

Exercise 2.45 *Create a list of the five most important things that you learned about high dimensions.*

Exercise 2.46 *Write a short essay whose purpose is to excite a college freshman to learn about high dimensions.*

3 Best-Fit Subspaces and Singular Value Decomposition (SVD)

3.1 Introduction

In this chapter, we examine the *Singular Value Decomposition* (SVD) of a matrix. Consider each row of an $n \times d$ matrix A as a point in d -dimensional space. The singular value decomposition finds the best-fitting k -dimensional subspace for $k = 1, 2, 3, \dots$, for the set of n data points. Here, “best” means minimizing the sum of the squares of the perpendicular distances of the points to the subspace, or equivalently, maximizing the sum of squares of the lengths of the projections of the points onto this subspace.⁴ We begin with a special case where the subspace is 1-dimensional, namely a line through the origin. We then show that the best-fitting k -dimensional subspace can be found by k applications of the best fitting line algorithm, where on the i^{th} iteration we find the best fit line perpendicular to the previous $i - 1$ lines. When k reaches the rank of the matrix, from these operations we get an exact decomposition of the matrix called the *singular value decomposition*.

In matrix notation, the singular value decomposition of a matrix A with real entries (we assume all our matrices have real entries) is the factorization of A into the product of three matrices, $A = UDV^T$, where the columns of U and V are orthonormal⁵ and the matrix D is diagonal with positive real entries. The columns of V are the unit length vectors defining the best fitting lines described above (the i^{th} column being the unit-length vector in the direction of the i^{th} line). The coordinates of a row of U will be the fractions of the corresponding row of A along the direction of each of the lines.

The SVD is useful in many tasks. Often a data matrix A is close to a low rank matrix and it is useful to find a good low rank approximation to A . For any k , the singular value decomposition of A gives the best rank- k approximation to A in a well-defined sense.

If \mathbf{u}_i and \mathbf{v}_i are columns of U and V respectively, then the matrix equation $A = UDV^T$ can be rewritten as

$$A = \sum_i d_{ii} \mathbf{u}_i \mathbf{v}_i^T.$$

Since \mathbf{u}_i is a $n \times 1$ matrix and \mathbf{v}_i is a $d \times 1$ matrix, $\mathbf{u}_i \mathbf{v}_i^T$ is an $n \times d$ matrix with the same dimensions as A . The i^{th} term in the above sum can be viewed as giving the components of the rows of A along direction \mathbf{v}_i . When the terms are summed, they reconstruct A .

⁴This equivalence is due to the Pythagorean Theorem. For each point, its squared length (its distance to the origin squared) is exactly equal to the squared length of its projection onto the subspace plus the squared distance of the point to its projection; therefore, maximizing the sum of the former is equivalent to minimizing the sum of the latter. For further discussion see Section 3.2.

⁵A set of vectors is orthonormal if each is of length one and they are pairwise orthogonal.

This decomposition of A can be viewed as analogous to writing a vector \mathbf{x} in some orthonormal basis $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d$. The coordinates of $\mathbf{x} = (\mathbf{x} \cdot \mathbf{v}_1, \mathbf{x} \cdot \mathbf{v}_2, \dots, \mathbf{x} \cdot \mathbf{v}_d)$ are the projections of \mathbf{x} onto the \mathbf{v}_i 's. For SVD, this basis has the property that for any k , the first k vectors of this basis produce the least possible total sum of squares error for that value of k .

In addition to the singular value decomposition, there is an eigenvalue decomposition. Let A be a square matrix. A vector \mathbf{v} such that $A\mathbf{v} = \lambda\mathbf{v}$ is called an eigenvector and λ the eigenvalue. When A is symmetric, the eigenvectors are orthogonal and A can be expressed as $A = VDV^T$ where the eigenvectors are the columns of V and D is a diagonal matrix with the corresponding eigenvalues on its diagonal. For a symmetric matrix A the singular values and eigenvalues are identical. If the singular values are distinct, then A 's singular vectors and eigenvectors are identical. If a singular value has multiplicity d greater than one, the corresponding singular vectors span a subspace of dimension d and any orthogonal basis of the subspace can be used as the eigenvectors or singular vectors.⁶

The singular value decomposition is defined for all matrices, whereas the more familiar eigenvector decomposition requires that the matrix A be square and certain other conditions on the matrix to ensure orthogonality of the eigenvectors. In contrast, the columns of V in the singular value decomposition, called the *right-singular vectors* of A , always form an orthogonal set with no assumptions on A . The columns of U are called the *left-singular vectors* and they also form an orthogonal set (see Section 3.6). A simple consequence of the orthonormality is that for a square and invertible matrix A , the inverse of A is $VD^{-1}U^T$.

Eigenvalues and eigenvectors satisfy $A\mathbf{v} = \lambda\mathbf{v}$. We will show that singular values and vectors satisfy a somewhat analogous relationship. Since $A\mathbf{v}_i$ is a $n \times 1$ matrix (vector), the matrix A cannot act on it from the left. But A^T , which is a $d \times n$ matrix, can act on this vector. Indeed, we will show that

$$A\mathbf{v}_i = d_{ii}\mathbf{u}_i \quad \text{and} \quad A^T\mathbf{u}_i = d_{ii}\mathbf{v}_i.$$

In words, A acting on \mathbf{v}_i produces a scalar multiple of \mathbf{u}_i and A^T acting on \mathbf{u}_i produces the same scalar multiple of \mathbf{v}_i . Note that $A^T A\mathbf{v}_i = d_{ii}^2 \mathbf{v}_i$. The i^{th} singular vector of A is the i^{th} eigenvector of the square symmetric matrix $A^T A$.

3.2 Preliminaries

Consider projecting a point $\mathbf{a}_i = (a_{i1}, a_{i2}, \dots, a_{id})$ onto a line through the origin. Then

$$a_{i1}^2 + a_{i2}^2 + \dots + a_{id}^2 = (\text{length of projection})^2 + (\text{distance of point to line})^2.$$

⁶When $d = 1$ there are actually two possible singular vectors, one the negative of the other. The subspace spanned is unique.

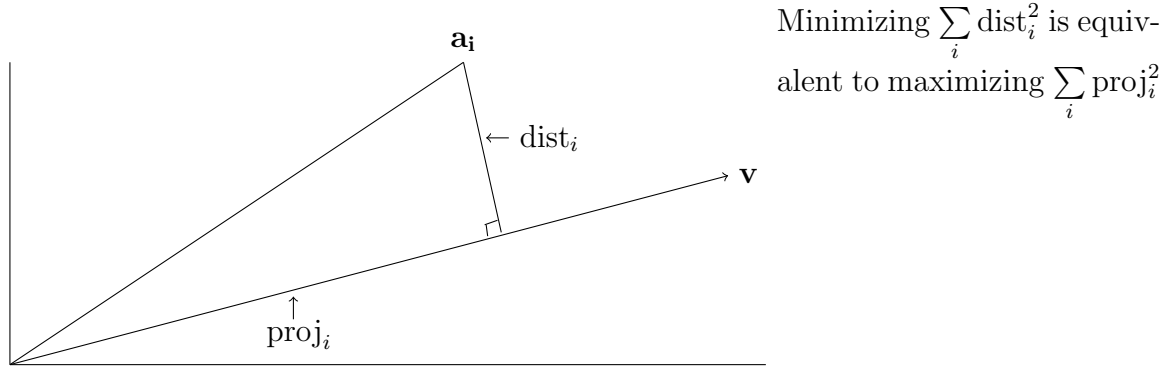


Figure 3.1: The projection of the point \mathbf{a}_i onto the line through the origin in the direction of \mathbf{v} .

This holds by the Pythagorean Theorem (see Figure 3.1). Thus

$$(\text{distance of point to line})^2 = a_{i1}^2 + a_{i2}^2 + \cdots + a_{id}^2 - (\text{length of projection})^2.$$

Since $\sum_{i=1}^n (a_{i1}^2 + a_{i2}^2 + \cdots + a_{id}^2)$ is a constant independent of the line, minimizing the sum of the squares of the distances to the line is equivalent to maximizing the sum of the squares of the lengths of the projections onto the line. Similarly for best-fit subspaces, maximizing the sum of the squared lengths of the projections onto the subspace minimizes the sum of squared distances to the subspace.

