

Optimal Recovery of Missing Values for Non-negative Matrix Factorization

Rebecca Chen and Lav R. Varshney

Abstract

We extend the approximation-theoretic technique of optimal recovery to the setting of imputing missing values in clustered data, specifically for non-negative matrix factorization (NMF), and develop an implementable algorithm. Under certain geometric conditions, we prove tight upper bounds on NMF relative error, which is the first bound of this type for missing values. We also give probabilistic bounds for the same geometric assumptions. Experiments on image data and biological data show that this theoretically-grounded technique performs as well as or better than other imputation techniques that account for local structure.

Index Terms

imputation, missing values, non-negative matrix factorization, optimal recovery

I. INTRODUCTION

MATRIX factorization is commonly used for clustering and dimensionality reduction in computational biology, imaging, and other fields. Non-negative matrix factorization (NMF) is particularly favored by engineers and biologists because non-negativity constraints preclude negative values that are difficult to interpret in biological processes [2], [3]. A recent tutorial article highlighted the interpretability and identifiability (or model uniqueness) of NMF, both of which are valuable for practical applications [4]. Without constraints on the factor matrices, latent factors are non-unique, but requiring non-negativity guarantees model uniqueness under certain assumptions. Furthermore, experimental results demonstrate that the latent factors are intuitive given the data [4]. In signal processing and statistical learning, NMF has been used for speech and audio separation, medical imaging, community detection, and topic modeling. In biology, NMF of gene count matrices can discover cell groups and lower-dimensional manifolds (latent factors) describing gene count ratios for different cell types. Due to channel noise, incomplete survey data, or biological limitations, however, data matrices are usually incomplete and matrix imputation is often necessary before further analysis [3]. In particular, Stein-O'Brien et al. argue that “newer MF algorithms that model missing data are essential for [single-cell RNA sequence] data” [5].

Imputation accuracy is commonly measured using root mean-squared error (RMSE) or similar error metrics. However, Tuikkala et al. argue that “the success of preprocessing methods should ideally be evaluated also in other terms, for example, based on clustering results and their biological interpretation, that are of more practical importance for the biologist” [6]. Here, we specifically consider imputation performance in the context of NMF and describe NMF error rather than RMSE (or some other prediction error) of imputed matrices. Previous analyses of information processing algorithms with missing data have considered high-dimensional regression [7] and subspace clustering [8]. Missing values for NMF have also been studied for the application of stock price prediction, but previous approaches lack theoretical guarantees [9].

Data often exhibits local structure, e.g., different groups of cells follow different gene expression patterns. Information about local structure can be used to improve imputation. We introduce a new imputation method based on *optimal recovery*, an approximation-theoretic approach for estimating linear functionals of a signal [10]–[12] previously applied in signal and image interpolation [13]–[15], to perform matrix imputation of clustered data. Characterizing optimal recovery for missing value imputation requires a new geometric analysis technique. Previous work on missing values take a statistical approach rather than a geometric one.

Our contributions include:

- A computationally efficient imputation algorithm that performs as well as or better than other modern imputation methods, as demonstrated on hyperspectral remote sensing data and biological data;
- A tight upper bound on the relative error of downstream analysis by NMF. This is the first such error bound for settings with missing values; and
- A probabilistic bound on NMF error after imputation.

The remainder of the paper is organized as follows. In Section II, we give background on missing data mechanisms, imputation algorithms, and NMF. In Section III, we introduce optimal recovery and apply it to NMF. In Section IV, we present an algorithm for optimal recovery imputation of clustered data and give a deterministic upper bound on algorithm performance.

This work was presented in part at the 2019 IEEE Data Science Workshop [1] and was supported in part by Air Force STTR Grant FA8650-16-M-1819 and in part by grant number 2018-182794 from the Chan Zuckerberg Initiative DAF, an advised fund of the Silicon Valley Community Foundation.

R. Chen and L. R. Varshney are with the Coordinated Science Laboratory and the Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign, Urbana, IL 61801 USA (e-mail: {rebecca9, varshney}@illinois.edu).

In Section V we give a probabilistic bound on the performance of our algorithm. In Section VI we give experimental results for both synthetic and real data, and we conclude in Section VII.

II. BACKGROUND

In this section, we describe the relationships between missingness patterns and the underlying data, which are referred to as *missing data mechanisms*. We then discuss prior work on practical imputation algorithms, and we present NMF as an analysis technique that is commonly performed after imputation.

A. Missingness mechanisms

Rubin originally described three mechanisms that may account for missing values in data: missing completely at random (MCAR), missing at random (MAR), and missing not at random (MNAR) [16]. When data is MCAR, the missing data is a random subset of all data, and the missing and observed values have similar distributions [17]. The MCAR condition is described in (1), where X refers to observed variables and Y refers to missing variables. This may occur if a researcher forgets to collect certain information for certain subjects, or if certain data samples are collected only for a random subset of test subjects. When data is MAR, the distribution of missing data is dependent on the observed data, (2). For example, in medical records, patients with normal blood pressure levels are more likely to have missing values for glucose levels than patients with high blood pressure. When data is MNAR, the distribution of missing data is dependent on the unobserved (missing) data, (3). For example, people with very high incomes may be less likely to report their incomes.

$$\text{MCAR: } \mathbb{P}(Y \text{ is missing} | X, Y) = \mathbb{P}(Y \text{ is missing}) \quad (1)$$

$$\text{MAR: } \mathbb{P}(Y \text{ is missing} | X, Y) = \mathbb{P}(Y \text{ is missing} | X) \quad (2)$$

$$\text{MNAR: } \mathbb{P}(Y \text{ is missing} | X, Y) = \mathbb{P}(Y \text{ is missing} | Y) \quad (3)$$

It is important to understand the missingness mechanism when analyzing data. When data is MCAR, the statistics of the complete cases (data points with no missing observations) will represent the statistics of the entire dataset, but the sample size will be much smaller. If data is MAR or MNAR, the complete cases may be a biased representation of the dataset. One can also estimate statistics such as means, variances, and covariances based on all available non-missing observations of a variable [18]. Then the sample size reduction may be less severe for certain variables. However, this may be a biased representation of the entire dataset when data is MAR or MNAR, and there is the additional problem of inconsistent sample sizes. Although some research has been done on MNAR imputation, this is generally a difficult problem, and most imputation methods assume the MAR or MCAR model.

Ding and Simonoff argue that missingness mechanisms are more nuanced than the three basic categories described by Rubin [19]. They claim that missingness can be dependent on any combination of the missing values, the observed predictors, and the response variable (e.g. a category label). In cases where the missingness pattern contains information about the response variable, the missingness is *informative* [20]. Ghorbani and Zou leverage informative missingness, using the missingness patterns themselves as an additional feature for data classification [21].

B. Imputation algorithms

Imputation is often necessary before specific downstream analysis, such as clustering or manifold-finding for classification. Two main categories of imputation are single imputation, in which missing values are imputed once, and multiple imputation, in which missing values are imputed multiple times with some built-in randomness. The variance in the multiple imputations of each missing observation reflects the uncertainty of the estimates, and *all* imputed datasets are used in the downstream analysis, which increases statistical power.

1) *Single imputation*: One of the simplest imputation techniques is *mean imputation*. Missing values of each variable are imputed with the mean value of that variable. Since all missing observations of a variable are imputed with the same value, variance is reduced, and other statistics may be skewed in the MAR and MNAR cases. The reduced variability in the imputed variable also decreases correlation with other variables [22].

