

CUR matrix decompositions for improved data analysis

Michael W. Mahoney* and Petros Drineas†

*Department of Mathematics, Stanford University, Stanford, California 94305, mmahoney@cs.stanford.edu, and †Department of Computer Science, Rensselaer Polytechnic Institute, Troy, New York 12180, drinep@cs.rpi.edu

Submitted to Proceedings of the National Academy of Sciences of the United States of America

Principal components analysis, and more generally the Singular Value Decomposition, are fundamental data analysis tools that express a data matrix in terms of a sequence of orthogonal or uncorrelated vectors of decreasing importance. Unfortunately, being linear combinations of up to all the data points, these vectors are notoriously difficult to interpret in terms of the data and processes generating the data. In this paper, we develop CUR matrix decompositions for improved data analysis. CUR decompositions are low-rank matrix decompositions that are explicitly expressed in terms of a small number of actual columns and/or actual rows of the data matrix. Since they are constructed from actual data elements, CUR decompositions are interpretable by practitioners of the field from which the data are drawn (to the extent that the original data are). We present a novel algorithm that preferentially chooses columns and rows that exhibit high “statistical leverage” and thus that (in a very precise statistical sense) exert a disproportionately large “influence” on the best low-rank fit of the data matrix. By selecting columns and rows in this manner, we obtain improved relative-error and constant-factor approximation guarantees in worst-case analysis (as opposed to the much coarser additive-error guarantees of prior work). In addition, since the construction involves computing quantities with a natural and widely-studied statistical interpretation, we can leverage ideas from diagnostic regression analysis to employ these matrix decompositions for exploratory data analysis.

Modern data sets are often represented by large matrices since an $m \times n$ real-valued matrix A provides a natural structure for encoding information about m objects, each of which is described by n features. Examples of such objects include documents, genomes, stocks, hyperspectral images, and web groups. Examples of the corresponding features are terms, environmental conditions, temporal resolution, frequency resolution, and individual web users. In many cases, an important step in the analysis of such data is to construct a compressed representation of A that may be easier to analyze and interpret in light of a corpus of field-specific knowledge. The most common such representation is obtained by truncating the Singular Value Decomposition (SVD) at some number $k \ll \min\{m, n\}$ terms. For example, Principal Components Analysis (PCA) is just this procedure applied to a suitably normalized data correlation matrix.

Recall the SVD of a general matrix $A \in \mathbb{R}^{m \times n}$. Given A , there exist orthogonal matrices $U = [u^1 u^2 \dots u^m] \in \mathbb{R}^{m \times m}$ and $V = [v^1 v^2 \dots v^n] \in \mathbb{R}^{n \times n}$, where $\{u^t\}_{t=1}^m \in \mathbb{R}^m$ and $\{v^t\}_{t=1}^n \in \mathbb{R}^n$ are such that

$$U^T A V = \Sigma = \mathbf{diag}(\sigma_1, \dots, \sigma_\rho),$$

where $\Sigma \in \mathbb{R}^{m \times n}$, $\rho = \min\{m, n\}$, $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_\rho \geq 0$, and $\mathbf{diag}(\cdot)$ represents a diagonal matrix with the specified elements on the diagonal. Equivalently, $A = U \Sigma V^T$. The three matrices U , V , and Σ constitute the SVD of A [1]—the σ_i are the singular values of A and the vectors u^i and v^i are the i -th left and the i -th right singular vectors, respectively—and $O(\min\{mn^2, m^2n\})$ time suffices to compute them.

The SVD is widely used in data analysis (often via methods such as PCA), in large part because the subspaces spanned by the vectors (typically obtained after truncating the SVD to some small number k of terms) provide the best rank- k approximation to the data matrix A .

If $k \leq r = \text{rank}(A)$ and we define $A_k = \sum_{t=1}^k \sigma_t u^t v^{tT}$, then

$$\|A - A_k\|_F^2 = \min_{X \in \mathbb{R}^{m \times n}: \text{rank}(X) \leq k} \|A - X\|_F^2,$$

i.e., the distance, as measured by the Frobenius norm $\|\cdot\|_F$, where $\|A\|_F^2 = \sum_{ij} A_{ij}^2$, between A and any rank k approximation to A , is minimized by A_k [1].

Note, though, that although the truncated SVD is widely-used, the vectors u^i and v^i themselves may lack any meaning in terms of the field from which the data are drawn. For example, the eigenvector

$$[(1/2) \text{ age} - (1/\sqrt{2}) \text{ height} + (1/2) \text{ income}],$$

being one of the significant uncorrelated “factors” or “features” from a data set of people’s features, is not particularly informative or meaningful. This fact should not be surprising. After all, the singular vectors are mathematical abstractions that can be calculated for *any* data matrix. They are not “things” with a “physical” reality.