Thus we have two interpretations of the best-fit subspace. The first is that it minimizes the sum of squared distances of the data points to it. This first interpretation and its use are akin to the notion of least-squares fit from calculus.⁷ The second interpretation of best-fit-subspace is that it maximizes the sum of projections squared of the data points on it. This says that the subspace contains the maximum content of data among all subspaces of the same dimension. The choice of the objective function as the sum of squared distances seems a bit arbitrary and in a way it is. But the square has many nice mathematical properties. The first of these, as we have just seen, is that minimizing the sum of squared distances is equivalent to maximizing the sum of squared projections.

3.3 Singular Vectors

We now define the *singular vectors* of an $n \times d$ matrix A . Consider the rows of A as n points in a d -dimensional space. Consider the best fit line through the origin. Let \mathbf{v} be a unit vector along this line. The length of the projection of \mathbf{a}_i , the i^{th} row of A , onto \mathbf{v} is $|\mathbf{a}_i \cdot \mathbf{v}|$. From this we see that the sum of the squared lengths of the projections is

⁷But there is a difference: here we take the perpendicular distance to the line or subspace, whereas, in the calculus notion, given n pairs, $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, we find a line $l = \{(x, y) | y = mx + b\}$ minimizing the vertical squared distances of the points to it, namely, $\sum_{i=1}^n (y_i - mx_i - b)^2$.

$|A\mathbf{v}|^2$. The best fit line is the one maximizing $|A\mathbf{v}|^2$ and hence minimizing the sum of the squared distances of the points to the line.

With this in mind, define the *first singular vector* \mathbf{v}_1 of A as

$$\mathbf{v}_1 = \arg \max_{|\mathbf{v}|=1} |A\mathbf{v}|.$$

Technically, there may be a tie for the vector attaining the maximum and so we should not use the article “the”; in fact, $-\mathbf{v}_1$ is always as good as \mathbf{v}_1 . In this case, we arbitrarily pick one of the vectors achieving the maximum and refer to it as “the first singular vector” avoiding the more cumbersome “one of the vectors achieving the maximum”. We adopt this terminology for all uses of $\arg \max$.

The value $\sigma_1(A) = |A\mathbf{v}_1|$ is called the *first singular value* of A . Note that $\sigma_1^2 = \sum_{i=1}^n (\mathbf{a}_i \cdot \mathbf{v}_1)^2$ is the sum of the squared lengths of the projections of the points onto the line determined by \mathbf{v}_1 .

If the data points were all either on a line or close to a line, intuitively, \mathbf{v}_1 should give us the direction of that line. It is possible that data points are not close to one line, but lie close to a 2-dimensional subspace or more generally a low dimensional space. Suppose we have an algorithm for finding \mathbf{v}_1 (we will describe one such algorithm later). How do we use this to find the best-fit 2-dimensional plane or more generally the best fit k -dimensional space?

The greedy approach begins by finding \mathbf{v}_1 and then finds the best 2-dimensional subspace containing \mathbf{v}_1 . The sum of squared distances helps. For every 2-dimensional subspace containing \mathbf{v}_1 , the sum of squared lengths of the projections onto the subspace equals the sum of squared projections onto \mathbf{v}_1 plus the sum of squared projections along a vector perpendicular to \mathbf{v}_1 in the subspace. Thus, instead of looking for the best 2-dimensional subspace containing \mathbf{v}_1 , look for a unit vector \mathbf{v}_2 perpendicular to \mathbf{v}_1 that maximizes $|A\mathbf{v}|^2$ among all such unit vectors. Using the same greedy strategy to find the best three and higher dimensional subspaces, defines $\mathbf{v}_3, \mathbf{v}_4, \dots$ in a similar manner. This is captured in the following definitions. There is no apriori guarantee that the greedy algorithm gives the best fit. But, in fact, the greedy algorithm does work and yields the best-fit subspaces of every dimension as we will show.

The *second singular vector*, \mathbf{v}_2 , is defined by the best fit line perpendicular to \mathbf{v}_1 .

$$\mathbf{v}_2 = \arg \max_{\substack{\mathbf{v} \perp \mathbf{v}_1 \\ |\mathbf{v}|=1}} |A\mathbf{v}|$$

The value $\sigma_2(A) = |A\mathbf{v}_2|$ is called the *second singular value* of A . The *third singular*

vector \mathbf{v}_3 and the *third singular value* are defined similarly by

$$\mathbf{v}_3 = \arg \max_{\substack{\mathbf{v} \perp \mathbf{v}_1, \mathbf{v}_2 \\ \|\mathbf{v}\|=1}} |A\mathbf{v}|$$

and

$$\sigma_3(A) = |A\mathbf{v}_3|,$$

and so on. The process stops when we have found singular vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$, singular values $\sigma_1, \sigma_2, \dots, \sigma_r$, and

$$\max_{\substack{\mathbf{v} \perp \mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r \\ \|\mathbf{v}\|=1}} |A\mathbf{v}| = 0.$$

The greedy algorithm found the \mathbf{v}_1 that maximized $|A\mathbf{v}|$ and then the best fit 2-dimensional subspace containing \mathbf{v}_1 . Is this necessarily the best-fit 2-dimensional subspace overall? The following theorem establishes that the greedy algorithm finds the best subspaces of every dimension.

Theorem 3.1 (The Greedy Algorithm Works) *Let A be an $n \times d$ matrix with singular vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$. For $1 \leq k \leq r$, let V_k be the subspace spanned by $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$. For each k , V_k is the best-fit k -dimensional subspace for A .*

Proof: The statement is obviously true for $k = 1$. For $k = 2$, let W be a best-fit 2-dimensional subspace for A . For any orthonormal basis $(\mathbf{w}_1, \mathbf{w}_2)$ of W , $|A\mathbf{w}_1|^2 + |A\mathbf{w}_2|^2$ is the sum of squared lengths of the projections of the rows of A onto W . Choose an orthonormal basis $(\mathbf{w}_1, \mathbf{w}_2)$ of W so that \mathbf{w}_2 is perpendicular to \mathbf{v}_1 . If \mathbf{v}_1 is perpendicular to W , any unit vector in W will do as \mathbf{w}_2 . If not, choose \mathbf{w}_2 to be the unit vector in W perpendicular to the projection of \mathbf{v}_1 onto W . This makes \mathbf{w}_2 perpendicular to \mathbf{v}_1 .⁸ Since \mathbf{v}_1 maximizes $|A\mathbf{v}|^2$, it follows that $|A\mathbf{w}_1|^2 \leq |A\mathbf{v}_1|^2$. Since \mathbf{v}_2 maximizes $|A\mathbf{v}|^2$ over all \mathbf{v} perpendicular to \mathbf{v}_1 , $|A\mathbf{w}_2|^2 \leq |A\mathbf{v}_2|^2$. Thus

$$|A\mathbf{w}_1|^2 + |A\mathbf{w}_2|^2 \leq |A\mathbf{v}_1|^2 + |A\mathbf{v}_2|^2.$$

Hence, V_2 is at least as good as W and so is a best-fit 2-dimensional subspace.

For general k , proceed by induction. By the induction hypothesis, V_{k-1} is a best-fit $k-1$ dimensional subspace. Suppose W is a best-fit k -dimensional subspace. Choose an orthonormal basis $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k$ of W so that \mathbf{w}_k is perpendicular to $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k-1}$. Then

$$|A\mathbf{w}_1|^2 + |A\mathbf{w}_2|^2 + \dots + |A\mathbf{w}_{k-1}|^2 \leq |A\mathbf{v}_1|^2 + |A\mathbf{v}_2|^2 + \dots + |A\mathbf{v}_{k-1}|^2$$

since V_{k-1} is an optimal $k-1$ dimensional subspace. Since \mathbf{w}_k is perpendicular to $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k-1}$, by the definition of \mathbf{v}_k , $|A\mathbf{w}_k|^2 \leq |A\mathbf{v}_k|^2$. Thus

$$|A\mathbf{w}_1|^2 + |A\mathbf{w}_2|^2 + \dots + |A\mathbf{w}_{k-1}|^2 + |A\mathbf{w}_k|^2 \leq |A\mathbf{v}_1|^2 + |A\mathbf{v}_2|^2 + \dots + |A\mathbf{v}_{k-1}|^2 + |A\mathbf{v}_k|^2,$$

proving that V_k is at least as good as W and hence is optimal. ■

⁸This can be seen by noting that \mathbf{v}_1 is the sum of two vectors that each are individually perpendicular to \mathbf{w}_2 , namely the projection of \mathbf{v}_1 to W and the portion of \mathbf{v}_1 orthogonal to W .