In *regression imputation*, a variable of interest is regressed on the other variables using the complete cases. Imputation puts points with missing values directly on the regression line. This method also underestimates variances, but it overestimates correlations. *Stochastic regression* attempts to add the variance back by distributing imputed points above and below the regression line using a normal distribution.

Bayesian imputation approaches also exist, including *Bayesian PCA* [23] and *maximum likelihood imputation* [24]. Bayesian methods are theoretically sound and assume that data samples are generated from some underlying joint distribution. In practice, these methods require numerical algorithms such as the Markov chain Monte Carlo (MCMC) method, which may be prohibitively time-consuming for large datasets.

2) *Multiple imputation*: Multiple imputation attempts to preserve the variance/covariance matrix of the data. Several imputations are randomly generated, resulting in multiple complete datasets. Imputed datasets are then analyzed and results are pooled. The different imputations introduce variance into the data, but the variance may still be an underestimate since the imputations assume correlation between the variables. One popular algorithm for multiple imputation is *multiple imputation by chained equations* (MICE) [25]. While MICE does not have the theoretical backing that maximum likelihood imputation has, MICE is flexible and can accommodate known interactions and independencies of real-world datasets [26]. A stepwise regression can be performed so that the missing variable is regressed on the best predictors.

C. Imputation with clustered data

When the underlying data is clustered, a data point should be imputed based on its cluster membership. Local imputation approaches outperform global ones when there is local structure in data. Global approaches generally perform some form of regression or mean matching across all samples [25], [27], whereas local approaches group subsets of similar samples. Popular imputation algorithms that utilize local structure include k-nearest neighbors (kNN), local least squares (LLSimpute), and bicluster Bayesian component analysis (biBPCA) [28]–[30]. The kNN imputation method finds the k closest neighbors of a sample with missing values (measured by some distance function) and fills in the missing values using an average of its neighbors. LLSimpute uses a multiple regression model to impute the missing values from k nearest neighbors. Rather than regressing on *all* variables, biBPCA performs linear regression using biclusters of a lower-dimensional space, i.e. coherent clusters consisting of correlated variables under correlated experimental conditions. Delalleau et al. develop an algorithm to train Gaussian mixtures with missing data using expectation-maximization (EM) [31]. By itself, MICE does not address clusters, but cluster-specific (group-wise) regression can be performed [32].

Tuikkala et al.’s clustering results on cDNA microarray datasets showed that “even when there are marked differences in the measurement-level imputation accuracies across the datasets, these differences become negligible when the methods are evaluated in terms of how well they can reproduce the original gene clusters or their biological interpretations” [6]. They use the Average Distance Between Partition (ADBP) to calculate clustering error, and they show that BPCA, LLS, and kNN give similar clustering results. Chiu et al. find that LLS-like algorithms performed better than kNN-like algorithms in terms of downstream clustering accuracy (measured using cluster pair proportions) [33]. De Souto et al. evaluate whether the effect of different imputation methods on clustering and classification are statistically significant [34]. They remove all genes with more than 10% missing values and compare classification using the corrected Rand index. They find that simple methods such as mean and median imputation perform as well as weighted kNN and BPCA.

D. Non-negative matrix factorization (NMF)

After imputation, downstream analysis such as NMF can be performed on data. Donoho and Stodden interpret NMF as the problem of finding cones in the positive orthant which contain clouds of data points [35]. Liu and Tan show that a rank-one NMF gives a good description of near-separable data and provide an upper bound on the relative reconstruction error [36]. Given that gene and protein expression data is often linearly separable on some manifold- or high-dimensional space [37], the bound given by rank-one NMF is valid. We extend these ideas to data with missing values and, for the first time, bound performance of downstream analysis of imputation. Loh and Wainwright have previously bounded linear regression error of data with missing values [7], but they do not consider imputation, and their proof is based on modifying the covariance matrix when data is missing. Tan and Févotte consider NMF with missing values, but they replace missing with zeros instead of performing imputation, and they do not provide error bounds [9]. In addition, their NMF algorithm requires hyperparameter tuning and makes some probabilistic assumptions on the data. Our proof is based on the geometry of NMF and is parameter-free. The nonnegativity of NMF lends itself to a geometric interpretation, which we describe in Section III.

III. OPTIMAL RECOVERY

In this section, we introduce our new approach to imputation based on approximation-theoretic ideas. Suppose we are given an unknown signal v that lies in some signal class C_k . The optimal recovery estimate \hat{v} minimizes the maximum error between \hat{v} and all signals in the feasible signal class. Given well-clustered non-negative data \mathbf{V} , we impute missing samples in \mathbf{V} so the maximum error is minimized over feasible clusters, regardless of the missingness pattern.

A. Application to clustered data

Let $\mathbf{V} \in \mathbb{R}_+^{F \times N}$ be a matrix of N sample points with F observations (N points in F -dimensional Euclidean space). Suppose the N data points lie in K disjoint clusters C_k (where $k = 1, 2, \dots, K$), and that these clusters are compact, convex spaces (e.g., the convex hull of the points belonging to C_k).

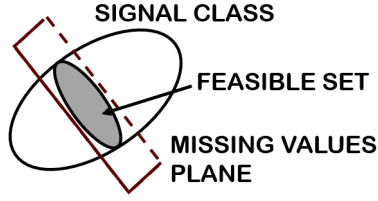


Fig. 1. Feasible set of estimators.

Now suppose there are missing values in \mathbf{V} . Let $\Omega \in \{0, 1\}^{F \times N}$ be a matrix of indicators with $\Omega_{ij} = 1$ if v_{ij} is observed and 0 otherwise. We make no assumptions on the missingness pattern, such as MCAR or MAR because we take a geometric approach rather than a statistical one. We define the projection operator of a matrix \mathbf{Y} onto an index set Ω by

$$[P_{\Omega}(\mathbf{Y})]_{ij} = \begin{cases} \mathbf{Y}_{ij} & \text{if } \Omega_{ij} = 1 \\ 0 & \text{if } \Omega_{ij} = 0 \end{cases}.$$

We use the subscripted vector $(\cdot)_{fo}$ to denote fully observed data points (columns), or data points with no missing values, and we use the subscripted vector $(\cdot)_{po}$ to denote partially observed data points. We use a subscripted matrix $(\cdot)_{fo}$ or $(\cdot)_{po}$ to denote the set of all fully observed or partially observed data columns in the matrix.

We can impute a partially observed vector v_{po} by observing where its observed samples intersect with the clusters C_1, \dots, C_k . Let the *missing values plane* be the restriction set over \mathbb{R}^F that satisfies the constraints on the observed values of v_{po} . We call this intersection the *feasible set* W :

$$W = \{\hat{v}_{po} \in C_k : P_{\Omega}(\hat{v}_{po}) = P_{\Omega}(v_{po})\} \text{ for some } k \in [K]. \quad (4)$$

Fig. 1 illustrates the feasible set of a three-dimensional vector with two missing samples when the signal class (convex space containing samples from k th cluster) covers an ellipsoid. If the signal had only one missing sample, the feasible set would be a line segment.