Nevertheless, data analysts often fall prey to a temptation for reification, *i.e.*, for assigning a “physical” meaning or interpretation to all large singular components. In certain special cases, *e.g.*, a data set consisting of points drawn from a multivariate normal distribution on the plane, as in Figure 1A, the principal components may be interpreted in terms of, *e.g.*, the directions of the axes of the ellipsoid from which the data are drawn. In most cases, however, *e.g.*, when the data are drawn from the union of two normals as in Figure 1B, such reification is not valid. In this and other examples it would be difficult to interpret these directions meaningfully in terms of processes generating the data. Although reification is certainly justified in some cases, such an interpretative claim cannot arise from the mathematics alone, but instead requires an intimate knowledge of the field from which the data are drawn [2].

To understand better the reification issues in modern biological applications, consider a synthetic data set introduced by Wall *et al.* [3] to model oscillatory and exponentially decaying patterns of gene expression from Cho *et al.* [4]. The data matrix consists of 14 expression level assays (columns of A) and 2000 genes (rows of A), corresponding to a 2000×14 matrix A . Genes have one of three types of transcriptional response: noise (1600 genes); noisy sine pattern (200 genes); and noisy exponential pattern (200 genes). Figures 1C and 1D present the “biological” data, *i.e.*, overlays of five noisy sine wave genes and five noisy exponential genes, respectively; Figure 1E presents first and second singular vectors of the data matrix, along with the original sine pattern and exponential pattern that generated

Reserved for Publication Footnotes

the data; and Figure 1F shows that the data cluster well in the space spanned by the top two singular vectors, which in this case account for 64% of the variance in the data. Note, though, that the top two singular vectors display both oscillatory and decaying properties, and thus they are not easily interpretable as “latent factors” or “fundamental modes” of the original “biological” processes generating the data.

This is problematic when one is interested in extracting insight from the output of data analysis algorithms. For example, biologists are typically more concerned with actual patients than eigenpatients, and researchers can more easily assay actual genes than eigengenes. Indeed, after describing the many uses of the vectors provided by the SVD and PCA in DNA microarray analysis, Kuruville *et al.* [5] bluntly conclude that “While very efficient basis vectors, the (singular) vectors themselves are completely artificial and do not correspond to actual (DNA expression) profiles. ... Thus, it would be interesting to try to find basis vectors for all experiment vectors, using actual experiment vectors and not artificial bases that offer little insight.”

These concerns about reification and interpreting linear combinations of data elements lie at the heart, conceptually and historically, of SVD-based data analysis methods. Recall that Spearman—a social scientist interested in models of human intelligence—invented factor analysis [6, 2]. He computed the first principal component of a battery of mental tests, much as depicted in the first singular vector of Figure 1B, and he invalidly reified it as an entity, calling it “g” or “general intelligence factor.” Subsequent principal components, such as the second singular vector of Figure 1B, were reified as so-called “group factors.” This provided the basis for ranking of individuals on a single intelligence scale, as well as such dubious social applications of data analysis as the 11+ examinations in Britain and the involuntary sterilization of imbeciles in Virginia [6, 2]. See Gould [2] for an enlightening and sobering discussion of invalid reification of singular vectors in a social scientific application of data analysis.

Main Contribution

We formulate and address this problem of constructing low-rank matrix approximations that depend on actual data elements. As with the SVD, the decomposition we desire (1) should have provable worst-case optimality and algorithmic properties; (2) should have a natural statistical interpretation associated with its construction; and (3) should perform well in practice. To this end, we develop and apply novel *CUR matrix decompositions*, *i.e.*, low-rank matrix decompositions that are explicitly expressed in terms of a small number of actual columns and/or actual rows of the original data matrix. Given an $m \times n$ matrix A , we decompose it as a product of three matrices, C , U , and R , where C consists of a small number of actual columns of A , R consists of a small number of actual rows of A , and U is a small carefully constructed matrix that guarantees that the product CUR is “close” to A . Of course, the extent to which $A \approx CUR$, and relatedly the extent to which CUR can be used in place of A or A_k in data analysis tasks, will depend sensitively on the choice of C and R , as well as on the construction of U .

To develop intuition about how such decompositions might behave, consider the previous pedagogical examples. In the data set consisting of the union of two normal distributions, one data point from each normal distribution (as opposed to a vector sitting between the axes of the two normals) could be chosen. Similarly, in the synthetic data set of Wall *et al.*, one could choose one sinusoid and one exponential function, as opposed to a linear combination of both. Finally, in the applications of Kuruville *et al.*, actual experimental DNA expression profiles, rather than artificial eigenprofiles, could be chosen. Thus, C and/or R can be used in place of the eigencolumns and eigenrows, but since they consist of actual data elements they will be interpretable in terms of the field from which the data are drawn (to the extent that the original data points and/or features are interpretable).