Note that the n -dimensional vector $A\mathbf{v}_i$ is a list of lengths (with signs) of the projections of the rows of A onto \mathbf{v}_i . Think of $|A\mathbf{v}_i| = \sigma_i(A)$ as the *component* of the matrix A along \mathbf{v}_i . For this interpretation to make sense, it should be true that adding up the squares of the components of A along each of the \mathbf{v}_i gives the square of the “whole content of A ”. This is indeed the case and is the matrix analogy of decomposing a vector into its components along orthogonal directions.

Consider one row, say \mathbf{a}_j , of A . Since $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$ span the space of all rows of A , $\mathbf{a}_j \cdot \mathbf{v} = 0$ for all \mathbf{v} perpendicular to $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$. Thus, for each row \mathbf{a}_j , $\sum_{i=1}^r (\mathbf{a}_j \cdot \mathbf{v}_i)^2 = |\mathbf{a}_j|^2$. Summing over all rows j ,

$$\sum_{j=1}^n |\mathbf{a}_j|^2 = \sum_{j=1}^n \sum_{i=1}^r (\mathbf{a}_j \cdot \mathbf{v}_i)^2 = \sum_{i=1}^r \sum_{j=1}^n (\mathbf{a}_j \cdot \mathbf{v}_i)^2 = \sum_{i=1}^r |A\mathbf{v}_i|^2 = \sum_{i=1}^r \sigma_i^2(A).$$

But $\sum_{j=1}^n |\mathbf{a}_j|^2 = \sum_{j=1}^n \sum_{k=1}^d a_{jk}^2$, the sum of squares of all the entries of A . Thus, the sum of squares of the singular values of A is indeed the square of the “whole content of A ”, i.e., the sum of squares of all the entries. There is an important norm associated with this quantity, the Frobenius norm of A , denoted $\|A\|_F$ defined as

$$\|A\|_F = \sqrt{\sum_{j,k} a_{jk}^2}.$$

Lemma 3.2 *For any matrix A , the sum of squares of the singular values equals the square of the Frobenius norm. That is, $\sum \sigma_i^2(A) = \|A\|_F^2$.*

Proof: By the preceding discussion. ■

The vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$ are called the *right-singular vectors*. The vectors $A\mathbf{v}_i$ form a fundamental set of vectors and we normalize them to length one by

$$\mathbf{u}_i = \frac{1}{\sigma_i(A)} A\mathbf{v}_i.$$

Later we will show that \mathbf{u}_i similarly maximizes $|\mathbf{u}^T A|$ over all \mathbf{u} perpendicular to $\mathbf{u}_1, \dots, \mathbf{u}_{i-1}$. These \mathbf{u}_i are called the *left-singular vectors*. Clearly, the right-singular vectors are orthogonal by definition. We will show later that the left-singular vectors are also orthogonal.

3.4 Singular Value Decomposition (SVD)

Let A be an $n \times d$ matrix with singular vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$ and corresponding singular values $\sigma_1, \sigma_2, \dots, \sigma_r$. The left-singular vectors of A are $\mathbf{u}_i = \frac{1}{\sigma_i} A\mathbf{v}_i$ where $\sigma_i \mathbf{u}_i$ is

a vector whose coordinates correspond to the projections of the rows of A onto \mathbf{v}_i . Each $\sigma_i \mathbf{u}_i \mathbf{v}_i^T$ is a rank one matrix whose rows are the “ \mathbf{v}_i components” of the rows of A , i.e., the projections of the rows of A in the \mathbf{v}_i direction. We will prove that A can be decomposed into a sum of rank one matrices as

$$A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T.$$

Geometrically, each point is decomposed in A into its components along each of the r orthogonal directions given by the \mathbf{v}_i . We will also prove this algebraically. We begin with a simple lemma that two matrices A and B are identical if $A\mathbf{v} = B\mathbf{v}$ for all \mathbf{v} .

Lemma 3.3 *Matrices A and B are identical if and only if for all vectors \mathbf{v} , $A\mathbf{v} = B\mathbf{v}$.*

Proof: Clearly, if $A = B$ then $A\mathbf{v} = B\mathbf{v}$ for all \mathbf{v} . For the converse, suppose that $A\mathbf{v} = B\mathbf{v}$ for all \mathbf{v} . Let \mathbf{e}_i be the vector that is all zeros except for the i^{th} component which has value one. Now $A\mathbf{e}_i$ is the i^{th} column of A and thus $A = B$ if for each i , $A\mathbf{e}_i = B\mathbf{e}_i$. ■

Theorem 3.4 *Let A be an $n \times d$ matrix with right-singular vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$, left-singular vectors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r$, and corresponding singular values $\sigma_1, \sigma_2, \dots, \sigma_r$. Then*

$$A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T.$$

Proof: We first show that multiplying both A and $\sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ by \mathbf{v}_j results in equality.

$$\sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T \mathbf{v}_j = \sigma_j \mathbf{u}_j = A\mathbf{v}_j$$

Since any vector \mathbf{v} can be expressed as a linear combination of the singular vectors plus a vector perpendicular to the \mathbf{v}_i , $A\mathbf{v} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T \mathbf{v}$ for all \mathbf{v} and by Lemma 3.3,

$$A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T. \quad \blacksquare$$

The decomposition $A = \sum_i \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ is called the *singular value decomposition*, *SVD*, of A . We can rewrite this equation in matrix notation as $A = UDV^T$ where \mathbf{u}_i is the i^{th} column of U , \mathbf{v}_i^T is the i^{th} row of V^T , and D is a diagonal matrix with σ_i as the i^{th} entry on its diagonal. For any matrix A , the sequence of singular values is unique and if the singular values are all distinct, then the sequence of singular vectors is unique up to signs. However, when some set of singular values are equal, the corresponding singular vectors span some subspace. Any set of orthonormal vectors spanning this subspace can be used as the singular vectors.

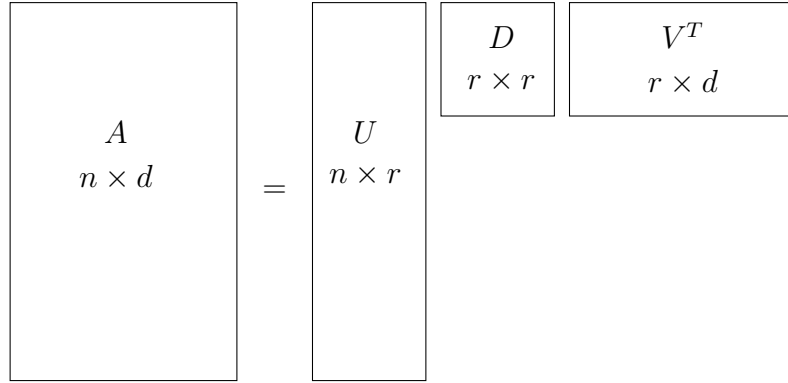


Figure 3.2: The SVD decomposition of an $n \times d$ matrix.

3.5 Best Rank- k Approximations

Let A be an $n \times d$ matrix and think of the rows of A as n points in d -dimensional space. Let

$$A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$

be the SVD of A . For $k \in \{1, 2, \dots, r\}$, let

$$A_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$

be the sum truncated after k terms. It is clear that A_k has rank k . We show that A_k is the best rank k approximation to A , where error is measured in the Frobenius norm. Geometrically, this says that $\mathbf{v}_1, \dots, \mathbf{v}_k$ define the k -dimensional space minimizing the sum of squared distances of the points to the space. To see why, we need the following lemma.

