ \mathbb{E} \left[\bar{\mathcal{L}}_\Omega(\mathbf{X}) \right] + \mathbb{E} \left[\bar{\mathcal{L}}_\Omega(\mathbf{A}) - \bar{\mathcal{L}}_\Omega(\mathbf{X}) \right] \\ &\leq 2 \sup_{\mathbf{A} \in G} \left\| \bar{\mathcal{L}}_\Omega(\mathbf{A}) - \mathbb{E}[\bar{\mathcal{L}}_\Omega(\mathbf{A})] \right\| - \frac{|\Omega|}{nd} \left(\sum_{(i,j)} (A_{i,j} - X_{i,j})^2 \right). \end{aligned}$$

Replacing \mathbf{A} by $\widehat{\mathbf{X}}^*$ which is the optimal solution to Eq. (11) and noting that $\bar{\mathcal{L}}_\Omega(\widehat{\mathbf{X}}^*) \geq \bar{\mathcal{L}}_\Omega(\mathbf{X})$, we have

$$\frac{|\Omega|}{nd} \left(\sum_{(i,j)} (A_{i,j} - X_{i,j})^2 \right) \leq 2 \sup_{\mathbf{A} \in G} \left\| \bar{\mathcal{L}}_\Omega(\mathbf{A}) - \mathbb{E}[\bar{\mathcal{L}}_\Omega(\mathbf{A})] \right\|.$$

Using Lemma 1, with probability at least $1 - \frac{C}{n+d}$, we have

$$\frac{|\Omega|}{nd} \|\widehat{\mathbf{X}}^* - \mathbf{X}\|_{\text{F}}^2 \leq 2 \left(C_0 \mu^2 \beta \sqrt{r} \sqrt{|\Omega|(n+d) + nd \log(n+d)} \right).$$

Table 1: The comparison results on matrix completion. The reconstruction error as well as classification accuracy are reported with 60% and 80% entries observed respectively.

Data	Observed Rate	AFASMC	OptSpace	LmaFit	NNLS	AFASMC	OptSpace	LmaFit	NNLS
Reconstruction Error					Test Accuracy (%)				
abalone	60%	0.13±0.01	0.38±0.01	0.14±0.00	0.14±0.00	78.5±1.2	71.8±0.8	78.3±1.3	78.4±1.2
	80%	0.07±0.00	0.23±0.01	0.09±0.00	0.07±0.00	79.7±0.7	76.3±1.2	79.5±0.6	79.5±0.9
letter	60%	0.18±0.00	0.33±0.01	0.24±0.00	0.23±0.00	98.5±0.6	92.1±3.7	97.0±1.3	94.3±1.0
	80%	0.11±0.00	0.29±0.00	0.17±0.00	0.17±0.01	99.3±0.2	94.2±0.8	99.1±0.4	94.6±1.1
image	60%	0.25±0.00	0.54±0.01	0.36±0.05	0.56±0.03	79.3±2.0	67.1±2.2	75.5±2.1	70.1±3.5
	80%	0.16±0.00	0.51±0.00	0.18±0.00	0.23±0.08	81.0±2.3	68.1±1.5	79.8±2.5	80.5±3.2
chess	60%	0.43±0.00	1.57±0.01	0.48±0.00	1.28±0.02	94.3±0.9	52.1±1.8	93.3±0.8	53.8±1.7
	80%	0.29±0.00	1.63±0.01	0.33±0.00	1.21±0.01	94.8±0.4	52.8±1.9	94.7±0.7	77.9±3.6
HillValley	60%	0.04±0.00	0.06±0.00	0.04±0.00	0.03±0.00	51.5±3.2	48.7±3.5	49.4±3.7	50.7±2.8
	80%	0.03±0.00	0.06±0.00	0.03±0.00	0.03±0.00	51.1±2.1	49.5±3.3	50.8±2.3	49.7±2.0
HTRU2	60%	0.29±0.00	0.59±0.00	0.63±0.00	0.29±0.00	97.2±0.2	90.8±0.3	97.1±0.2	97.1±0.2
	80%	0.16±0.00	0.39±0.00	0.45±0.00	0.17±0.00	97.4±0.2	94.3±0.2	97.3±0.2	97.4±0.2