All k for which (4) is satisfied are possible clusters from which the true v originated. Since W cannot be empty, there must be at least one C_k that has non-empty intersection with the set of all points satisfying the $P_{\Omega}(v_{po})$ constraint. The optimal recovery estimator \hat{v}_{po}^* minimizes the maximum error over the feasible set of estimates:

$$\hat{v}_{po}^* = \arg \min_{\hat{v}_{po} \in C_k} \max_{v \in C_k} \|\hat{v}_{po} - v\|, \quad (5)$$

where $\|\cdot\|$ denotes some norm or error function. If we use the ∞ -norm, \hat{v}_{po}^* is the Chebyshev center of the feasible set.

If W contains estimators belonging to more than one C_k , W can be partitioned into K disjoint sets, W_k , defined as

$$W_k = \{\hat{v}_{po} \in C_k : P_{\Omega}(\hat{v}_{po}) = P_{\Omega}(v_{po})\}, \quad k \in [K]. \quad (6)$$

Feasible clusters are those for which W_k is not empty, and we can find (5) over the C_k for which the corresponding W_k covers the largest volume: $k = \arg \max_k |W_k|$.

B. Application to non-negative matrix factorization

Let $\mathbf{V} \in \mathbb{R}_+^{F \times N}$ be a matrix of N sample points with F non-negative observations. Suppose the columns in \mathbf{V} are generated from K clusters. There exist $\mathbf{W} \in \mathbb{R}_+^{F \times K}$ and $\mathbf{H} \in \mathbb{R}_+^{K \times N}$ such that $\mathbf{V} = \mathbf{W}\mathbf{H}$. This is the NMF of \mathbf{V} [38]. We use the conical interpretation of NMF [35], [36], described as follows.

Suppose the N data points originate from K cones. We define a circular cone $C(u, \alpha)$ by a direction vector u and an angle α :

$$C(u, \alpha) := \left\{ x \in \mathbb{R}^F \setminus \{0\} : \frac{x \cdot u}{\|x\|_2} \geq \cos \alpha \right\}, \quad (7)$$

or equivalently,

$$C(u, \alpha) := \{x \in \mathbb{R}^F \setminus \{0\} : (x \cdot u)^2 - (x \cdot x) \cos^2(\alpha) \geq 0\}. \quad (8)$$

In a three-dimensional space, this conical hull is sometimes called an ice cream cone [4]. We truncate the circular cones to be in the non-negative orthant P so that we have $C(u, \alpha) \cap P$. We can consider u_k to be the dictionary entry corresponding to C_k and all x 's belonging to C_k as noisy versions of u_k . We call the angle between cones $\beta_{ij} := \arccos(u_i \cdot u_j)$. Assume the columns of \mathbf{V} are in K well-separated cones, that is,

$$\min_{i,j \in [K], i \neq j} \beta_{ij} > \max_{i,j \in [K], i \neq j} \{\max\{\alpha_i + 3\alpha_j, 3\alpha_i + \alpha_j\}\}. \quad (9)$$

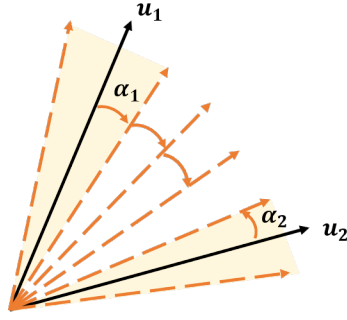


Fig. 2. Geometric assumption for greedy clustering.

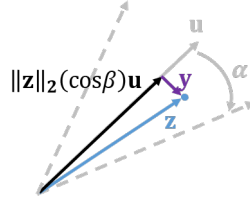


Fig. 3. Decomposition of vectors in a circular cone.

This implies that the distance between any two points originating from the same cluster is less than the distance between any two points in different clusters, which is a common assumption used to guarantee clustering performance [36], [39], [40] (see Fig. 2). We can then partition \mathbf{V} into k sets, denoted $\mathbf{V}_k := \{\mathbf{v}_n \in C_k \cap P\}$, and rewrite \mathbf{V}_k as the sum of a rank-one matrix \mathbf{A}_k (parallel to u_k) and a perturbation matrix \mathbf{E}_k (orthogonal to u_k). For any vector $\mathbf{z} \in \mathbf{V}_k$, $\mathbf{z} = \|\mathbf{z}\|_2(\cos \beta)\mathbf{u}_k + \mathbf{y}$, where $\|\mathbf{y}\|_2 = \|\mathbf{z}\|_2(\sin \beta) \leq \|\mathbf{z}\|_2(\sin \alpha_k)$. We use this rank-one approximation to find error bounds [36] (see Fig. 3).

If \mathbf{V} contains missing values, we can use the optimal recovery estimator to impute \mathbf{V} . Assuming the columns in \mathbf{V} come from K circular cones defined as (7), there is a pair of factor matrices $\mathbf{W}^* \in \mathbb{R}_+^{F \times K}$, $\mathbf{H}^* \in \mathbb{R}_+^{K \times N}$, such that

$$\frac{\|\mathbf{V} - \mathbf{W}^* \mathbf{H}^*\|_F}{\|\mathbf{V}\|_F} \leq \max_{k \in [K]} \{\sin \alpha_k\}. \quad (10)$$

Since the error is bounded by $\sin \alpha_k$, we choose our optimal recovery estimator to minimize α_k . This is equivalent to maximizing the inequality in (8):

$$\hat{v}_{po}^* = \arg \max_{\hat{v}_{po} \in C_k} \{(\hat{v}_{po} \cdot u_k)^2 - (\hat{v}_{po} \cdot \hat{v}_{po}) \cos^2(\alpha_k)\}. \quad (11)$$

We can solve (11) analytically using the Lagrangian with known values of v_{po} as equality constraints. We can also solve (11) numerically using projected gradient descent.

Generally, u_k is not known beforehand, but we can find u_k given W_k . Given an ellipse in \mathbb{R}^3 , we reconstruct its cone by drawing lines from its limit points to the origin. Then it is straightforward to find the center of the cone. (Note that while this volume minimization problem is NP-hard, there are efficient and accurate algorithms when certain assumptions are met, which have been used with NMF [4].) Liu and Tan propose the following optimization problem (in the absence of missing values) over the optimal size angle and basis vector for each cluster [36]. We write the data points in each cluster as $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_M] \in \mathbb{R}_+^{F \times M}$ where $M \in \mathbb{N}_+$:

$$\begin{aligned} & \underset{(0, \pi/2)}{\text{minimize}} && \alpha \\ & \text{subject to} && \mathbf{x}_m^T \mathbf{u} \geq \cos \alpha, \quad m \in [M], \\ & && \mathbf{u} \geq 0, \quad \|\mathbf{u}\|_2 = 1, \quad \alpha \geq 0. \end{aligned} \quad (12)$$

Of course, we also do not know C_k or W_k , so we use a clustering algorithm to find the vectors belonging to each C_k (see Sec. IV).

IV. ALGORITHM AND ERROR BOUND

Now we consider clustering and NMF with missing values. If the geometric assumption (9) holds, a greedy clustering algorithm [36, Alg. 1] returns the correct clustering of fully observed data. Here we show that a greedy algorithm also guarantees correct clustering of partially observed data under certain conditions.