Lemma 3.5 *The rows of A_k are the projections of the rows of A onto the subspace V_k spanned by the first k singular vectors of A .*

Proof: Let \mathbf{a} be an arbitrary row vector. Since the \mathbf{v}_i are orthonormal, the projection of the vector \mathbf{a} onto V_k is given by $\sum_{i=1}^k (\mathbf{a} \cdot \mathbf{v}_i) \mathbf{v}_i^T$. Thus, the matrix whose rows are the projections of the rows of A onto V_k is given by $\sum_{i=1}^k A \mathbf{v}_i \mathbf{v}_i^T$. This last expression simplifies to

$$\sum_{i=1}^k A \mathbf{v}_i \mathbf{v}_i^T = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T = A_k.$$

■

Theorem 3.6 For any matrix B of rank at most k

$$\|A - A_k\|_F \leq \|A - B\|_F$$

Proof: Let B minimize $\|A - B\|_F^2$ among all rank k or less matrices. Let V be the space spanned by the rows of B . The dimension of V is at most k . Since B minimizes $\|A - B\|_F^2$, it must be that each row of B is the projection of the corresponding row of A onto V : Otherwise replace the row of B with the projection of the corresponding row of A onto V . This still keeps the row space of B contained in V and hence the rank of B is still at most k . But it reduces $\|A - B\|_F^2$, contradicting the minimality of $\|A - B\|_F$. ■

Since each row of B is the projection of the corresponding row of A , it follows that $\|A - B\|_F^2$ is the sum of squared distances of rows of A to V . Since A_k minimizes the sum of squared distance of rows of A to any k -dimensional subspace, from Theorem 3.1, it follows that $\|A - A_k\|_F \leq \|A - B\|_F$. ■

In addition to the Frobenius norm, there is another matrix norm of interest. Consider an $n \times d$ matrix A and a large number of vectors where for each vector \mathbf{x} we wish to compute $A\mathbf{x}$. It takes time $O(nd)$ to compute each product $A\mathbf{x}$ but if we approximate A by $A_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ and approximate $A\mathbf{x}$ by $A_k \mathbf{x}$ it requires only k dot products of d -dimensional vectors, followed by a sum of k n -dimensional vectors, and takes time $O(kd + kn)$, which is a win provided $k \ll \min(d, n)$. How is the error measured? Since \mathbf{x} is unknown, the approximation needs to be good for every \mathbf{x} . So we take the maximum over all \mathbf{x} of $|(A_k - A)\mathbf{x}|$. Since this would be infinite if $|\mathbf{x}|$ could grow without bound, we restrict the maximum to $|\mathbf{x}| \leq 1$. Formally, we define a new norm of a matrix A by

$$\|A\|_2 = \max_{|\mathbf{x}| \leq 1} |A\mathbf{x}|.$$

This is called the 2-norm or the spectral norm. Note that it equals $\sigma_1(A)$.

As an application consider a large database of documents that form rows of an $n \times d$ matrix A . There are d terms and each document is a d -dimensional vector with one component for each term, which is the number of occurrences of the term in the document. We are allowed to “preprocess” A . After the preprocessing, we receive queries. Each query \mathbf{x} is an d -dimensional vector which specifies how important each term is to the query. The desired answer is an n -dimensional vector which gives the similarity (dot product) of the query to each document in the database, namely $A\mathbf{x}$, the “matrix-vector” product. Query time is to be much less than preprocessing time, since the idea is that we need to answer many queries for the same database. There are many other applications where one performs many matrix vector products with the same matrix. This technique is applicable to these situations as well.

3.6 Left Singular Vectors

The left singular vectors are also pairwise orthogonal. Intuitively if \mathbf{u}_i and \mathbf{u}_j , $i < j$, were not orthogonal, one would suspect that the right singular vector \mathbf{v}_j had a component of \mathbf{v}_i

which would contradict that \mathbf{v}_i and \mathbf{v}_j were orthogonal. Let i be the smallest integer such that \mathbf{u}_i is not orthogonal to all other \mathbf{u}_j . Then to prove that \mathbf{u}_i and \mathbf{u}_j are orthogonal, we add a small component of \mathbf{v}_j to \mathbf{v}_i , normalize the result to be a unit vector

$$\mathbf{v}'_i = \frac{\mathbf{v}_i + \epsilon \mathbf{v}_j}{|\mathbf{v}_i + \epsilon \mathbf{v}_j|}$$

and show that $|A\mathbf{v}'_i| > |A\mathbf{v}_i|$, a contradiction.

Theorem 3.7 *The left singular vectors are pairwise orthogonal.*

Proof: Let i be the smallest integer such that \mathbf{u}_i is not orthogonal to some other \mathbf{u}_j . Without loss of generality assume that $\mathbf{u}_i^T \mathbf{u}_j = \delta > 0$. If $\mathbf{u}_i^T \mathbf{u}_j < 0$ then just replace \mathbf{u}_j with $-\mathbf{u}_j$. Clearly $j > i$ since i was selected to be the smallest such index. For $\epsilon > 0$, let

$$\mathbf{v}'_i = \frac{\mathbf{v}_i + \epsilon \mathbf{v}_j}{|\mathbf{v}_i + \epsilon \mathbf{v}_j|}.$$

Notice that \mathbf{v}'_i is a unit-length vector.

$$A\mathbf{v}'_i = \frac{\sigma_i \mathbf{u}_i + \epsilon \sigma_j \mathbf{u}_j}{\sqrt{1 + \epsilon^2}}$$

has length at least as large as its component along \mathbf{u}_i which is

$$\mathbf{u}_i^T \left(\frac{\sigma_i \mathbf{u}_i + \epsilon \sigma_j \mathbf{u}_j}{\sqrt{1 + \epsilon^2}} \right) > (\sigma_i + \epsilon \sigma_j \delta) \left(1 - \frac{\epsilon^2}{2} \right) > \sigma_i - \frac{\epsilon^2}{2} \sigma_i + \epsilon \sigma_j \delta - \frac{\epsilon^3}{2} \sigma_j \delta > \sigma_i,$$

for sufficiently small ϵ , a contradiction since $\mathbf{v}_i + \epsilon \mathbf{v}_j$ is orthogonal to $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{i-1}$ since $j > i$ and σ_i is defined to be the maximum of $|A\mathbf{v}|$ over such vectors. \blacksquare

Next we prove that A_k is the best rank k , 2-norm approximation to A . We first show that the square of the 2-norm of $A - A_k$ is the square of the $(k+1)^{st}$ singular value of A . This is essentially by definition of A_k ; that is, A_k represents the projections of the rows in A onto the space spanned by the top k singular vectors, and so $A - A_k$ is the remaining portion of those rows, whose top singular value will be σ_{k+1} .

Lemma 3.8 $\|A - A_k\|_2^2 = \sigma_{k+1}^2$.

Proof: Let $A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ be the singular value decomposition of A . Then $A_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ and $A - A_k = \sum_{i=k+1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$. Let \mathbf{v} be the top singular vector of $A - A_k$. Express \mathbf{v} as a

linear combination of $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$. That is, write $\mathbf{v} = \sum_{j=1}^r c_j \mathbf{v}_j$. Then

$$\begin{aligned} |(A - A_k)\mathbf{v}| &= \left| \sum_{i=k+1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T \sum_{j=1}^r c_j \mathbf{v}_j \right| = \left| \sum_{i=k+1}^r c_i \sigma_i \mathbf{u}_i \mathbf{v}_i^T \mathbf{v}_i \right| \\ &= \left| \sum_{i=k+1}^r c_i \sigma_i \mathbf{u}_i \right| = \sqrt{\sum_{i=k+1}^r c_i^2 \sigma_i^2}, \end{aligned}$$

since the \mathbf{u}_i are orthonormal. The \mathbf{v} maximizing this last quantity, subject to the constraint that $|\mathbf{v}|^2 = \sum_{i=1}^r c_i^2 = 1$, occurs when $c_{k+1} = 1$ and the rest of the c_i are zero. Thus, $\|A - A_k\|_2^2 = \sigma_{k+1}^2$ proving the lemma. ■

Finally, we prove that A_k is the best rank k , 2-norm approximation to A :

Theorem 3.9 *Let A be an $n \times d$ matrix. For any matrix B of rank at most k*

$$\|A - A_k\|_2 \leq \|A - B\|_2.$$

Proof: If A is of rank k or less, the theorem is obviously true since $\|A - A_k\|_2 = 0$. Assume that A is of rank greater than k . By Lemma 3.8, $\|A - A_k\|_2^2 = \sigma_{k+1}^2$. The null space of B , the set of vectors \mathbf{v} such that $B\mathbf{v} = 0$, has dimension at least $d - k$. Let $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k+1}$ be the first $k + 1$ singular vectors of A . By a dimension argument, it follows that there exists a $\mathbf{z} \neq 0$ in

$$\text{Null}(B) \cap \text{Span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k+1}\}.$$

Scale \mathbf{z} to be of length one.

$$\|A - B\|_2^2 \geq |(A - B)\mathbf{z}|^2.$$

Since $B\mathbf{z} = 0$,

$$\|A - B\|_2^2 \geq |A\mathbf{z}|^2.$$

Since \mathbf{z} is in the $\text{Span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k+1}\}$

$$|A\mathbf{z}|^2 = \left| \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^T \mathbf{z} \right|^2 = \sum_{i=1}^n \sigma_i^2 (\mathbf{v}_i^T \mathbf{z})^2 = \sum_{i=1}^{k+1} \sigma_i^2 (\mathbf{v}_i^T \mathbf{z})^2 \geq \sigma_{k+1}^2 \sum_{i=1}^{k+1} (\mathbf{v}_i^T \mathbf{z})^2 = \sigma_{k+1}^2.$$

It follows that $\|A - B\|_2^2 \geq \sigma_{k+1}^2$ proving the theorem. ■

For a square symmetric matrix A and eigenvector \mathbf{v} , $A\mathbf{v} = \lambda\mathbf{v}$. We now prove the analog for singular values and vectors we discussed in the introduction.