Prior CUR Matrix Decompositions

Within the numerical linear algebra community, Stewart developed the quasi-Gram-Schmidt method and applied it to a matrix and its transpose to obtain a CUR matrix decomposition [7, 8]. Similarly, Goreinov, Tyrtshnikov, and Zamarashkin developed a CUR matrix decomposition (a *pseudoskeleton approximation*) and related the choice of columns and rows to a “maximum uncorrelatedness” concept [9, 10].

Within the theoretical computer science community, much work has followed that of Frieze, Kannan, and Vempala [11], who randomly sample columns of A according to a probability distribution that depends on the Euclidean norms of those columns. If the number of chosen columns is polynomial in k and $1/\epsilon$ (for some error parameter ϵ), then worst-case additive-error guarantees of the form

$$\|A - P_C A\|_F \leq \|A - A_k\|_F + \epsilon \|A\|_F \quad [1]$$

can be obtained, with high probability. Here $P_C A$ denotes the projection of A on the subspace spanned by the columns of C . Subsequently, Drineas, Kannan, and Mahoney [12] constructed an additive-error CUR matrix decomposition by choosing columns and rows simultaneously.¹ In two passes over the matrix A , they randomly construct a matrix C of columns, a matrix R of rows, and a matrix U such that

$$\|A - CUR\|_F \leq \|A - A_k\|_F + \epsilon \|A\|_F, \quad [2]$$

with high probability.

The additive-error algorithms of [11] and [12] were motivated by resource-constrained computational environments, *e.g.*, where the input matrices are extremely large or where only a very small fraction of the data is actually available, and in those applications they are appropriate. For example, this additive-error CUR matrix decomposition has been successfully applied to applications such as hyperspectral image analysis, recommendation system analysis, and DNA SNP analysis [13, 14]. These additive-error matrix decompositions are, however, quite coarse in worst-case. Moreover, the insights provided by their sampling probabilities into the data are limited—the probabilities are often uniform due to data preprocessing, or they may correspond, *e.g.*, simply to the degree of a node if the data matrix is derived from a graph.

Statistical leverage and improved matrix decompositions

To construct C (similarly R), we will compute an “importance score” for each column of A , and we will randomly sample a small number of columns from A using that score as an importance sampling probability distribution. This importance score (see Equation (3) below) depends on the matrix A , and it has a natural interpretation as capturing the “statistical leverage” or “influence” of a given column on the best low-rank fit of the data matrix. By preferentially choosing columns that exert a disproportionately large influence on the best low-rank fit (as opposed to procedures that sample columns that have larger Euclidean norm, or empirical variance, as in prior work), we will ensure that CUR is nearly as good as A_k at capturing the dominant part of the spectrum of A . In addition, by choosing “high-statistical-leverage” or “highly-influential” columns, we can leverage ideas from diagnostic regression analysis to apply CUR matrix decompositions as a tool for exploratory data analysis.

To motivate our choice of importance sampling scores, recall that we can express the j -th column of A (denoted by A^j) *exactly* as

$$A^j = \sum_{\xi=1}^r (\sigma_{\xi} u_{\xi}^j) v_{\xi}^j,$$

¹ The algorithms of [11] and [12] provide worst-case *additive-error* guarantees since the “additional error” in (1) and (2) is $\epsilon \|A\|_F$, which is independent of the “base error” of $\|A - A_k\|_F$. This should be contrasted with a *relative-error* guarantee of the form (4), in which the “additional error” is $\epsilon \|A - A_k\|_F$, for $\epsilon > 0$ arbitrarily small, and thus the total error is bounded by $(1 + \epsilon) \|A - A_k\|_F$. It should also be contrasted with a *constant-factor* guarantee of the form (5), in which the “additional error” is $\gamma \|A - A_k\|_F$, for some constant γ .

where $r = \text{rank}(A)$ and where v_j^ξ is the j -th coordinate of the ξ -th right singular vector. That is, the j -th column of A is a linear combination of all the left singular vectors and singular values, and the elements of the j -th row of V are the coefficients. Thus, we can approximate A^j as a linear combination of the top k left singular vectors and corresponding singular values as

$$A^j \approx \sum_{\xi=1}^k (\sigma_\xi u^\xi) v_j^\xi.$$

Since we seek columns of A that are simultaneously correlated with the span of all top k right singular vectors, we then compute the *normalized statistical leverage scores*:

$$\pi_j = \frac{1}{k} \sum_{\xi=1}^k (v_j^\xi)^2, \quad [3]$$

for all $j = 1, \dots, n$. With this normalization, it is straight-forward to show that $\pi_j \geq 0$ and that $\sum_{j=1}^n \pi_j = 1$, and thus that these scores form a probability distribution over the n columns.