Further applying $\sqrt{nd} \leq n + d$, we have

$$\frac{1}{nd} \|\widehat{\mathbf{X}}^* - \mathbf{X}\|_F^2 \leq 2 \left(C_0 \mu^2 \beta \sqrt{\frac{r(n+d)}{|\Omega|}} \sqrt{1 + \frac{(n+d) \log(n+d)}{|\Omega|}} \right).$$

4 EXPERIMENTS

In this section, we experimentally investigate the proposed method.

4.1 Settings

We perform experiments on 6 benchmarks, namely *abalone*, *letter*, *image*, *chess*, *HillValley* and *HTRU2*. The number of entries in the matrix varies from 22,960 to 143,184. For each dataset, we randomly separate the set into two subsets, one with 70% examples for training, and the other one with 30% examples for testing. We repeat the random partition 10 times and report the average results.

In the experiments, we examine the performance both on matrix completion and the classification after active queries. The proposed supervised matrix completion algorithm AFASMC is compared with following methods: OptSpace [17]—a low-rank matrix completion method based on spectral techniques and manifold optimization; LmaFit [41]—a low-rank factorization model based on the nonlinear successive over-relaxation (SOR) algorithm; NNLS [36]—an accelerated proximal gradient algorithm for low-rank matrix completion.

Also our active feature acquisition method AFASMC is compared with the following methods: QBC [5]—an active matrix completion using Query by Committee strategy; Stability [5]—an active matrix completion method based on committee stability; EM Inference [26]—it selects the instances with maximum expected utility; Random—randomly select features.

For AFASMC, the parameters λ_1 and λ_2 are fixed to 1 as default on all datasets. For other methods, parameters are set or tuned as suggested in the corresponding literature. We employ the linear SVM with default parameters as the classifier for all baselines.

4.2 Results on matrix completion

Firstly, we examine the effectiveness of the proposed method for supervised matrix completion. The performances are evaluated

with the matrix reconstruct error as well as the classification accuracy. For each dataset, we compare all the methods under different missing rates. The results are reported in Table 1. The first row of each dataset corresponds to the case where 60% entries of the training set are observed, while the second row corresponds to the case with 80% entries observed. From the table we can see that our proposed method AFASMC can achieve the best performance in terms of both reconstruction and classification. The only exception is on HillValley with 60% observed entries, where AFASMC is outperformed by NNLS on the reconstruction error with tiny margin, but still achieves the best performance on the test accuracy.

4.3 Results on classification performance

In this subsection we examine the performance of active feature acquisition. The feature matrix is initialized with 60% observed entries for each dataset, while the 40% entries are randomly missing. Then active selection is performed iteratively based on the variance criterion. After each query, we perform matrix completion, and then train a linear SVM on the training data. The accuracy of the classifier on the test set is record.

Figure 1 plots the performance curves of compared methods as the number of queried features increases. Note that the performance of EM Inference is unbearably poor on the abalone, letter and HillValley datasets, and its curves are not plotted on these three datasets to avoid the poor visualization of other curves. Also it can be seen that the initial points are different because the methods are employing different matrix completion methods. It can be observed that the proposed approach AFASMC achieves the best performance in most cases. The performance of EM inference is not stable. It achieves decent performance on image and chess, but loses its edge on the others. The QBC and Stability methods perform similarly and are less competitive to AFASMC in most cases. Lastly, as expected, the Random method is not effective compared to the active methods. We also present in Table 2 the AUC results after different percentages of entries queried. It can be observed that the proposed approach outperforms the others in most cases.