Lemma 3.10 (Analog of eigenvalues and eigenvectors)

$$A\mathbf{v}_i = \sigma_i \mathbf{u}_i \text{ and } A^T \mathbf{u}_i = \sigma_i \mathbf{v}_i.$$

Proof: The first equation follows from the definition of left singular vectors. For the second, note that from the SVD, we get $A^T \mathbf{u}_i = \sum_j \sigma_j \mathbf{v}_j \mathbf{u}_j^T \mathbf{u}_i$, where since the \mathbf{u}_j are orthonormal, all terms in the summation are zero except for $j = i$. ■

3.7 Power Method for Singular Value Decomposition

Computing the singular value decomposition is an important branch of numerical analysis in which there have been many sophisticated developments over a long period of time. The reader is referred to numerical analysis texts for more details. Here we present an “in-principle” method to establish that the approximate SVD of a matrix A can be computed in polynomial time. The method we present, called the *power method*, is simple and is in fact the conceptual starting point for many algorithms. Let A be a matrix whose SVD is $\sum_i \sigma_i \mathbf{u}_i \mathbf{v}_i^T$. We wish to work with a matrix that is square and symmetric. Let $B = A^T A$. By direct multiplication, using the orthogonality of the \mathbf{u}_i ’s that was proved in Theorem 3.7,

$$\begin{aligned} B &= A^T A = \left(\sum_i \sigma_i \mathbf{v}_i \mathbf{u}_i^T \right) \left(\sum_j \sigma_j \mathbf{u}_j \mathbf{v}_j^T \right) \\ &= \sum_{i,j} \sigma_i \sigma_j \mathbf{v}_i (\mathbf{u}_i^T \cdot \mathbf{u}_j) \mathbf{v}_j^T = \sum_i \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T. \end{aligned}$$

The matrix B is square and symmetric, and has the same left and right-singular vectors. In particular, $B \mathbf{v}_j = (\sum_i \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T) \mathbf{v}_j = \sigma_j^2 \mathbf{v}_j$, so \mathbf{v}_j is an eigenvector of B with eigenvalue σ_j^2 . If A is itself square and symmetric, it will have the same right and left-singular vectors, namely $A = \sum_i \sigma_i \mathbf{v}_i \mathbf{v}_i^T$ and computing B is unnecessary.

Now consider computing B^2 .

$$B^2 = \left(\sum_i \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T \right) \left(\sum_j \sigma_j^2 \mathbf{v}_j \mathbf{v}_j^T \right) = \sum_{ij} \sigma_i^2 \sigma_j^2 \mathbf{v}_i (\mathbf{v}_i^T \mathbf{v}_j) \mathbf{v}_j^T$$

When $i \neq j$, the dot product $\mathbf{v}_i^T \mathbf{v}_j$ is zero by orthogonality.⁹ Thus, $B^2 = \sum_{i=1}^r \sigma_i^4 \mathbf{v}_i \mathbf{v}_i^T$. In computing the k^{th} power of B , all the cross product terms are zero and

$$B^k = \sum_{i=1}^r \sigma_i^{2k} \mathbf{v}_i \mathbf{v}_i^T.$$

If $\sigma_1 > \sigma_2$, then the first term in the summation dominates, so $B^k \rightarrow \sigma_1^{2k} \mathbf{v}_1 \mathbf{v}_1^T$. This means a close estimate to \mathbf{v}_1 can be computed by simply taking the first column of B^k and normalizing it to a unit vector.

3.7.1 A Faster Method

A problem with the above method is that A may be a very large, sparse matrix, say a $10^8 \times 10^8$ matrix with 10^9 nonzero entries. Sparse matrices are often represented by just

⁹The “outer product” $\mathbf{v}_i \mathbf{v}_j^T$ is a matrix and is not zero even for $i \neq j$.

a list of nonzero entries, say a list of triples of the form (i, j, a_{ij}) . Though A is sparse, B need not be and in the worse case may have all 10^{16} entries nonzero¹⁰ and it is then impossible to even write down B , let alone compute the product B^2 . Even if A is moderate in size, computing matrix products is costly in time. Thus, a more efficient method is needed.

Instead of computing B^k , select a random vector \mathbf{x} and compute the product $B^k \mathbf{x}$. The vector \mathbf{x} can be expressed in terms of the singular vectors of B augmented to a full orthonormal basis as $\mathbf{x} = \sum_{i=1}^d c_i \mathbf{v}_i$. Then

$$B^k \mathbf{x} \approx (\sigma_1^{2k} \mathbf{v}_1 \mathbf{v}_1^T) \left(\sum_{i=1}^d c_i \mathbf{v}_i \right) = \sigma_1^{2k} c_1 \mathbf{v}_1.$$

Normalizing the resulting vector yields \mathbf{v}_1 , the first singular vector of A . The way $B^k \mathbf{x}$ is computed is by a series of matrix vector products, instead of matrix products. $B^k \mathbf{x} = A^T A \dots A^T A \mathbf{x}$, which can be computed right-to-left. This consists of $2k$ vector times sparse matrix multiplications.

To compute k singular vectors, one selects a random vector \mathbf{r} and finds an orthonormal basis for the space spanned by $\mathbf{r}, A\mathbf{r}, \dots, A^{k-1}\mathbf{r}$. Then compute A times each of the basis vectors, and find an orthonormal basis for the space spanned by the resulting vectors. Intuitively, one has applied A to a subspace rather than a single vector. One repeatedly applies A to the subspace, calculating an orthonormal basis after each application to prevent the subspace collapsing to the one dimensional subspace spanned by the first singular vector. The process quickly converges to the first k singular vectors.

An issue occurs if there is no significant gap between the first and second singular values of a matrix. Take for example the case when there is a tie for the first singular vector and $\sigma_1 = \sigma_2$. Then, the above argument fails. We will overcome this hurdle. Theorem 3.11 below states that even with ties, the power method converges to some vector in the span of those singular vectors corresponding to the “nearly highest” singular values. The theorem assumes it is given a vector \mathbf{x} which has a component of magnitude at least δ along the first right singular vector \mathbf{v}_1 of A . We will see in Lemma 3.12 that a random vector satisfies this condition with fairly high probability.

Theorem 3.11 *Let A be an $n \times d$ matrix and \mathbf{x} a unit length vector in \mathbf{R}^d with $|\mathbf{x}^T \mathbf{v}_1| \geq \delta$, where $\delta > 0$. Let V be the space spanned by the right singular vectors of A corresponding to singular values greater than $(1 - \varepsilon) \sigma_1$. Let \mathbf{w} be the unit vector after $k = \frac{\ln(1/\varepsilon\delta)}{2\varepsilon}$ iterations of the power method, namely,*

$$\mathbf{w} = \frac{(A^T A)^k \mathbf{x}}{\|(A^T A)^k \mathbf{x}\|}.$$

¹⁰E.g., suppose each entry in the first row of A is nonzero and the rest of A is zero.

Then \mathbf{w} has a component of at most ε perpendicular to V .

Proof: Let

$$A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$

be the SVD of A . If the rank of A is less than d , then for convenience complete $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r\}$ into an orthonormal basis $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d\}$ of d -space. Write \mathbf{x} in the basis of the \mathbf{v}_i 's as

$$\mathbf{x} = \sum_{i=1}^d c_i \mathbf{v}_i.$$

Since $(A^T A)^k = \sum_{i=1}^d \sigma_i^{2k} \mathbf{v}_i \mathbf{v}_i^T$, it follows that $(A^T A)^k \mathbf{x} = \sum_{i=1}^d \sigma_i^{2k} c_i \mathbf{v}_i$. By hypothesis, $|c_1| \geq \delta$.