Algorithm 1: Greedy Clustering with Missing Values

Data: Data matrix $\mathbf{V} \in \mathbb{R}_+^{F \times N}$, $K \in \mathbb{N}$, $\Omega \in \{0, 1\}^{F \times N}$

Result: Cone indices $J \in \{0, 1, \dots, K\}^N$; $\alpha \in (0, \pi/2)^K$; $u \in \mathbb{R}_+^{F \times K}$

- 1 Partition columns in \mathbf{V} into subsets \mathbf{V}_{f_o} and \mathbf{V}_{p_o} , where \mathbf{V}_{f_o} contains data columns for which $\sum_i r_{ij} = F$, and \mathbf{V}_{p_o} contains remaining columns.;
 - 2 Normalize \mathbf{V}_{f_o} so that all columns have unit ℓ_2 -norm. Let \mathbf{V}'_{f_o} be the normalized matrix ;
 - 3 Cluster items in \mathbf{V}'_{f_o} using greedy clustering [36, Alg. 1] to obtain cluster indices J and run Alg. 3 on \mathbf{V}'_{f_o} to get u_1, \dots, u_k from W^* . ;
 - 4 **for** $v_{p_o} \in \mathbf{V}_{p_o}$ **do**
 - 5 Let Ω_j correspond to observed entries of v_{p_o} . Find $k = \arg \max_{j \in [K]} \cos^{-1} \left(\frac{P_\Omega(\mathbf{z}_j) \cdot P_\Omega(\mathbf{v})}{\|P_\Omega(\mathbf{z}_j)\| \|P_\Omega(\mathbf{v})\|} \right)$. If this condition is maximized by more than one k , choose one at random. Add the index of v_{p_o} to J_k . ;
 - 6 **end**
 - 7 **for** $k \in [K]$ **do**
 - 8 $\alpha_k = \max_{v_{p_o}} \cos^{-1} \left(\frac{P_\Omega(v_{p_o}) \cdot P_\Omega(u_k)}{\|P_\Omega(v_{p_o})\| \|P_\Omega(u_k)\|} \right)$;
 - 9 **end**
 - 10 Return cone indices J , u , α ;
-

Lemma 1 (Greedy clustering with missing values). *Let Ω indicate the missing values of v_{p_o} . Let α_k be the defining angle of C_k and $P_\Omega(\alpha_k)$ be the defining angle of the cone resulting from projecting C_k onto the missing value plane from Ω . If, for exactly one k ,*

$$\arccos \left(\frac{P_\Omega(v_{p_o}) \cdot P_\Omega(u_k)}{\|P_\Omega(v_{p_o})\| \|P_\Omega(u_k)\|} \right) \leq P_\Omega(\alpha_k) \quad (13)$$

then v_{p_o} originated from the corresponding C_k . If α_k are identical for all k , Alg. 1 will cluster v_{p_o} correctly.

Proof. The result follows directly. □

Now consider feasibility of imputing data points using the $\hat{\alpha}$ and \hat{u} from Alg. 2. Clearly, the missing values plane for each point intersects the original corresponding cone defined by the true u and α of the cone. We know the \hat{u} fall somewhere within the original cones, but if the $\hat{\alpha}$ are too small, the new cones may not intersect with the missing values plane.

Lemma 2 (Feasibility of imputation algorithm). *The estimator in (5) is able to find an imputation within the feasible set given $\alpha_1, \dots, \alpha_K$ and u_1, \dots, u_k returned by Alg. 1.*

Proof. Let vector v_{p_o} be a partially observed version of $v_{f_o} \in \mathbf{V}$. We define the angle between v_{p_o} and cluster center u_k in the F -dimensional space:

$$\gamma_k = \arccos \left(\frac{P_\Omega(v_{p_o}) \cdot u_k}{\|P_\Omega(v_{p_o})\| \|u_k\|} \right), \quad (14)$$

and between v_{p_o} and the projected cluster center in the projected $(F - f)$ -dimensional space:

$$\hat{\gamma}_k = \arccos \left(\frac{P_\Omega(v_{p_o}) \cdot P_\Omega(u_k)}{\|P_\Omega(v_{p_o})\| \|P_\Omega(u_k)\|} \right), \quad (15)$$

where Ω is the observed values indicator corresponding to v_{p_o} . Then $\gamma_k \leq \hat{\gamma}_k$ since $P_\Omega(v_{p_o}) \cdot u_k = P_\Omega(v_{p_o}) \cdot P_\Omega(u_k)$ and $\|u_k\| \geq \|P_\Omega(u_k)\|$. Thus $\hat{\gamma}_k$ is large enough that an imputation on the missing values plane is feasible for each v_{p_o} . Since $\alpha_k = \max \gamma_k$, all partially observed points labeled as belonging to C_k can be imputed. □

Algorithm 2: Rank-1 NMF with Missing Values

Data: Partially observed data $\mathbf{V} \in \mathbb{R}_+^{F \times N}$, $\Omega \in \{0, 1\}^{F \times N}$, $K \in \mathbb{N}$

Result: $\hat{\mathbf{W}}^* \in \mathbb{R}_+^{F \times K}$ and $\hat{\mathbf{H}}^* \in \mathbb{R}_+^{K \times N}$

- 1 Cluster data using Alg. 1 ;
 - 2 Impute data using (5) ;
 - 3 Perform rank-1 NMF on imputed data using [36, Alg. 2] ;
-

We extend bound (10) on the relative NMF error to missing values (Alg. 3). Note that the original bound allows for overlapping cones and does not assume (9) holds. It only requires all points be within α_k of u_k , which essentially allows the normalized perturbation matrix \mathbf{E}_k to be upper-bounded by $\sin \alpha_k$. If the missing entries of each v_{p_o} are imputed using Alg. 1,

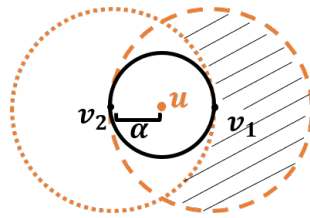


Fig. 4. Geometric proof of relative NMF error bound.

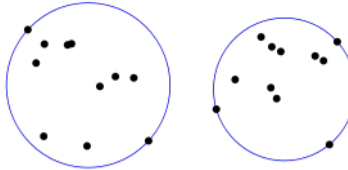


Fig. 5. Minimum covering sphere in two dimensions.

then the perturbation from the original u_k , which we denote $\hat{\mathbf{E}}_k$, will be at most $2\mathbf{E}_k$. We can prove this using a worst-case scenario.

Theorem 1 (Rank-1 NMF with missing values). *Suppose \mathbf{V} is drawn from K cones and missing values are introduced to get \mathbf{V}_{po} . If Alg. 2 correctly clusters data points and Alg. 1 is used to perform imputation, then*

$$\frac{\|\mathbf{V} - \mathbf{W}_{po}^* \mathbf{H}_{po}^*\|_F}{\|\mathbf{V}\|_F} \leq \max_{k \in [K]} \{\sin 2\alpha_k\}, \quad (16)$$

where \mathbf{W}_{po}^* and \mathbf{H}_{po}^* are found by Alg. 3.