Our main algorithm for choosing columns from a matrix—we will call it `COLUMNSELECT`—takes as input *any* $m \times n$ matrix A , a rank parameter k , and an error parameter ϵ , and then performs the following steps:

1. First, compute v^1, \dots, v^k (the top k right singular vectors of A) and the normalized statistical leverage scores of Equation (3).
2. Second, keep the j -th column of A with probability $p_j = \min\{1, c\pi_j\}$, for all $j \in \{1, \dots, n\}$, where $c = O(k \log k / \epsilon^2)$.
3. Third, return the matrix C consisting of the selected columns of A .

With this procedure, the matrix C contains c' columns, where $c' \leq c$ in expectation and where c' is tightly concentrated around its expectation. The computation of the column leverage scores uses the top k right singular vectors of A , and this computation is the bottleneck in the running time of `COLUMNSELECT`. It can be performed in time linear in the number of nonzero elements of the matrix A times a low-degree polynomial in the rank parameter k [1]. We have proven that, with probability at least 99%, this choice of columns satisfies

$$\|A - P_C A\|_F \leq (1 + \epsilon/2) \|A - A_k\|_F, \quad [4]$$

where P_C denotes a projection matrix onto the column space of C .² See [15] for the proof of Equation (4), which depends crucially on the use of Equation (3). In some applications, this restricted CUR decomposition, $A \approx P_C A = CX$, where $X = C^+ A$, is of interest.

In other applications, one wants such a CUR matrix decomposition in terms of both columns and rows simultaneously. Our main algorithm computing a CUR matrix decomposition—we will call it `ALGORITHMCUR`—is illustrated in Figure S0 in Supplementary Materials. This algorithm takes as input *any* $m \times n$ matrix A , a rank parameter k , and an error parameter ϵ , and then it performs the following steps:

1. Run `COLUMNSELECT` on A with $c = O(k \log k / \epsilon^2)$ to choose columns of A and construct the matrix C .
2. Run `COLUMNSELECT` on A^T with $r = O(k \log k / \epsilon^2)$ to choose rows of A (columns of A^T) and construct the matrix R .
3. Define the matrix U as $U = C^+ AR^+$, where X^+ denotes a Moore-Penrose generalized inverse of the matrix X [17].

As with our algorithm for selecting columns, the running time of `ALGORITHMCUR` is dominated by computation of the column and row leverage scores. For this choice of C , U , and R , we will prove that

$$\|A - CUR\|_F \leq (2 + \epsilon) \|A - A_k\|_F, \quad [5]$$

with probability at least 98%.³ First, since $U = C^+ AR^+$, it immediately follows that

$$\|A - CUR\|_F = \|A - CC^+ AR^+ R\|_F.$$

Adding and subtracting $CC^+ A$ and applying the triangle inequality for the Frobenius norm, we get

$$\begin{aligned} \|A - CUR\|_F &\leq \|A - CC^+ A\|_F + \|CC^+ A - CC^+ AR^+ R\|_F \\ &\leq \|A - CC^+ A\|_F + \|A - AR^+ R\|_F \\ &= \|A - P_C A\|_F + \|A - AP_R\|_F. \end{aligned} \quad [6]$$

Inequality (6) follows since CC^+ is a projection matrix and thus does not increase the Frobenius norm, and the last equality follows since $P_C = CC^+$ and similarly $P_R = R^+ R$. Since `ALGORITHMCUR` chooses columns and rows by calling `COLUMNSELECT` on A and A^T , respectively, Equation (5) follows by two applications of Equation (4). Note that $r = c$ for `ALGORITHMCUR`, and also that the failure probability for `ALGORITHMCUR` is at most twice the failure probability of `COLUMNSELECT`, since the latter algorithm may fail when applied to columns or when applied to rows, independently.

Although one might like to fix a rank parameter k and choose k columns and/or rows deterministically according to some criterion—*e.g.*, such as to define a parallelepiped of maximal volume over all $\binom{n}{k}$ such parallelepipeds, or to span a subspace that “captures” a maximal amount of variance from A over all $\binom{n}{k}$ such subspaces—most such criteria would lead to intractable combinatorial optimization problems [9, 10, 18]. Thus, `ALGORITHMCUR` takes advantage of oversampling (choosing slightly more than k columns) and randomness as computational resources in order to obtain its strong provable approximation guarantees.

Note that the quantities in Equation (3) are, up to scaling, equal to the diagonal elements of the so-called “hat matrix,” *i.e.*, the projection matrix onto the span of the top k right singular vectors of A [19, 20]. As such, they have a natural statistical interpretation as a “leverage score” or “influence score” associated with each of the data points [19, 21, 20]. In particular, π_j quantifies the amount of “leverage” or “influence” exerted by the j -th column of A on its optimal low-rank approximation. Furthermore, these quantities have been widely-used for outlier identification in diagnostic regression analysis [22, 23]. As such, using these scores to select columns not only is crucial for our improved worst-case bounds but also aids in exploratory data analysis.