Suppose that $\sigma_1, \sigma_2, \dots, \sigma_m$ are the singular values of A that are greater than or equal to $(1 - \varepsilon) \sigma_1$ and that $\sigma_{m+1}, \dots, \sigma_d$ are the singular values that are less than $(1 - \varepsilon) \sigma_1$. Now

$$|(A^T A)^k \mathbf{x}|^2 = \left| \sum_{i=1}^d \sigma_i^{2k} c_i \mathbf{v}_i \right|^2 = \sum_{i=1}^d \sigma_i^{4k} c_i^2 \geq \sigma_1^{4k} c_1^2 \geq \sigma_1^{4k} \delta^2.$$

The component of $|(A^T A)^k \mathbf{x}|^2$ perpendicular to the space V is

$$\sum_{i=m+1}^d \sigma_i^{4k} c_i^2 \leq (1 - \varepsilon)^{4k} \sigma_1^{4k} \sum_{i=m+1}^d c_i^2 \leq (1 - \varepsilon)^{4k} \sigma_1^{4k}$$

since $\sum_{i=1}^d c_i^2 = |\mathbf{x}|^2 = 1$. Thus, the component of \mathbf{w} perpendicular to V has squared length at most $\frac{(1-\varepsilon)^{4k} \sigma_1^{4k}}{\sigma_1^{4k} \delta^2}$ and so its length is at most

$$\frac{(1 - \varepsilon)^{2k} \sigma_1^{2k}}{\delta \sigma_1^{2k}} = \frac{(1 - \varepsilon)^{2k}}{\delta} \leq \frac{e^{-2k\varepsilon}}{\delta} = \varepsilon$$

since $k = \frac{\ln(1/\varepsilon\delta)}{2\varepsilon}$. ■

Lemma 3.12 *Let $\mathbf{y} \in \mathbf{R}^n$ be a random vector with the unit variance spherical Gaussian as its probability density. Normalize \mathbf{y} to be a unit length vector by setting $\mathbf{x} = \mathbf{y}/|\mathbf{y}|$. Let \mathbf{v} be any unit length vector. Then*

$$\text{Prob} \left(|\mathbf{x}^T \mathbf{v}| \leq \frac{1}{20\sqrt{d}} \right) \leq \frac{1}{10} + 3e^{-d/96}.$$

Proof: Proving for the unit length vector \mathbf{x} that $\text{Prob}\left(|\mathbf{x}^T \mathbf{v}| \leq \frac{1}{20\sqrt{d}}\right) \leq \frac{1}{10} + 3e^{-d/96}$ is equivalent to proving for the unnormalized vector \mathbf{y} that $\text{Prob}(|\mathbf{y}| \geq 2\sqrt{d}) \leq 3e^{-d/96}$ and $\text{Prob}(|\mathbf{y}^T \mathbf{v}| \leq \frac{1}{10}) \leq 1/10$. That $\text{Prob}(|\mathbf{y}| \geq 2\sqrt{d})$ is at most $3e^{-d/96}$ follows from Theorem (2.9) with \sqrt{d} substituted for β . The probability that $|\mathbf{y}^T \mathbf{v}| \leq \frac{1}{10}$ is at most $1/10$ follows from the fact that $\mathbf{y}^T \mathbf{v}$ is a random, zero mean, unit variance Gaussian with density is at most $1/\sqrt{2\pi} \leq 1/2$ in the interval $[-1/10, 1/10]$, so the integral of the Gaussian over the interval is at most $1/10$. ■

3.8 Singular Vectors and Eigenvectors

For a square matrix B , if $B\mathbf{x} = \lambda\mathbf{x}$, then \mathbf{x} is an *eigenvector* of B and λ is the corresponding *eigenvalue*. We saw in Section 3.7, if $B = A^T A$, then the right singular vectors \mathbf{v}_j of A are eigenvectors of B with eigenvalues σ_j^2 . The same argument shows that the left singular vectors \mathbf{u}_j of A are eigenvectors of AA^T with eigenvalues σ_j^2 .

The matrix $B = A^T A$ has the property that for any vector \mathbf{x} , $\mathbf{x}^T B \mathbf{x} \geq 0$. This is because $B = \sum_i \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T$ and for any \mathbf{x} , $\mathbf{x}^T \mathbf{v}_i \mathbf{v}_i^T \mathbf{x} = (\mathbf{x}^T \mathbf{v}_i)^2 \geq 0$. A matrix B with the property that $\mathbf{x}^T B \mathbf{x} \geq 0$ for all \mathbf{x} is called *positive semi-definite*. Every matrix of the form $A^T A$ is positive semi-definite. In the other direction, any positive semi-definite matrix B can be decomposed into a product $A^T A$, and so its eigenvalue decomposition can be obtained from the singular value decomposition of A . The interested reader should consult a linear algebra book.

3.9 Applications of Singular Value Decomposition

3.9.1 Centering Data

Singular value decomposition is used in many applications and for some of these applications it is essential to first center the data by subtracting the centroid of the data from each data point.¹¹ If you are interested in the statistics of the data and how it varies in relationship to its mean, then you would center the data. On the other hand, if you are interested in finding the best low rank approximation to a matrix, then you do not center the data. The issue is whether you are finding the best fitting subspace or the best fitting affine space. In the latter case you first center the data and then find the best fitting subspace. See Figure 3.3.

We first show that the line minimizing the sum of squared distances to a set of points, if not restricted to go through the origin, must pass through the centroid of the points. This implies that if the centroid is subtracted from each data point, such a line will pass through the origin. The best fit line can be generalized to k dimensional “planes”. The operation of subtracting the centroid from all data points is useful in other contexts as well. We give it the name “centering data”.

¹¹The centroid of a set of points is the coordinate-wise average of the points.

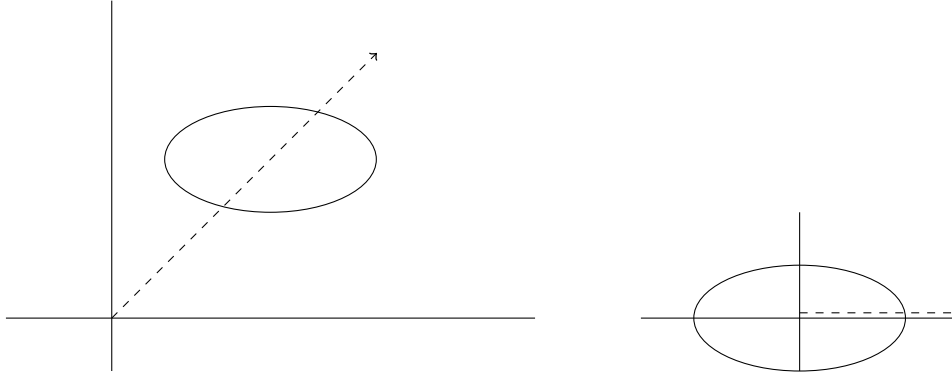


Figure 3.3: If one wants statistical information relative to the mean of the data, one needs to center the data. If one wants the best low rank approximation, one would not center the data.

Lemma 3.13 *The best-fit line (minimizing the sum of perpendicular distances squared) of a set of data points must pass through the centroid of the points.*

Proof: Subtract the centroid from each data point so that the centroid is $\mathbf{0}$. After centering the data let ℓ be the best-fit line and assume for contradiction that ℓ does not pass through the origin. The line ℓ can be written as $\{\mathbf{a} + \lambda \mathbf{v} | \lambda \in \mathbf{R}\}$, where \mathbf{a} is the closest point to $\mathbf{0}$ on ℓ and \mathbf{v} is a unit length vector in the direction of ℓ , which is perpendicular to \mathbf{a} . For a data point \mathbf{a}_i , let $\text{dist}(\mathbf{a}_i, \ell)$ denote its perpendicular distance to ℓ . By the Pythagorean theorem, we have $|\mathbf{a}_i - \mathbf{a}|^2 = \text{dist}(\mathbf{a}_i, \ell)^2 + (\mathbf{v} \cdot \mathbf{a}_i)^2$, or equivalently, $\text{dist}(\mathbf{a}_i, \ell)^2 = |\mathbf{a}_i - \mathbf{a}|^2 - (\mathbf{v} \cdot \mathbf{a}_i)^2$. Summing over all data points:

$$\begin{aligned} \sum_{i=1}^n \text{dist}(\mathbf{a}_i, \ell)^2 &= \sum_{i=1}^n (|\mathbf{a}_i - \mathbf{a}|^2 - (\mathbf{v} \cdot \mathbf{a}_i)^2) = \sum_{i=1}^n (|\mathbf{a}_i|^2 + |\mathbf{a}|^2 - 2\mathbf{a}_i \cdot \mathbf{a} - (\mathbf{v} \cdot \mathbf{a}_i)^2) \\ &= \sum_{i=1}^n |\mathbf{a}_i|^2 + n|\mathbf{a}|^2 - 2\mathbf{a} \cdot \left(\sum_i \mathbf{a}_i \right) - \sum_{i=1}^n (\mathbf{v} \cdot \mathbf{a}_i)^2 = \sum_i |\mathbf{a}_i|^2 + n|\mathbf{a}|^2 - \sum_i (\mathbf{v} \cdot \mathbf{a}_i)^2, \end{aligned}$$

where we used the fact that since the centroid is $\mathbf{0}$, $\sum_i \mathbf{a}_i = \mathbf{0}$. The above expression is minimized when $\mathbf{a} = \mathbf{0}$, so the line $\ell' = \{\lambda \mathbf{v} : \lambda \in \mathbf{R}\}$ through the origin is a better fit than ℓ , contradicting ℓ being the best-fit line. \blacksquare