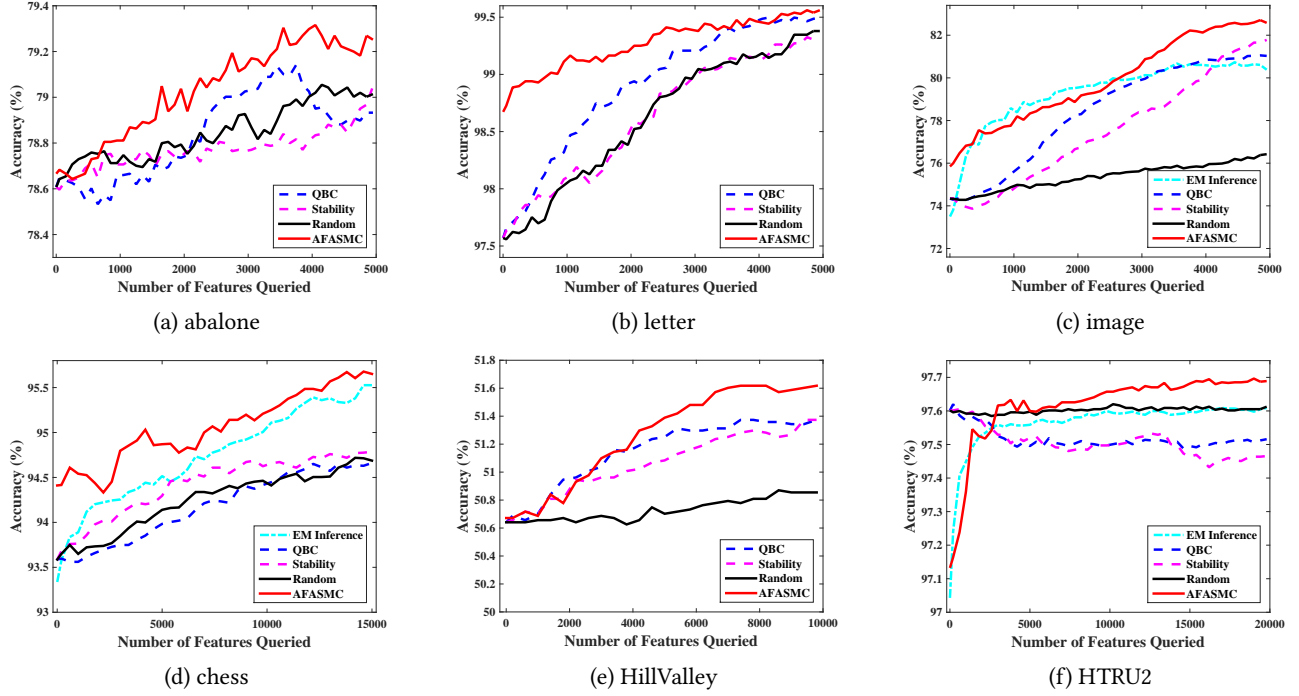


Figure 1: The accuracy curves of compared methods with the number of queried features increasing.

Table 2: The AUC results. The best performances based on paired t -tests at 95% significance level are bold.