Diagnostic Data Analysis Applications

Implementation of `ALGORITHMCUR` is straightforward. We have applied this algorithm to data analysis problems in several application domains: Internet term-document data analysis (see Figure 2 and Figures S2 to S5 in Supplementary Materials); genetics (see Figure 3); and social science (see Figure S1 in Supplementary Materials). In practice, we typically only need to sample a number of columns and/or rows that is a small constant, *e.g.*, between 2 and 4, times the input rank parameter k . In addition, not only can we perform common data analysis tasks (such as clustering and classification) for which the basis provided by truncating the SVD is often used, but we can also use the normalized leverage scores to explore the data and identify whether there are any disproportionately “important” data elements. (Note that since CUR decompositions are low-rank approximations that use information in the top k singular subspaces, their domain of applicability is not expected to be broader than that of the SVD.)

²The randomness and failure probability in `COLUMNSELECT` and `ALGORITHMCUR` are over the choices made by the algorithm and not over the input data matrix. The quality-of-approximation bound (4) holds for *any* input matrix A , regardless of how A is constructed; its proof relies on matrix perturbation theory [15, 16]. The arbitrarily-chosen failure probability can be set to any $\delta > 0$ by repeating the algorithm $O(\log(1/\delta))$ times and taking the best of the results.

³This can be improved to $1 + \epsilon$ using a somewhat more complicated algorithm [15].

Internet term-document data are a common application of SVD-based techniques, often via latent semantic analysis [24, 25]. The Open Directory Project (ODP) [26] is a multilingual open content directory of World Wide Web links. TechTC (Technion Repository of Text Categorization Datasets) is a publicly-available benchmark set of term-document matrices from ODP with varying categorization difficulty [27]. Each matrix of the TechTC data set encodes information from ca. 150 documents from two different ODP categories. To illustrate our method, we focused on four data sets such that the documents clustered well into two classes when projected in a low-dimensional space spanned by the top few left singular vectors.⁴ (See Table S1 in Supplementary Materials for a description of the data and also Figure 2 and Figures S3 to S5 in Supplementary Materials for more details.)

For example, consider the collection of 139 documents (each described with respect to 15,170 terms) on two topics: US:Florida (id:11346) and US:Indiana:Evansville (id:10567). The first topic has 71 documents, and the second has 68; the topic names are descriptive; and the sparsity of the associated document-term matrix is 2.1%. Projecting the documents on the top two eigenterms and then running k -means clustering on the projected data leads to a clustering which has a Pearson correlation coefficient of 0.85 with the (provided) ground truth, is illustrated in Figure 2A. This class separation implies that the documents are semantically well-represented by low-rank approximation via the SVD, even though they are not numerically low-rank. The singular vectors, however, are dense; they contain negative entries; and they are not easily interpretable in terms of natural languages or the ODP hierarchy.

Not only does a CUR matrix decomposition capture the Frobenius norm of the matrix (data not shown), but it can be used to cluster the documents. Figure 2B shows the 139 documents projected on the best rank-2 approximation to the subspace spanned by the top five “highest-leverage” terms. (Two was chosen as the dimensionality of the low-dimensional space since the documents belong to one of two categories—the slowly decaying spectrum provides no guidance in this case. Five was chosen as the number of selected columns based on an analysis of the leverage scores, as described below, and other choices yielded worse results.) In this case, the class separation is quite pronounced—the Pearson correlation coefficient is 0.94, which is improved since CUR provides a low-dimensional space that respects the sparsity in the data. Of course, the data are much more axis-aligned in this low-dimensional space, largely because this space is the span of a small number of actual/interpretable terms, each of which is extremely sparse.

As an example of how we can leverage ideas from diagnostic regression analysis to explore the data, consider Figure 2C, which shows the statistical leverage scores of all 15,170 terms. The leverage scores of the top five terms—florida (.099482), evansville (.042291), south (.026892), miami (.016890), and information (.011792), as seen in Table S1 in Supplementary Materials, are orders of magnitude larger than the uniform leverage scores, equal to $1/n$, where $n = 15,170$ here. This, coupled with the obvious relevance of these terms to the task at hand, suggests that these five terms are particularly important or influential in this low-dimensional clustering/classification problem. Further evidence supporting this intuition follows from Figure 2D, which shows the information gain (IG) statistic for each of the 15,170 terms, again with the top five terms highlighted. (Recall that the IG for the i -th term is defined as $IG_i = |f_{i,1} - f_{i,2}|$, where $f_{i,1}$ is the frequency of the i -th term within the documents of the first category, and $f_{i,2}$ is the frequency of the i -th term within the documents of the second category. In particular, note that it is a supervised metric, *i.e.*, the topic of each document is known prior to computing it.)