A statement analogous to Lemma 3.13 holds for higher dimensional objects. Define an *affine space* as a subspace translated by a vector. So an affine space is a set of the form

$$\{\mathbf{v}_0 + \sum_{i=1}^k c_i \mathbf{v}_i | c_1, c_2, \dots, c_k \in \mathbf{R}\}.$$

Here, \mathbf{v}_0 is the translation and $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$ form an orthonormal basis for the subspace.

Lemma 3.14 *The k dimensional affine space which minimizes the sum of squared perpendicular distances to the data points must pass through the centroid of the points.*

Proof: We only give a brief idea of the proof, which is similar to the previous lemma. Instead of $(\mathbf{v} \cdot \mathbf{a}_i)^2$, we will now have $\sum_{j=1}^k (\mathbf{v}_j \cdot \mathbf{a}_i)^2$, where the $\mathbf{v}_j, j = 1, 2, \dots, k$ are an orthonormal basis of the subspace through the origin parallel to the affine space. ■

3.9.2 Principal Component Analysis

The traditional use of SVD is in Principal Component Analysis (PCA). PCA is illustrated by a movie recommendation setting where there are n customers and d movies. Let matrix A with elements a_{ij} represent the amount that customer i likes movie j . One hypothesizes that there are only k underlying basic factors that determine how much a given customer will like a given movie, where k is much smaller than n or d . For example, these could be the amount of comedy, drama, and action, the novelty of the story, etc. Each movie can be described as a k -dimensional vector indicating how much of these basic factors the movie has, and each customer can be described as a k -dimensional vector indicating how important each of these basic factors is to that customer. The dot-product of these two vectors is hypothesized to determine how much that customer will like that movie. In particular, this means that the $n \times d$ matrix A can be expressed as the product of an $n \times k$ matrix U describing the customers and a $k \times d$ matrix V describing the movies. Finding the best rank k approximation A_k by SVD gives such a U and V . One twist is that A may not be exactly equal to UV , in which case $A - UV$ is treated as noise. Another issue is that SVD gives a factorization with negative entries. Nonnegative matrix factorization (NMF) is more appropriate in some contexts where we want to keep entries nonnegative. NMF is discussed in Chapter 9

In the above setting, A was available fully and we wished to find U and V to identify the basic factors. However, in a case such as movie recommendations, each customer may have seen only a small fraction of the movies, so it may be more natural to assume that we are given just a few elements of A and wish to estimate A . If A was an arbitrary matrix of size $n \times d$, this would require $\Omega(nd)$ pieces of information and cannot be done with a few entries. But again hypothesize that A was a small rank matrix with added noise. If now we also assume that the given entries are randomly drawn according to some known distribution, then there is a possibility that SVD can be used to estimate the whole of A . This area is called collaborative filtering and one of its uses is to recommend movies or to target an ad to a customer based on one or two purchases. We do not describe it here.

3.9.3 Clustering a Mixture of Spherical Gaussians

Clustering is the task of partitioning a set of points into k subsets or clusters where each cluster consists of nearby points. Different definitions of the quality of a clustering lead to different solutions. Clustering is an important area which we will study in detail in Chapter 7. Here we will see how to solve a particular clustering problem using singular value decomposition.

$$\begin{array}{c} \text{customers} \end{array} \left(\begin{array}{c} \\ \\ \\ \end{array} \right) \begin{array}{c} A \end{array} = \begin{array}{c} \begin{array}{c} \text{factors} \\ \\ \\ \end{array} \end{array} \left(\begin{array}{c} \\ \\ \\ \end{array} \right) \begin{array}{c} U \end{array} \left(\begin{array}{c} \begin{array}{c} \text{movies} \\ \\ \\ \end{array} \\ \\ \\ \end{array} \right) \begin{array}{c} V \end{array} \left(\begin{array}{c} \\ \\ \\ \end{array} \right)$$

Figure 3.4: Customer-movie data

Mathematical formulations of clustering tend to have the property that finding the highest quality solution to a given set of data is NP-hard. One way around this is to assume stochastic models of input data and devise algorithms to cluster data generated by such models. Mixture models are a very important class of stochastic models. A mixture is a probability density or distribution that is the weighted sum of simple component probability densities. It is of the form

$$f = w_1 p_1 + w_2 p_2 + \cdots + w_k p_k,$$

where p_1, p_2, \dots, p_k are the basic probability densities and w_1, w_2, \dots, w_k are positive real numbers called mixture weights that add up to one. Clearly, f is a probability density and integrates to one.

The *model fitting problem* is to fit a mixture of k basic densities to n independent, identically distributed samples, each sample drawn according to the same mixture distribution f . The class of basic densities is known, but various parameters such as their means and the component weights of the mixture are not. Here, we deal with the case where the basic densities are all spherical Gaussians. There are two equivalent ways of thinking of the hidden sample generation process when only the samples are given:

1. Pick each sample according to the density f on \mathbf{R}^d .
2. Pick a random i from $\{1, 2, \dots, k\}$ where probability of picking i is w_i . Then, pick a sample according to the density p_i .

One approach to the model-fitting problem is to break it into two subproblems:

1. First, cluster the set of samples into k clusters C_1, C_2, \dots, C_k , where C_i is the set of samples generated according to p_i (see (2) above) by the hidden generation process.
2. Then fit a single Gaussian distribution to each cluster of sample points.

The second problem is relatively easier and indeed we saw the solution in Chapter 2, where we showed that taking the empirical mean (the mean of the sample) and the empirical standard deviation gives us the best-fit Gaussian. The first problem is harder and this is what we discuss here.

If the component Gaussians in the mixture have their centers very close together, then the clustering problem is unresolvable. In the limiting case where a pair of component densities are the same, there is no way to distinguish between them. What condition on the inter-center separation will guarantee unambiguous clustering? First, by looking at 1-dimensional examples, it is clear that this separation should be measured in units of the standard deviation, since the density is a function of the number of standard deviation from the mean. In one dimension, if two Gaussians have inter-center separation at least six times the maximum of their standard deviations, then they hardly overlap. This is summarized in the question: How many standard deviations apart are the means? In one dimension, if the answer is at least six, we can easily tell the Gaussians apart. What is the analog of this in higher dimensions?

We discussed in Chapter 2 distances between two sample points from the same Gaussian as well the distance between two sample points from two different Gaussians. Recall from that discussion that if

- If \mathbf{x} and \mathbf{y} are two independent samples from the same spherical Gaussian with standard deviation¹² σ then

$$|\mathbf{x} - \mathbf{y}|^2 \approx 2(\sqrt{d} \pm O(1))^2 \sigma^2.$$

- If \mathbf{x} and \mathbf{y} are samples from different spherical Gaussians each of standard deviation σ and means separated by distance Δ , then

$$|\mathbf{x} - \mathbf{y}|^2 \approx 2(\sqrt{d} \pm O(1))^2 \sigma^2 + \Delta^2.$$

To ensure that points from the same Gaussian are closer to each other than points from different Gaussians, we need

$$2(\sqrt{d} - O(1))^2 \sigma^2 + \Delta^2 > 2(\sqrt{d} + O(1))^2 \sigma^2.$$

Expanding the squares, the high order term $2d$ cancels and we need that

$$\Delta > cd^{1/4},$$

for some constant c . While this was not a completely rigorous argument, it can be used to show that a distance based clustering approach (see Chapter 2 for an example) requires an

¹²Since a spherical Gaussian has the same standard deviation in every direction, we call it the standard deviation of the Gaussian.

inter-mean separation of at least $cd^{1/4}$ standard deviations to succeed, thus unfortunately not keeping with mnemonic of a constant number of standard deviations separation of the means. Here, indeed, we will show that $\Omega(1)$ standard deviations suffice provided the number k of Gaussians is $O(1)$.