Data	Algorithms	Percentage of queried entries						
		5%	10%	20%	30%	40%	50%	80%
abalone	Random	0.859±0.010	0.859±0.010	0.860±0.010	0.861±0.009	0.861±0.009	0.862±0.009	0.866±0.009
	QBC	0.858±0.010	0.858±0.010	0.859±0.009	0.859±0.009	0.861±0.009	0.864±0.009	0.868±0.009
	Stability	0.859±0.010	0.859±0.009	0.859±0.009	0.860±0.009	0.861±0.008	0.862±0.008	0.865±0.009
	AFASMC	0.862±0.009	0.862±0.009	0.863±0.009	0.865±0.009	0.866±0.010	0.867±0.010	0.867±0.010
letter	Random	0.999±0.001	0.999±0.000	0.999±0.000	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000
	QBC	0.999±0.000	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000
	Stability	0.999±0.001	0.999±0.001	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000
	AFASMC	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000
image	Random	0.806±0.012	0.808±0.012	0.810±0.013	0.814±0.014	0.819±0.012	0.824±0.013	0.846±0.013
	QBC	0.804±0.015	0.811±0.014	0.826±0.012	0.838±0.015	0.845±0.015	0.851±0.013	0.858±0.012
	Stability	0.805±0.013	0.807±0.014	0.818±0.015	0.829±0.016	0.835±0.015	0.846±0.016	0.857±0.011
	AFASMC	0.822±0.017	0.832±0.015	0.846±0.012	0.852±0.013	0.855±0.012	0.857±0.012	0.857±0.012
chess	Random	0.983±0.004	0.984±0.004	0.986±0.003	0.987±0.003	0.988±0.002	0.989±0.002	0.991±0.002
	QBC	0.983±0.004	0.984±0.003	0.986±0.004	0.988±0.003	0.990±0.002	0.991±0.002	0.992±0.002
	Stability	0.985±0.004	0.986±0.004	0.987±0.003	0.987±0.003	0.988±0.002	0.990±0.002	0.991±0.002
	AFASMC	0.987±0.004	0.988±0.003	0.989±0.004	0.991±0.003	0.991±0.003	0.992±0.002	0.992±0.002
HillValley	Random	0.452±0.094	0.451±0.094	0.449±0.093	0.446±0.094	0.445±0.090	0.442±0.090	0.437±0.078
	QBC	0.449±0.078	0.443±0.075	0.436±0.075	0.435±0.073	0.434±0.072	0.434±0.072	0.434±0.072
	Stability	0.445±0.085	0.442±0.086	0.439±0.085	0.438±0.070	0.434±0.072	0.434±0.072	0.434±0.072
	AFASMC	0.454±0.082	0.447±0.085	0.446±0.073	0.451±0.077	0.451±0.078	0.450±0.078	0.465±0.042
HTRU2	Random	0.971±0.002	0.971±0.002	0.972±0.002	0.972±0.002	0.972±0.002	0.972±0.002	0.973±0.002
	QBC	0.971±0.002	0.968±0.003	0.969±0.002	0.974±0.002	0.975±0.002	0.975±0.002	0.976±0.002
	Stability	0.970±0.002	0.971±0.002	0.969±0.002	0.969±0.003	0.971±0.003	0.971±0.003	0.975±0.002
	AFASMC	0.972±0.002	0.971±0.002	0.973±0.002	0.975±0.002	0.975±0.002	0.976±0.002	0.976±0.002

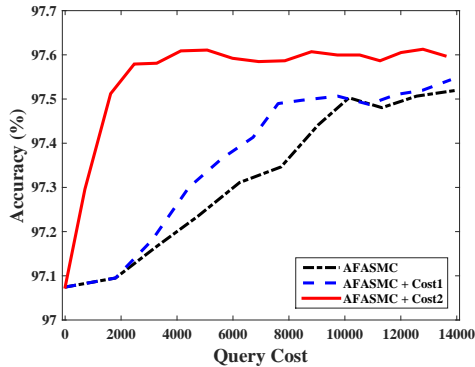


Figure 2: Results of cost-aware selections on HTRU2.

4.4 Study with varied acquisition costs

As discussed previously, the acquisition costs of different features may be diverse. In this subsection, we examine the performance of the proposed strategies for cost-effective feature acquisition. We compare the two optional methods: AFASMC+Cost1, which simply divides the informativeness by the cost; and AFASMC+Cost2, which balances the informativeness and cost via bi-objective optimization. We specify the acquisition cost of each feature dimension as a random integer in $\{1, \dots, 10\}$. Due to space limit, we present the results on the largest dataset HTRU2 as an example.

We record the accuracy after each query, and plot the performance curves in Figure 2. Note that the curve of the original AFASMC is also presented for reference. It can be observed that both the two strategies for considering the acquisition cost can achieve better performance than the original AFASMC. When comparing AFASMC+Cost1 and AFASMC+Cost2, the method with bi-objective optimization achieves a significantly better performance.