The truncated SVD has also been applied to gene expression data [28, 29, 5, 3] in systems biology applications (where one wants to understand the relationship between actual genes) and in clinical applications (where one wants to classify tissue samples from actual patients with and without a disease). See Figure 3A for a raster plot

of a typical such data set [30], consisting of $m = 31$ patients with three different cancer types (gastrointestinal stromal tumor (GIST), leiomyosarcoma (LEIO), and synovial sarcoma (SARC)) with respect to $n = 5520$ genes.⁵ For this data set, the top two eigen-genes/eigenpatients capture 37% of the variance of the matrix, and reliable clustering can be performed in the low-dimensional space provided by the truncated SVD, as illustrated in Figure 3D. That is, there is clear separation between the three different cancer types, and running k -means on the projected data perfectly separates the three different classes of patients.

Figure 3B shows the statistical leverage scores of the 5520 genes, with red stars indicating the top 12 genes whose leverage score is well above the average uniform leverage score of approximately $0.18 \cdot 10^{-3}$. (Here, 12 was chosen by trial-and-error, and applying the IG metric to select the 12 genes that are most differentiating between the three cancer types leads to selecting 8 out of these 12 “high-leverage” genes.) Again, the leverage scores are quite nonuniform; random sampling according to the probabilities of Equation (3) easily recovers the highest scoring genes; and these genes are the most interesting and useful data points for the application of interest. Figure 3E illustrates a raster plot of those 12 genes—notice that we select genes that are overexpressed for the gastrointestinal stromal tumor and the synovial sarcoma. The fact that our leverage scores pick these “highly-influential” genes explains why the three different cancer types are reliably separated in the space spanned by these 12 genes, as demonstrated by Figure 3F. That is, running k -means using only these 12 genes leads to a perfect clustering of the 31 patients to their respective cancer types since these actual genes exert a disproportionately large influence on the best low-rank fit of the data matrix (presumably in part for reasons implicit in the data preparation).

Finally, and crucially for medical interpretation of this analysis in clinical settings, some of the selected genes are well-known to have expression patterns that correlate with various cancers, and thus they can be further studied in the lab. Most notably, PRKCCQ has been associated with gastrointestinal stromal tumors; PRAME has been known to confer growth or survival advantages in human cancers; BCHE has been associated with colorectal carcinomas and colonic adenocarcinomas; and SFRP1 and CRABP1 have been associated with a wide range of cancers. There is no need for the reification of the artificial singular directions that offer little insight into the underlying biology [5].

From the perspective of diagnostic regression analysis, the extreme nonuniformity in “importance” of individual data elements is quite surprising. In those applications, one is taught to be wary of data points with leverage greater than 2 or 3 times the average value [20], investigating them to see whether they are errors or outliers. Moreover, there is no *a priori* reason that the nonuniformity in this statistic should correlate with nonuniformity in a supervised metric like information gain. In our experience, however, this nonuniformity is not uncommon. (It is, of course, far from ubiquitous. For example, Congressional roll call data [31, 32] are much more homogeneous; if $k = 2$ and n is the number of representatives, the highest leverage score for any representative is only $1.38k/n$.) Most often, this phenomenon arises in cases where SVD-based methods are used for computational convenience, rather than since the statistical assump-

⁴As is common with term-document data, the TechTC matrices are very sparse, and they are not numerically low-rank. For example, the top 2.5% and 5% of the non-zero singular values of these matrices capture (on average) 5.5% and 12.5% of the Frobenius norm, respectively. For data sets in which a low-dimensional space provided by the SVD failed to capture the category separation, CUR matrix decompositions performed correspondingly poorly.

⁵The original data contained 46 patients from six different cancer types; here we consider a subset of 31 patients with one of three cancer types. Linear dimensionality reduction techniques such as SVD could accurately separate patients suffering from gastrointestinal stromal tumor, leiomyosarcoma, and synovial sarcoma. Including patients suffering from the remaining three cancer types (liposarcoma, malignant fibrous histiocytoma, and schwannoma) led to several inaccurate classifications, even using all 5520 genes. This is consistent with the finding that hierarchical clustering in all 5520 genes resulted to five (instead of six) clusters, with several incorrect assignments [30]. To illustrate the use of low-rank CUR matrix decompositions, we focused on the subset of the data that were amenable to analysis with the truncated SVD. Similar, but noisier, results are obtained on the full data set.

tions underlying its use are satisfied by the data. This suggests the use of these “leverage scores” more generally in modern massive data set analysis. Intuitively, conditioned on being reliable, more “outlier-like” data points may be the most important and informative.