Proof. Suppose there are two points v_1 and v_2 in a cone, as indicated by the solid circle in Fig. 4. Then u will be at an angle α from both v_1 and v_2 . Now suppose v_2 contains missing values. Then the new v_1 will be the only vector in the cone, \hat{v}_2 is imputed using (11), where $\hat{u} = v_1$, and \hat{v}_2 is at an angular distance $\sin 2\alpha$ from \hat{u} . (One can check that if there are more than two points in the cone, this distance cannot increase.) A worst-case imputation places \hat{v}_2 at an angle 2α away from v_1 (suppose the optimizer places \hat{v}_2 at an angle greater than 2α from v_1 , but this is a contradiction since then v_2 would be a better estimate than the optimum). The dashed circle in Fig. 4 represents points at an angle 2α from v_1 . Any \hat{v}_2 outside the dotted circle is at an angle greater than 2α from v_2 . So the shaded region indicates when the error may be greater than $\sin 2\alpha$. But the missing values of v_2 allow for “movement” only along the axes. Since the intersection of a hyperplane with a cone is a finite-dimensional ellipsoid [41], [42], which is compact [43], v_2 cannot “travel” via imputation to the shaded region without crossing a feasible region less than 2α from \hat{u} . Hence the theorem holds and is tight. \square

V. PROBABILISTIC ERROR

We now make some probabilistic assumptions on our data and missingness patterns to calculate the expected maximum error of optimal recovery imputation. First, consider a cone C in an F -dimensional space defined by u and α . Let us ignore the length of the vectors in C and preserve only the angles of the vectors from u . We can then represent vectors of an F -dimensional cone as points in an $(F - 1)$ -dimensional ball. For example, a 3-dimensional cone can be represented as points in a circle, as in Fig. 4.

Let there be N points $\{x_1, \dots, x_N\} \in \mathbb{R}^F$, drawn uniformly at random from K F -dimensional balls, labeled B_1, \dots, B_K . Let $d(x_i, x_j)$ be the Euclidean distance between x_i and x_j . We assume there is at least one data point in each ball, and that the distance between any two points in a ball B_k is less than the distance between any point in B_k and a point not in B_k . That is, for any $i, j \in [N], i \neq j$,

$$\max_{i, j \in B_k} d(x_i, x_j) < \min_{i \in B_k, j \notin B_k} d(x_i, x_j) \quad \text{for all } k = 1, \dots, K. \quad (17)$$

This is equivalent to the geometric assumption in (9), and we can correctly cluster any points drawn from such balls using Alg. 1. After obtaining the clusters, we can compute the minimum covering sphere (MCS) on the points in each cluster [44] (Fig. 5). This gives us K balls with N_k points in each ball.

Now suppose that we have partially observed entries in our data. Let the missingness of a point be a Bernoulli random variable with parameter γ . That is, x is fully observed with probability γ and partially observed with probability $1 - \gamma$. There is now some uncertainty about the position of partially observed data points, so we will find the MCS for only the fully observed

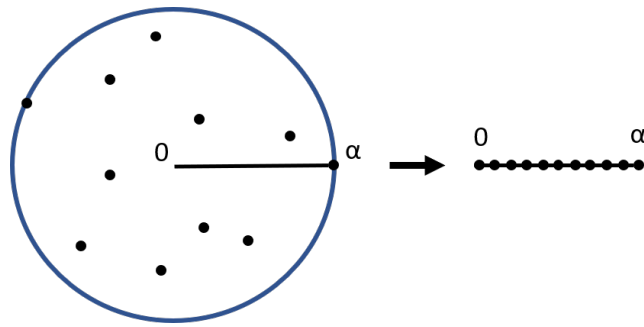


Fig. 6. Assumption that points are uniformly random on the radius.

points. This is analogous to step 3 in Algorithm 2. By calculating the expected change in the radius of the MCS, we can calculate the expected change in its corresponding cone.

Theorem 2 (Probabilistic bound on NMF error). *Given the setting described above, and assuming that the N points are drawn uniformly at random from the K balls, then after imputing with Alg. 1, we can tighten the bound in (16) to*

$$\frac{\|\mathbf{V} - \mathbf{W}_{po}^* \mathbf{H}_{po}^*\|_F}{\|\mathbf{V}\|_F} \leq \max_{k \in [K]} \{\sin \alpha_k\}. \quad (18)$$

Proof. If the N points are drawn uniformly at random from the K balls, then $\mathbb{E}[N_k] = N/k$, and the expected number of fully observed and partially observed points in each cluster is

$$\mathbb{E}[|X_{k,fo}|] = \gamma N_k \quad \text{and} \quad \mathbb{E}[|X_{k,po}|] = (1 - \gamma) N_k. \quad (19)$$

Clearly, the volume of the MCS can only decrease as $|X_{k,fo}|$ decreases. Let R_{max} be the radius of MCS if there were no missing values, and let \hat{R} be the radius of the MCS of only the fully observed points. Then $\hat{R} < R_{max}$ only if any $x \in X_{po}$ originally lay on the surface of $MCS_{k,fo}$. Suppose the points are randomly distributed along the radius of the F -ball and we pick points to be partially observed uniformly at random. Let

$$N_{po} = \lceil (1 - \gamma)N \rceil. \quad (20)$$

Assume x_i are i.i.d. and uniformly distributed (without loss of generality) on $[0, 1]$. This matches the assumption in the probabilistic analysis in [36] that the angles are drawn uniformly at random on $[0, \alpha]$ (see Fig. 6). Assuming a continuous distribution, almost surely no two points have exactly the same radius, and the probability of picking the ℓ outermost points is

$$\mathbb{P}(\ell) = \frac{\binom{N-\ell}{N_{po}-\ell}}{\binom{N}{N_{po}}}, \quad \text{where } \ell = 0, 1, \dots, N_{po}. \quad (21)$$

This gives us

$$\mathbb{E}[\ell] = \sum_{\ell=1}^{N_{po}} \ell \cdot \mathbb{P}[\ell] \quad (22)$$

$$= \sum_{\ell=1}^{N_{po}} \ell \cdot \frac{\binom{N-\ell}{N_{po}-\ell}}{\binom{N}{N_{po}}} \quad (23)$$

$$= \frac{1}{\binom{N}{N_{po}}} \sum_{\ell=1}^{N_{po}} \ell \cdot \binom{N-\ell}{N_{po}-\ell} \quad (24)$$

$$= \frac{\binom{N-1}{N_{po}-1} N(N+1)}{\binom{N}{N_{po}} (N - N_{po} + 1)(N - N_{po} + 2)}, \quad (25)$$

where N_{po} is dependent on γ , as defined in (20).

The radius of the resulting MCS is dependent on the distribution of points along the radius. We can determine \hat{R} using order statistics. If we assume uniform distribution between 0 and 1, and order the points x_1, \dots, x_n so that x_1 is closest to the center of the sphere and x_n is farthest, the radius of the n th point, R_n , is given by the beta distribution

$$R_n \sim B(n, 1), \quad (26)$$

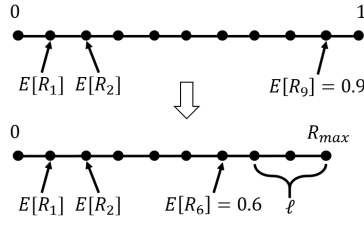


Fig. 7. Example of $\mathbb{E}[\hat{R}]$ with $N = 9$ and $\ell = 3$.

and

$$\mathbb{E}[R_n] = \frac{n}{n+1}. \quad (27)$$

Thus if ℓ of the outermost points are chosen to be missing,

$$\mathbb{E}[\hat{R}] = R_{max} - (\ell/N)R_{max} = \left(\frac{N-\ell}{N}\right)R_{max}. \quad (28)$$

We illustrate with an example in Fig. 7. We can substitute $\mathbb{E}[\ell]$ for ℓ , and since $\mathbb{E}[\ell]$ is a function of γ , we have derived the expected radius of the MCS as a function of missingness:

$$\mathbb{E}[\hat{R}] = \left(\frac{N-\mathbb{E}[\ell]}{N}\right)R_{max}. \quad (29)$$