4.5 Study on the variance computation

In Section 3.2, when calculating the informativeness based on the variance, we count in all previous iterations of active learning. As discussed before, it is more important to capture the change of an entry within recent iterations. An entry with large variance in the beginning iterations may have been well recovered from recent queries. To examine this idea, we perform experiments to compare the results of calculating the variance with different iterations. Specifically, we use the values of an entry during the last m iterations to calculate the variance, and set m to 2, 4, 8, 16, respectively. Again, for space limit, we report the results on the largest dataset HTRU2 as an example.

The performance curves are plotted in Figure 3. We also plot the curve of counting all iterations as the original AFASMC method. It can be seen that $m = 4$ is the best choice, while counting too few or too many iterations may degrade the performance. This observation is consistent with our conjecture, that the variance computing should emphasize more on recent iterations. Note that we set $m = T$ as default on all datasets to perform the experiments. All the results of AFASMC in previous sections are obtained by counting all iterations. It is thus expected to further improve the

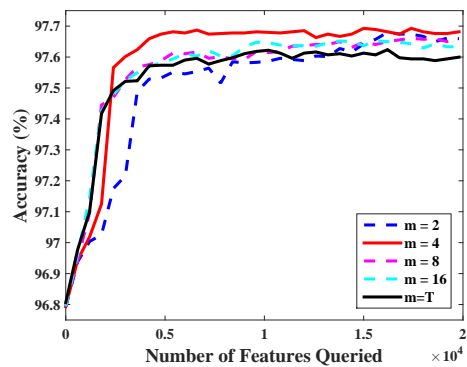


Figure 3: Comparison of variance computation on HTRU2.

performance of the proposed approach by tuning the number of iterations for evaluating the informativeness.

5 CONCLUSION

In this paper, we studied the problem of learning from data with missing features. Since the acquisition of ground-truth feature values is usually expensive, our target was to train an effective classification model with the least acquisition cost. We proposed a unified framework to jointly perform matrix completion and active feature acquisition. On one hand, missing values of the feature matrix are recovered by supervised matrix completion, which exploits the feature correlations with a low-rank regularizer, and the label supervision is utilized by minimizing the empirical classification error. On the other hand, the missing entries are actively queried based on a novel selection criterion, which simultaneously evaluates potential contribution of a feature on both recovering other entries and improving the classification model. Moreover, a bi-objective optimization method was introduced to handle the case where acquisition costs vary for different features. Extensive experimental results validated the superiority of our approach on matrix completion as well as classification performance. In the future, we plan to extend our approach and theoretical analysis to perform active querying both for missing features and class labels.

ACKNOWLEDGMENT

This research was partially supported by National Key R&D Program of China (2018YFB1004300), NSFC (61503182, 61732006), JiangsuSF (BK20150754), the Collaborative Innovation Center of Novel Software Technology and Industrialization, the International Research Center for Neurointelligence (WPI-IRCIN) at The University of Tokyo Institutes for Advanced Study. Authors want to thank Bo-Jian Hou for proofreading.

REFERENCES

- [1] Aniruddha Bhargava, Ravi Ganti, and Rob Nowak. 2017. Active positive semidefinite matrix completion: Algorithms, theory and applications. In *International Conference on Artificial Intelligence and Statistics*. 1349–1357.
- [2] Ricardo Cabral, Fernando De la Torre, Joao Paulo Costeira, and Alexandre Bernardino. 2015. Matrix completion for weakly-supervised multi-label image classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 37, 1 (2015), 121–135.