Conclusion

Although the SVD *per se* cannot be blamed for its misapplication, the desire for interpretability in data analysis is sufficiently strong so as to argue for *interpretable low-rank matrix decompositions*. Even when an immediate application of the truncated SVD is not appropriate, the low-rank matrix approximation thereby obtained is a funda-

mental building block of some of the most widely-used data analysis methods, including PCA, Multidimensional Scaling, Factor Analysis, and many of the recently-developed techniques to extract non-linear structure from data. Thus, although here we have focused on the use of CUR matrix decompositions for the improved interpretability of SVD-based data analysis methods, we expect that their promise is much more general.

ACKNOWLEDGMENTS. We would like to thank Deepak Agarwal, Christos Boutsidis, Gene Golub, Ravi Kannan, Kenneth Kidd, Jon Kleinberg, and Michael Jordan for useful discussions and constructive comments. PD has been supported by the NSF through a CAREER award and a NeTS-NBD award.

- G.H. Golub and C.F. Van Loan. *Matrix Computations*. Johns Hopkins University Press, Baltimore, 1996.
- S.J. Gould. *The Mismeasure of Man*. W. W. Norton and Company, New York, 1996.
- M.E. Wall, A. Rechtsteiner, and L.M. Rocha. Singular value decomposition and principal component analysis. In D.P. Berrar, W. Dubitzky, and M. Granzow, editors, *A Practical Approach to Microarray Data Analysis*, pages 91–109. Kluwer, 2003.
- R.J. Cho, M.J. Campbell, E.A. Winzeler, L. Steinmetz, A. Conway, L. Wodicka, T.G. Wolfsberg, A.E. Gabriellian, D. Landsman, D.J. Lockhart, and R.W. Davis. A genome-wide transcriptional analysis of the mitotic cell cycle. *Molecular Cell*, 2:65–73, 1998.
- F.G. Kuruvilla, P.J. Park, and S.L. Schreiber. Vector algebra in the analysis of genome-wide expression data. *Genome Biology*, 3:research0011.1–0011.11, 2002.
- S.F. Blinkhorn. Burt and the early history of factor analysis. In N.J. Mackintosh, editor, *Cyril Burt: Fraud or Framed?* Oxford University Press, 1995.
- G.W. Stewart. Four algorithms for the efficient computation of truncated QR approximations to a sparse matrix. *Numerische Mathematik*, 83:313–323, 1999.
- M.W. Berry, S.A. Pulatova, and G.W. Stewart. Computing sparse reduced-rank approximations to sparse matrices. Technical Report UMIACS TR-2004-32 CMSC TR-4589, University of Maryland, College Park, MD, 2004.
- S.A. Goreinov, E.E. Tyrtshnikov, and N.L. Zamarashkin. A theory of pseudoskeleton approximations. *Linear Algebra and Its Applications*, 261:1–21, 1997.
- S.A. Goreinov and E.E. Tyrtshnikov. The maximum-volume concept in approximation by low-rank matrices. *Contemporary Mathematics*, 280:47–51, 2001.
- A. Frieze, R. Kannan, and S. Vempala. Fast Monte-Carlo algorithms for finding low-rank approximations. *Journal of the ACM*, 51(6):1025–1041, 2004.
- P. Drineas, R. Kannan, and M.W. Mahoney. Fast Monte Carlo algorithms for matrices III: Computing a compressed approximate matrix decomposition. *SIAM Journal on Computing*, 36:184–206, 2006.
- M.W. Mahoney, M. Maggioni, and P. Drineas. Tensor-CUR decompositions for tensor-based data. In *Proceedings of the 12th Annual ACM SIGKDD Conference*, pages 327–336, 2006.
- P. Paschou, M. W. Mahoney, A. Javed, J. R. Kidd, A. J. Pakstis, S. Gu, K. K. Kidd, and P. Drineas. Intra- and interpopulation genotype reconstruction from tagging SNPs. *Genome Research*, 17(1):96–107, 2007.
- P. Drineas, M.W. Mahoney, and S. Muthukrishnan. Relative-Error CUR Matrix Decompositions. *SIAM Journal on Matrix Analysis and Applications*, 30:844–881, 2008.
- G.W. Stewart and J.G. Sun. *Matrix Perturbation Theory*. Academic Press, New York, 1990.
- M.Z. Nashed, editor. *Generalized Inverses and Applications*. Academic Press, New York, 1976.
- M.R. Garey and D.S. Johnson. *Computers and Intractability: a Guide to the Theory of NP-Completeness*. W.H. Freeman and Company, New York, 1979.
- D.C. Hoaglin and R.E. Welsch. The hat matrix in regression and ANOVA. *The American Statistician*, 32(1):17–22, 1978.
- S. Chatterjee and A.S. Hadi. *Sensitivity Analysis in Linear Regression*. John Wiley & Sons, New York, 1988.
- S. Chatterjee and A.S. Hadi. Influential observations, high leverage points, and outliers in linear regression. *Statistical Science*, 1(3):379–393, 1986.
- P.F. Velleman and R.E. Welsch. Efficient computing of regression diagnostics. *The American Statistician*, 35(4):234–242, 1981.
- S. Chatterjee, A.S. Hadi, and B. Price. *Regression Analysis by Example*. John Wiley & Sons, New York, 2000.
- S.T. Deerwester, S.T. Dumais, G.W. Furnas, T.K. Landauer, and R. Harshman. Indexing by latent semantic analysis. *Journal of the American Society for Information Science*, 41(6):391–407, 1990.
- T.K. Landauer, D. Laham, and M. Derr. From paragraph to graph: latent semantic analysis for information visualization. *Proc. Natl. Acad. Sci. USA*, 101:5214–5219, 2004.
- Open Directory Project. <http://www.dmoz.org/>.
- E. Gabrilovich and S. Markovitch. Text categorization with many redundant features: using aggressive feature selection to make SVMs competitive with C4.5. In *Proceedings of the 21th International Conference on Machine Learning*, pages 41–48, 2004.
- O. Alter, P.O. Brown, and D. Botstein. Singular value decomposition for genome-wide expression data processing and modeling. *Proc. Natl. Acad. Sci. USA*, 97(18):10101–10106, 2000.
- N.S. Holter, M. Mitra, A. Maritan, M. Cieplak, J.R. Banavar, and N.V. Fedoroff. Fundamental patterns underlying gene expression profiles: Simplicity from complexity. *Proc. Natl. Acad. Sci. USA*, 97:8409–8414, 2000.
- T. Nielsen, R.B. West, S.C. Linn, O. Alter, M.A. Knowling, J. O’Connell, S. Zhu, M. Fero, G. Sherlock, J.R. Pollack, P.O. Brown, D. Botstein, and M. van de Rijn. Molecular characterisation of soft tissue tumours: a gene expression study. *Lancet*, 359(9314):1301–1307, 2002.
- K.T. Poole and H. Rosenthal. Patterns of congressional voting. *American Journal of Political Science*, 35:228–278, 1991.
- M.A. Porter, P.J. Mucha, M.E.J. Newman, and C.M. Warmbrand. A network analysis of committees in the U.S. House of Representatives. *Proc. Natl. Acad. Sci. USA*, 102(20):7057–7062, 2005.