Now we reverse the arrow in Fig. 6. Due to the random distribution of points in the sphere, removing the ℓ outermost points does not change the expected center u of the MCS. Transitioning from spheres back to cones, we get

$$\mathbb{E}[\hat{\alpha}] = \left(\frac{N-\mathbb{E}[\ell]}{N}\right)\alpha. \quad (30)$$

Thus

$$\alpha - \mathbb{E}[\hat{\alpha}] = \frac{\mathbb{E}[\ell]}{N} \cdot \alpha, \quad (31)$$

and the normalized Frobenius distance between $\mathbf{W}_{fo}^* \mathbf{H}_{fo}^*$ and $\mathbf{W}^* \mathbf{H}^*$ for a single cone is:

$$\frac{\|\mathbf{W}_{fo}^* \mathbf{H}_{fo}^* - \mathbf{W}^* \mathbf{H}^*\|_F}{\|\mathbf{W}^* \mathbf{H}^*\|_F} \leq \sin\left(\frac{\mathbb{E}[\ell]}{N} \cdot \alpha\right). \quad (32)$$

If we assume $v_n \in \mathbf{V}$ are MCAR, the statistical mean of \mathbf{V}_{fo} is the same as that of \mathbf{V} . Since v_n are uniformly distributed, the range of v_n remains centered on the mean, so the expected center of the MCS does not change. Thus the maximum difference between a point $v \in C_k$ and its imputed point \hat{v} is $\sin \alpha_k$, and the theorem follows. \square

A. MCS with a different assumption

If instead we assume points are uniformly distributed in the volume of the ball, we find the change in radius as follows. First, calculate the volume of a F -dimensional ball of radius $R = 1$:

$$V_F(R) = \frac{\pi^{F/2}}{\Gamma(F/2+1)} R^F. \quad (33)$$

Then we calculate radius \hat{R} of an F -dimensional ball as:

$$\hat{R}_F(\hat{V}) = \frac{\Gamma(F/2+1)^{1/F}}{\sqrt{\pi}} \hat{V}^{1/F}, \quad (34)$$

where volume $\hat{V} = \left(\frac{1-\ell}{N}\right) V_F(1)$.

The probability that a point x is in MCS_{po} is

$$\mathbb{P}(x \in \text{MCS}_{po}) = \frac{V(\hat{R})}{V(R_{max})}. \quad (35)$$

Thus the expected radius given a missing parameter γ is given by

$$\mathbb{E}[\hat{R}] = \hat{R}_F\left(\frac{1-\mathbb{E}[\ell]}{N} V_F(1)\right), \quad (36)$$

where $\mathbb{E}[\ell]$ is a function of γ , and the expected NMF error is

$$\mathbb{E} \left[\frac{\|\mathbf{V} - \mathbf{W}^* \mathbf{H}^*\|_F}{\|\mathbf{W}^* \mathbf{H}^*\|_F} \right] = \sin \left(\mathbb{E}[\hat{R}] \cdot \alpha \right). \quad (37)$$

B. Minimum covering spherical cap for normalized data

If the data is normalized such that each vector has an L_2 norm of 1, all the points will fall on the surface of a sphere. Let there be N points $\{x_1, \dots, x_N\} \in \mathbb{R}^F$, drawn at random from K F -dimensional spherical caps of a radius R F -ball, labeled C_1, \dots, C_K . Let $d(x_i, x_j)$ be some distance between x_i and x_j . Assume there is at least one data point in each spherical cap, and that (9) holds.

The area of an F -dimensional spherical cap is

$$A(R, h) = \frac{1}{2} A_F R^{F-1} I_{2rh-h^2/r^2} \left(\frac{F-1}{2}, \frac{1}{2} \right), \quad (38)$$

where $0 \leq h \leq R$, $A_n = 2\pi^{n/2}/\Gamma[n/2]$ is the area of the unit n -ball, h is the height of the cap, which can be calculated as a function of the angle α between the center and the edge of the cap, and $I_x(a, b)$ is the regularized incomplete beta function. Using the same style of analysis from the previous section, we can find the expected angle $\mathbb{E}[\alpha^{po}]$ given a parameter γ for partially observed points. Thus,

$$\mathbb{E} \left[\frac{\|\mathbf{V} - \mathbf{W}^* \mathbf{H}^*\|_F}{\|\mathbf{W}^* \mathbf{H}^*\|_F} \right] = \sin \left(\mathbb{E}[\alpha^{po}] \right). \quad (39)$$

VI. EXPERIMENTAL RESULTS

To test our algorithm, we first generate conical data satisfying the geometric assumption, using $N = 10000$, $F = 160$, and $K = 40$. We choose squared length of each v as a Poisson random variable with parameter 1, and we choose the angles of v uniformly. We then let \mathbf{V} be partially-observed with Bernoulli parameter ξ to obtain \mathbf{V}_{po} . That is,

$$\Omega(i, j) \stackrel{i.i.d.}{\sim} \text{Bern}(\xi). \quad (40)$$

We run tests using $\xi \in \{0.4, 0.55, 0.7, 0.8, 0.9\}$ and find imputation relative error for NMF:

$$E[\mathbf{V}, \mathbf{W}_{po}^* \mathbf{H}_{po}^*] = \frac{\|\mathbf{V} - \mathbf{W}_{po}^* \mathbf{H}_{po}^*\|_F}{\|\mathbf{V}\|_F}. \quad (41)$$

Fig. 8 shows relative error of our optimal recovery imputation with different values of α when we enforce correct clustering. The error for all α values and missingness percentages lies within the bound given by (16). Note that because our data is drawn uniformly at random, the error does not approach the worst-case bound.

In the next experiment, we impute the conical data with $\alpha = 0.1$ with other local imputation algorithms, including kNNimpute [45] with Euclidean, cosine, and Chebyshev (L_∞) distances and iterated local least squares (itrLLS) [46]. We perform two tests with optimal recovery: one with enforced correct clusterings and one without prior knowledge of the correct clusterings. We use $\alpha = 0.1$ and do not enforce correct clustering for Alg. 3 as before (see Fig. 9). We find $k = 8$ neighbors gives us the best results. Optimal recovery performs much better than other methods when clusters are known, and it performs similarly to other methods when they are not.

Following [36], the next experiment tests a subset of the hyperspectral imaging data set from Pavia [47]. We crop the 103 images to have 2000 pixels per image, set $K = 9$, corresponding to the different imagery categories, and introduce missing values in the same proportions as before (see Fig. 10). We also run tests with mice protein data [48] (see Fig. 11). The original dataset contains 1077 measurements with 77 proteins. We remove the 9 proteins that had missing measurements, then introduce missing values. We find $k = 5$ neighbors gives us the best results for kNNimpute on these datasets. On the mouse data, we also test bicluster BPCA [30] in addition to the other methods. The conical and Pavia test data were not sufficiently well-conditioned to run bicluster BPCA. See Tab. I for a comparison of run times. Our results demonstrate that optimal recovery performs similarly to kNN methods when clusters are not known beforehand. When clusters are known, optimal recovery performs similarly to more advanced methods (itrLLSimpute and biBPCA) in a fraction of the time.

VII. CONCLUSION

We have extended classical approximation-theoretic *optimal recovery* to the setting of imputing missing values, specifically for NMF. Fu et al. remark that “separability-based methods have many attractive features, such as identifiability, solvability, provable noise robustness, and the existence of lightweight greedy algorithms” [4]. We showed that imputation using optimal recovery minimizes relative NMF error under certain separability assumptions, and provided a straightforward algorithm for implementation.