- [3] Jian-Feng Cai, Emmanuel Candès, and Zuowei Shen. 2010. A singular value thresholding algorithm for matrix completion. *SIAM Journal on Optimization* 20, 4 (2010), 1956–1982.
- [4] Emmanuel Candès and Benjamin Recht. 2009. Exact matrix completion via convex optimization. *Foundations of Computational mathematics* 9, 6 (2009), 717.
- [5] Shayok Chakraborty, Jiayu Zhou, Vineeth Balasubramanian, Sethuraman Panchanathan, Ian Davidson, and Jieping Ye. 2013. Active matrix completion. In *IEEE International Conference on Data Mining*. 81–90.
- [6] Yudong Chen, Srinadh Bhojanapalli, Sujay Sanghavi, and Rachel Ward. 2014. Coherent matrix completion. In *International Conference on Machine Learning*. 674–682.
- [7] David Cohn, Zoubin Ghahramani, and Michael Jordan. 1995. Active learning with statistical models. In *Advances in Neural Information Processing Systems*, Vol. 8. 705–712.
- [8] Mark Davenport, Yaniv Plan, Ewout van den Berg, and Mary Wootters. 2012. 1-bit matrix completion. *arXiv* 1209.3672 (2012).
- [9] Amit Dhurandhar and Karthik Sankaranarayanan. 2015. Improving classification performance through selective instance completion. *Machine Learning* 100, 2-3 (2015), 425–447.
- [10] Brian Eriksson, Laura Balzano, and Robert Nowak. 2011. High-rank matrix completion and subspace clustering with missing data. *arXiv* 1112.5629 (2011).
- [11] Yoav Freund, Sebastian Seung, Eli Shamir, and Naftali Tishby. 1997. Selective sampling using the query by committee algorithm. *Machine Learning* 28, 2-3 (1997), 133–168.
- [12] Andrew Goldberg, Xiaojin Zhu, Ben Recht, Jun-Ming Xu, and Robert Nowak. 2010. Transduction with matrix completion: Three birds with one stone. In *Advances in Neural Information Processing Systems*. 757–765.
- [13] Bo-Jian Hou, Lijun Zhang, and Zhi-Hua Zhou. 2017. Learning with feature evolvable streams. In *Advances In Neural Information Processing Systems*. 1416–1426.
- [14] Sheng-Jun Huang, Miao Xu, Ming-Kun Xie, Masashi Sugiyama, Gang Niu, and Songcan Chen. 2018. Active Feature Acquisition with Supervised Matrix Completion. *arXiv* 1802.05380 (2018).
- [15] Sheng-Jun Huang, Rong Jin, and Zhi-Hua Zhou. 2014. Active learning by querying informative and representative examples. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 10 (2014), 1936–1949.
- [16] Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi. 2013. Low-rank matrix completion using alternating minimization. In *ACM Symposium on Theory of Computing*. 665–674.
- [17] Raghunandan Keshavan, Andrea Montanari, and Sewoong Oh. 2010. Matrix completion from a few entries. *IEEE Transactions on Information Theory* 56, 6 (2010), 2980–2998.
- [18] Franz Király, Louis Theran, and Ryota Tomioka. 2015. The algebraic combinatorial approach for low-rank matrix completion. *Journal of Machine Learning Research* 16 (2015), 1391–1436.
- [19] Akshay Krishnamurthy and Aarti Singh. 2013. Low-rank matrix and tensor completion via adaptive sampling. In *Advances In Neural Information Processing Systems*. 836–844.
- [20] Chee-Peng Lim, Jenn-Hwai Leong, and Mei-Ming Kuan. 2005. A hybrid neural network system for pattern classification tasks with missing features. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 27, 4 (2005), 648–653.
- [21] Huan Liu and Hiroshi Motoda. 1998. *Feature extraction, construction and selection: A data mining perspective*. Vol. 453. Springer Science & Business Media.
- [22] Yong Luo, Tongliang Liu, Dacheng Tao, and Chao Xu. 2015. Multiview matrix completion for multilabel image classification. *IEEE Transaction on Image Processing* 24, 8 (2015), 2355–2368.
- [23] Charalampos Mavroforakis, Dóra Erdős, Mark Crovella, and Evimaria Terzi. 2017. Active positive-definite matrix completion. In *SIAM International Conference on Data Mining*. 264–272.