Captions

Figure 1 Applying the SVD to data matrices A . (A) 1000 points on the plane, corresponding to a 1000×2 matrix A , (and the two principal components) drawn from a multivariate normal distribution. (B) 1000 points on the plane (and the two principal components) drawn from a more complex distribution, in this case the union of two multivariate normal distributions. (C-F) A synthetic data set considered by Wall, et al. [3] to model oscillatory and exponentially decaying patterns of gene expression from Cho, et al. [4], as described in the text. (C) Overlays of five noisy sine wave genes. (D) Overlays of five noisy exponential genes. (E) The first and second singular vectors of the data matrix (which account for 64% of the variance in the data), along with the original sine pattern and exponential pattern that generated the data. (F) Projection of the synthetic data on its top two singular vectors. Although the data cluster well in the low-dimensional space, the top two singular vectors are completely artificial and do not offer insight into the oscillatory and exponentially decaying patterns that generated the data.

Figure 2 Application of ALGORITHM CUR to Internet term-document data. The matrix consists of 139 documents from TechTC on two topics: US:Indiana:Evansville (id:10567) and US:Florida (id:11346). (A) A projection of the documents on the two-dimensional space spanned by the top two eigenterms. (B) A projection of the documents on the best two-dimensional approximation to the space spanned by the five actual terms with the highest “statistical leverage” scores of Equation (3). (C) The statistical leverage scores of the 15,170 terms, with the top five terms highlighted. (D) The information gain statistic of the 15,170 terms, with the top five terms highlighted.

Figure 3 Application of ALGORITHM CUR to cancer microarray data: gene expression patterns of patients with soft tissue tumors analyzed with cDNA microarrays. (A) Raster plot of the data. (B) Statistical leverage scores for each gene. Red stars indicate the 12 genes with the highest leverage scores, and the red dashed line indicates the uniform leverage scores. (C) Information gain for each gene, with the highest scoring genes marked. (D) A projection of the 31 patients in the subspace spanned by the top three eigengenes of the full data matrix. (E) Raster plot of the 31 patients with respect to the top 12 genes from Panel B. A (+) sign in front of the name of a gene indicates that this gene was also selected by the (supervised) information gain metric. (F) A 3D plot of the 31 patients in the span of the genes of Panel E.