Future work aims to extend optimal recovery to other settings of missing values in modern data science. Different applications (and therefore models) require different NMF algorithms and model assumptions. Thus different error bounds will be necessary

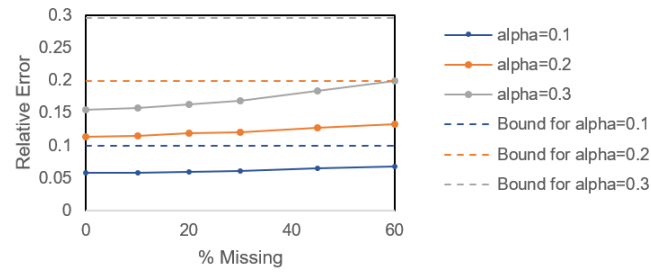


Fig. 8. Relative NMF error of imputed conical data with correct clustering.

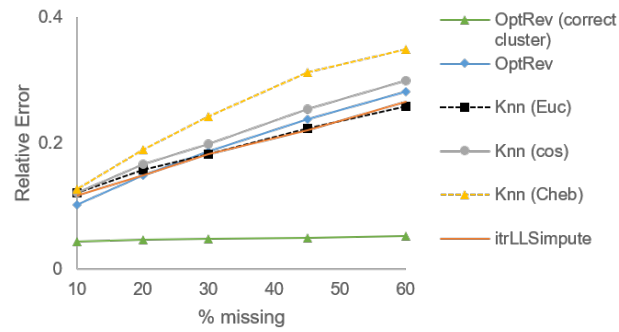


Fig. 9. Relative NMF error for Conical data.

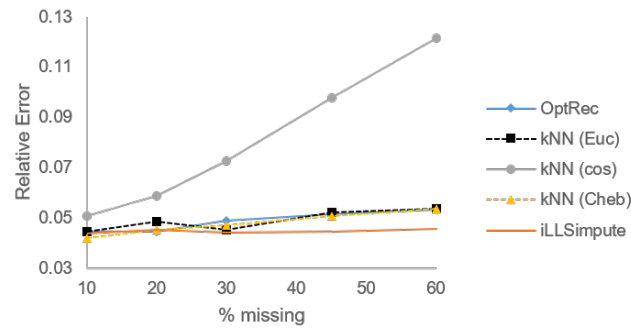


Fig. 10. Relative NMF error for Pavia data.

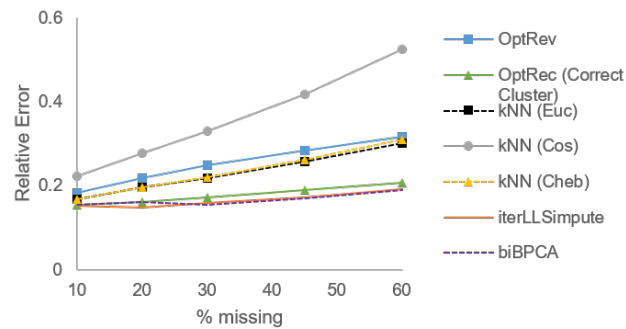


Fig. 11. Relative NMF error for Mouse data.

TABLE I
AVERAGE IMPUTATION TIMES FOR MOUSE DATA IN SECONDS.

% Missing	10	20	30	45	60
OptRec	0.51	0.48	0.63	0.91	1.41
OptRec (Correct Clusters)	0.48	0.48	0.58	0.85	1.36
kNN (Euc)	0.11	0.17	0.29	0.34	0.51
kNN (Cos)	0.10	0.15	0.19	0.27	0.41
kNN (Cheb)	0.08	0.16	0.23	0.35	0.52
itrLLSimpute	43.9	30.4	25.7	20.5	14.5
biBPCA	5000+				

[4]. On the experimental side, we plan to test our imputation algorithm on single-cell RNA sequencing data along with different clustering algorithms; these experiments will inform more specific heuristic refinements. We also aim to extend our algorithm to the scenario where complete cases for each cluster are not available.