- [24] Raghu Meka, Prateek Jain, and Inderjit Dhillon. 2009. Matrix completion from power-law distributed samples. In *Advances in Neural Information Processing Systems*. 1258–1266.
- [25] Prem Melville, Foster Provost, and Raymond Mooney. 2005. An expected utility approach to active feature-value acquisition. In *IEEE International Conference on Data Mining*. 745–748.
- [26] Seungwhan Moon, Calvin McCarter, and Yu-Hsin Kuo. 2014. Active learning with partially featured data. In *International Conference on World Wide Web*. 1143–1148.
- [27] Sahand Negahban and Martin Wainwright. 2012. Restricted strong convexity and weighted matrix completion: Optimal bounds with noise. *Journal of Machine Learning Research* 13, May (2012), 1665–1697.
- [28] Chao Qian, Yang Yu, and Zhi-Hua Zhou. 2015. Subset selection by Pareto optimization. In *Advances in Neural Information Processing Systems*. 1774–1782.
- [29] Jason Rennie and Nathan Srebro. 2005. Fast maximum margin matrix factorization for collaborative prediction. In *International Conference on Machine Learning*. 713–719.
- [30] Natali Ruchansky, Mark Crovella, and Evimaria Terzi. 2015. Matrix completion with queries. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 1025–1034.
- [31] Karthik Sankaranarayanan and Amit Dhurandhar. 2013. Intelligently querying incomplete instances for improving classification performance. In *Conference on Information and Knowledge Management*. 2169–2178.
- [32] Burr Settles. 2012. Active learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning* 6, 1 (2012), 1–114.
- [33] Amit Singer and Mihai Cucuringu. 2010. Uniqueness of low-rank matrix completion by rigidity theory. *SIAM Journal on Matrix Analysis Applications* 31, 4 (2010), 1621–1641.
- [34] Leilei Sun, Chonghui Guo, Chuanren Liu, and Hui Xiong. 2017. Fast affinity propagation clustering based on incomplete similarity matrix. *Knowledge and Information System* 51, 3 (2017), 941–963.
- [35] Dougal Sutherland, Barnabás Póczos, and Jeff Schneider. 2013. Active learning and search on low-rank matrices. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 212–220.
- [36] Kim-Chuan Toh and Sangwoon Yun. 2010. An accelerated proximal gradient algorithm for nuclear norm regularized linear least squares problems. *Pacific Journal of Optimization* 6, 615-640 (2010), 15.
- [37] Simon Tong and Daphne Koller. 2001. Support vector machine active learning with applications to text classification. *Journal of Machine Learning Research* 2 (2001), 45–66.
- [38] Paul Tseng. 2008. *On accelerated proximal gradient methods for convex-concave optimization*. Technical Report. University of Washington, Seattle.
- [39] Duy Vu, Mikhail Bilenko, Maytal Saar-Tszechansky, and Prem Melville. 2008. Intelligent information acquisition for improved clustering. *Folia Veterinaria* (2008).
- [40] Kilian Weinberger and Lawrence Saul. 2006. Unsupervised learning of image manifolds by semidefinite programming. *International Journal of Computer Vision* 70, 1 (2006), 77–90.
- [41] Zaiwen Wen, Wotao Yin, and Yin Zhang. 2012. Solving a low-rank factorization model for matrix completion by a nonlinear successive over-relaxation algorithm. *Mathematical Programming Computation* 4, 4 (2012), 333–361.
- [42] Miao Xu, Rong Jin, and Zhi-Hua Zhou. 2013. Speedup matrix completion with side information: Application to multi-label learning. In *Advances In Neural Information Processing Systems*. 2301–2309.
- [43] Jinfeng Yi, Tianbao Yang, Rong Jin, Anil Jain, and Mehrdad Mahdavi. 2012. Robust ensemble clustering by matrix completion. In *IEEE International Conference on Data Mining*. 1176–1181.
- [44] Guangxiang Zeng, Ping Luo, Enhong Chen, Hui Xiong, Hengshu Zhu, and Qi Liu. 2015. Convex matrix completion: A trace-ball optimization perspective. In *SIAM International Conference on Data Mining*. 334–342.