REFERENCES

- [1] R. Chen and L. R. Varshney, "Non-negative matrix factorization of clustered data with missing values," in *Proc. IEEE Data Sci. Workshop*, Jun. 2019.
- [2] Q. Qi, Y. Zhao, M. Li, and R. Simon, "Non-negative matrix factorization of gene expression profiles: a plug-in for BRB-ArrayTools," *Bioinformatics*, vol. 25, no. 4, pp. 545–547, Feb. 2009.
- [3] Y. Li and A. Ngom, "The non-negative matrix factorization toolbox for biological data mining," *Source Code for Biology and Medicine*, vol. 8, no. 10, Sep. 2013.
- [4] X. Fu, K. Huang, N. D. Sidiropoulos, and W.-K. Ma, "Nonnegative matrix factorization for signal and data analytics: Identifiability, algorithms, and applications," *IEEE Signal Process. Mag.*, vol. 36, no. 2, pp. 59–80, March 2019.
- [5] G. L. Stein-O'Brien, R. Arora, A. C. Culhane, A. V. Favorov, L. X. Garmire, C. S. Greene, L. A. Goff, Y. Li, A. Ngom, M. F. Ochs, Y. Xu, and E. J. Fertig, "Enter the matrix: factorization uncovers knowledge from omics," *Trends in Genetics*, vol. 34, no. 10, pp. 790–805, Oct. 2018.
- [6] J. Tuikkala *et al.*, "Missing value imputation improves clustering and interpretation of gene expression microarray data," *BMC Bioinformatics*, vol. 9, no. 202, Apr. 2008.
- [7] P.-L. Loh and M. J. Wainwright, "Corrupted and missing predictors: Minimax bounds for high-dimensional linear regression," in *Proc. 2012 IEEE Int. Symp. Inf. Theory (ISIT)*, Jul. 2012.
- [8] Z. Charles, A. Jalali, and R. Willett, "Subspace clustering with missing and corrupted data," *arXiv:1707.02461 [stat.ML]*, Jan. 2018.
- [9] V. Tan and C. Févotte, "Automatic relevance determination in nonnegative matrix factorization with the β -divergence," *IEEE Trans. Pattern Analysis Machine Intelligence*, vol. 35, no. 7, pp. 1592–1602, 2013.
- [10] M. Golomb and H. F. Weinberger, "Optimal approximation and error bounds," in *On Numerical Approximation*, R. E. Langer, Ed. Madison: University of Wisconsin Press, 1959, pp. 117–190.
- [11] C. A. Micchelli and T. J. Rivlin, "A survey of optimal recovery," in *Optimal Estimation in Approximation Theory*, C. A. Micchelli and T. J. Rivlin, Eds. New York: Plenum Press, 1976, pp. 1–54.
- [12] —, "Lectures on optimal recovery," in *Numerical Analysis Lancaster 1984*, ser. Lecture Notes in Mathematics, P. R. Turner, Ed. Berlin: Springer-Verlag, 1985, vol. 1129, pp. 21–93.
- [13] R. G. Shenoy and T. W. Parks, "An optimal recovery approach to interpolation," *IEEE Journal of Signal Processing*, vol. 40, no. 8, pp. 1987–1996, Aug. 1992.
- [14] D. L. Donoho, "Statistical estimation and optimal recovery," *The Annals of Statistics*, vol. 22, no. 1, pp. 238–270, Mar. 1994.
- [15] D. D. Muresan and T. W. Parks, "Adaptively quadratic (AQua) image interpolation," *IEEE Trans. Image Process.*, vol. 13, no. 5, pp. 690–698, May 2004.
- [16] D. B. Rubin, "Inference and missing data," *Biometrika*, vol. 63, no. 3, pp. 581–592, 1976.
- [17] K. Bhaskaran and L. Smeeth, "What is the difference between missing completely at random and missing at random?" *International Journal of Epidemiology*, vol. 43, no. 4, pp. 1336–1339, 2014.
- [18] M. Kolar and E. P. Xing, "Estimating sparse precision matrices from data with missing values," in *Proc. 29th Int. Conf. Machine Learning (ICML '12)*, Jun. 2012, pp. 551–558.
- [19] Y. Ding and J. S. Simonoff, "An investigation of missing data methods for classification trees applied to binary response data," *J. Machine Learning Research*, vol. 11, pp. 131–170, Jan. 2010.
- [20] P. J. García-Laencina and J.-L. Sancho-Gómez and A. R. Figueiras-Vidal, "Pattern classification with missing data: a review," *Neural Computing and Applications*, vol. 19, no. 2, pp. 263–282, 2010.
- [21] A. Ghorbani and J. Y. Zou, "Embedding for informative missingness: Deep learning with incomplete data," in *Proc. 2018 56th Ann. Allerton Conf. Commun., Control, and Comput.*, Oct. 2018, pp. 437–445.
- [22] C. K. Enders, *Applied Missing Data Analysis*. The Guilford Press, 2010.
- [23] S. Oba *et al.*, "A Bayesian missing value estimation method for gene expression profile data," *Bioinformatics*, vol. 19, no. 6, pp. 2088–2096, Nov. 2003.
- [24] K. Messer and L. Natarajan, "Maximum likelihood, multiple imputation and regression calibration for measurement error adjustment," *Statistics in Medicine*, vol. 27, no. 30, pp. 6332–6350, Dec. 2008.
- [25] S. van Buuren and K. Groothuis-Oudshoorn, "Mice: Multivariate imputation by chained equations in R," *Journal of Statistical Software*, vol. 45, no. 3, Dec. 2011.
- [26] M. J. Azur *et al.*, "Multiple imputation by chained equations: What is it and how does it work?" *International Journal of Methods in Psychiatric Research*, vol. 20, no. 1, pp. 40–49, 2011.
- [27] R. J. A. Little and D. B. Rubin, *Statistical Analysis with Missing Data*. Wiley, 2002.
- [28] T. Hastie, R. Tibshirani, G. Sherlock, M. Eisen, P. Brown, and D. Botstein, "Imputing missing data for gene expression arrays," Division of Biostatistics, Stanford University, Technical Report, Oct. 1999.
- [29] H. Kim, G. Golub, and H. Park, "Missing value estimation for DNA microarray gene expression data: local least squares imputation," *Bioinformatics*, vol. 21, no. 2, pp. 187–198, Jan. 2005.
- [30] F. Meng, C. Cai, and H. Yan, "A bicluster-based Bayesian principal component analysis method for microarray missing value estimation," *IEEE Journal of Biomedical and Health Informatics*, vol. 18, no. 3, pp. 863–871, May 2014.
- [31] A. C. Olivier Delalleau and Y. Bengio, "Efficient EM training of Gaussian mixtures with missing data," *arXiv:1209.0521 [cs.LG]*, Jan. 2018.

- [32] A. Robitzsch, S. Grund, and T. Henke, "Miceadds: Some additional multiple imputation functions, especially for mice," 2018, R package version 3.0-16. [Online]. Available: <https://cran.r-project.org/web/packages/miceadds/index.html>
- [33] C.-C. Chiu *et al.*, "Missing value imputation for microarray data: a comprehensive comparison study and a web tool," *BMC Systems Biology*, vol. 7, no. Suppl 6:S12, Dec. 2013.
- [34] M. C. de Souto, P. A. Jaskowiak, and I. G. Costa, "Impact of missing data imputation methods on gene expression clustering and classification," *BMC Bioinformatics*, vol. 16, no. 64, Feb. 2015.
- [35] D. Donoho and V. Stodden, "When does non-negative matrix factorization give a correct decomposition into parts?" in *Advances in Neural Information Processing Systems 16*, S. Thrun, L. K. Saul, and B. Schölkopf, Eds. MIT Press, 2004, pp. 1141–1148.
- [36] Z. Liu and V. Y. F. Tan, "Rank-one NMF-based initialization for NMF and relative error bounds under a geometric assumption," *IEEE Trans. Sign. Process.*, vol. 65, no. 18, pp. 4717–4731, Sep. 2017.
- [37] R. Clarke, H. W. Resson, A. Wang, J. Xuan, M. C. Liu, E. A. Gehan, and Y. Wang, "The properties of high-dimensional data spaces: implications for exploring gene and protein expression data," *Nature Reviews Cancer*, vol. 8, no. 1, pp. 37–49, Jan. 2008.
- [38] D. D. Lee and H. S. Seung, "Learning the parts of objects by non-negative matrix factorization," *Nature*, vol. 401, no. 18, pp. 788–791, Oct. 1999.
- [39] Y. Bu, S. Zou, and V. V. Veeravalli, "Linear-complexity exponentially-consistent tests for universal outlying sequence detection," in *2017 IEEE Int. Symp. Inf. Theory (ISIT)*, Jun. 2017, pp. 988–992.
- [40] A. Y. Ng, M. I. Jordan, and Y. Weiss, "On spectral clustering: Analysis and an algorithm," in *Proc. 14th Int. Conf. Neur. Informat. Process. Syst. (NIPS)*. MIT Press, 2001, pp. 849–856.
- [41] T. L. Heath, *Apollonius of Perga: Treatise on Conic Sections (Edited in Modern Notation)*. Cambridge University Press, 1986.
- [42] M. S. Handlin, "Conic sections beyond \mathbb{R}^2 ," May 2013, notes.
- [43] Y. N. Kiselev, "Approximation of convex compact sets by ellipsoids. Ellipsoids of best approximation," *Proc. Steklov Institute of Mathematics*, vol. 262, no. 1, pp. 96–120, Sep. 2008.
- [44] T. H. Hopp and C. P. Reeve, "An algorithm for computing the minimum covering sphere in any dimension," National Institute of Standards and Technology, Gaithersburg, Maryland, NISTIR 5831, May 1996.
- [45] A. W.-C. Liew, N.-F. Law, and H. Yan, "Missing value imputation for gene expression data: computational techniques to recover missing data from available information." *Briefings in Bioinformatics*, vol. 12, no. 5, pp. 498–513, Sep. 2011.
- [46] Z. Cai, M. Heydari, and G. Lin, "Iterated local least squares microarray missing value imputation," *J. Bioinform. Comput. Biol.*, vol. 4, no. 5, pp. 935–957, Oct. 2006.
- [47] "Hyperspectral remote sensing scenes," accessed: 2019-10-29. [Online]. Available: http://www.ehu.es/ccwintco/index.php/Hyperspectral_Remote_Sensing_Scenes
- [48] C. Higuera, K. Gardiner, and K. Cios, "Self-organizing feature maps identify proteins critical to learning in a mouse model of down syndrome," *PLoS ONE*, vol. 10, no. 6: e0129126, 2015.