The central idea is the following. Suppose we can find the subspace spanned by the k centers and project the sample points to this subspace. The projection of a spherical Gaussian with standard deviation σ remains a spherical Gaussian with standard deviation σ (Lemma 3.15). In the projection, the inter-center separation remains the same. So in the projection, the Gaussians are distinct provided the inter-center separation in the whole space is at least $ck^{1/4}\sigma$ which is less than $cd^{1/4}\sigma$ for $k \ll d$. Interestingly, we will see that the subspace spanned by the k -centers is essentially the best-fit k -dimensional subspace that can be found by singular value decomposition.

Lemma 3.15 *Suppose p is a d -dimensional spherical Gaussian with center $\boldsymbol{\mu}$ and standard deviation σ . The density of p projected onto a k -dimensional subspace V is a spherical Gaussian with the same standard deviation.*

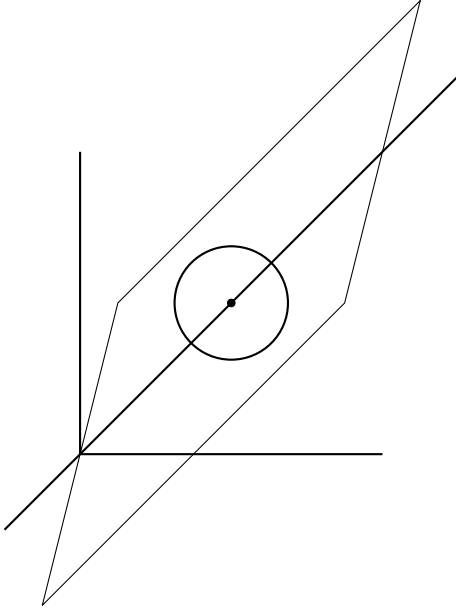
Proof: Rotate the coordinate system so V is spanned by the first k coordinate vectors. The Gaussian remains spherical with standard deviation σ although the coordinates of its center have changed. For a point $\mathbf{x} = (x_1, x_2, \dots, x_d)$, we will use the notation $\mathbf{x}' = (x_1, x_2, \dots, x_k)$ and $\mathbf{x}'' = (x_{k+1}, x_{k+2}, \dots, x_n)$. The density of the projected Gaussian at the point (x_1, x_2, \dots, x_k) is

$$ce^{-\frac{|\mathbf{x}' - \boldsymbol{\mu}'|^2}{2\sigma^2}} \int_{\mathbf{x}''} e^{-\frac{|\mathbf{x}'' - \boldsymbol{\mu}''|^2}{2\sigma^2}} d\mathbf{x}'' = c'e^{-\frac{|\mathbf{x}' - \boldsymbol{\mu}'|^2}{2\sigma^2}}.$$

This implies the lemma. ■

We now show that the top k singular vectors produced by the SVD span the space of the k centers. First, we extend the notion of best fit to probability distributions. Then we show that for a single spherical Gaussian whose center is not the origin, the best fit 1-dimensional subspace is the line through the center of the Gaussian and the origin. Next, we show that the best fit k -dimensional subspace for a single Gaussian whose center is not the origin is any k -dimensional subspace containing the line through the Gaussian's center and the origin. Finally, for k spherical Gaussians, the best fit k -dimensional subspace is the subspace containing their centers. Thus, the SVD finds the subspace that contains the centers.

Recall that for a set of points, the best-fit line is the line passing through the origin that maximizes the sum of squared lengths of the projections of the points onto the line. We extend this definition to probability densities instead of a set of points.



1. The best fit 1-dimension subspace to a spherical Gaussian is the line through its center and the origin.
2. Any k -dimensional subspace containing the line is a best fit k -dimensional subspace for the Gaussian.
3. The best fit k -dimensional subspace for k spherical Gaussians is the subspace containing their centers.

Figure 3.5: Best fit subspace to a spherical Gaussian.

Definition 3.1 If p is a probability density in d space, the best fit line for p is the line in the \mathbf{v}_1 direction where

$$\mathbf{v}_1 = \arg \max_{|\mathbf{v}|=1} E \left[(\mathbf{v}^T \mathbf{x})^2 \right].$$

■

For a spherical Gaussian centered at the origin, it is easy to see that any line passing through the origin is a best fit line. Our next lemma shows that the best fit line for a spherical Gaussian centered at $\boldsymbol{\mu} \neq 0$ is the line passing through $\boldsymbol{\mu}$ and the origin.

Lemma 3.16 Let the probability density p be a spherical Gaussian with center $\boldsymbol{\mu} \neq 0$. The unique best fit 1-dimensional subspace is the line passing through $\boldsymbol{\mu}$ and the origin. If $\boldsymbol{\mu} = 0$, then any line through the origin is a best-fit line.

Proof: For a randomly chosen \mathbf{x} (according to p) and a fixed unit length vector \mathbf{v} ,

$$\begin{aligned} E_{\mathbf{x} \sim p} [(\mathbf{v}^T \mathbf{x})^2] &= E_{\mathbf{x} \sim p} \left[(\mathbf{v}^T (\mathbf{x} - \boldsymbol{\mu}) + \mathbf{v}^T \boldsymbol{\mu})^2 \right] \\ &= E_{\mathbf{x} \sim p} \left[(\mathbf{v}^T (\mathbf{x} - \boldsymbol{\mu}))^2 + 2 (\mathbf{v}^T \boldsymbol{\mu}) (\mathbf{v}^T (\mathbf{x} - \boldsymbol{\mu})) + (\mathbf{v}^T \boldsymbol{\mu})^2 \right] \\ &= E_{\mathbf{x} \sim p} \left[(\mathbf{v}^T (\mathbf{x} - \boldsymbol{\mu}))^2 \right] + 2 (\mathbf{v}^T \boldsymbol{\mu}) E_{\mathbf{x} \sim p} [\mathbf{v}^T (\mathbf{x} - \boldsymbol{\mu})] + (\mathbf{v}^T \boldsymbol{\mu})^2 \\ &= E_{\mathbf{x} \sim p} \left[(\mathbf{v}^T (\mathbf{x} - \boldsymbol{\mu}))^2 \right] + (\mathbf{v}^T \boldsymbol{\mu})^2 \\ &= \sigma^2 + (\mathbf{v}^T \boldsymbol{\mu})^2 \end{aligned}$$

where the fourth line follows from the fact that $E[\mathbf{v}^T(\mathbf{x} - \boldsymbol{\mu})] = 0$, and the fifth line follows from the fact that $E[(\mathbf{v}^T(\mathbf{x} - \boldsymbol{\mu}))^2]$ is the variance in the direction \mathbf{v} . The best fit line \mathbf{v} maximizes $E_{\mathbf{x} \sim p}[(\mathbf{v}^T \mathbf{x})^2]$ and therefore maximizes $(\mathbf{v}^T \boldsymbol{\mu})^2$. This is maximized when \mathbf{v} is aligned with the center $\boldsymbol{\mu}$. To see uniqueness, just note that if $\boldsymbol{\mu} \neq 0$, then $\mathbf{v}^T \boldsymbol{\mu}$ is strictly less when \mathbf{v} is not aligned with the center. ■

We now extend Definition 3.1 to k -dimensional subspaces.

Definition 3.2 *If p is a probability density in d -space then the best-fit k -dimensional subspace V_k is*

$$V_k = \underset{\substack{V \\ \dim(V)=k}}{\operatorname{argmax}} E_{\mathbf{x} \sim p} (|\operatorname{proj}(\mathbf{x}, V)|^2),$$

where $\operatorname{proj}(\mathbf{x}, V)$ is the orthogonal projection of \mathbf{x} onto V . ■

Lemma 3.17 *For a spherical Gaussian with center $\boldsymbol{\mu}$, a k -dimensional subspace is a best fit subspace if and only if it contains $\boldsymbol{\mu}$.*

Proof: If $\boldsymbol{\mu} = \mathbf{0}$, then by symmetry any k -dimensional subspace is a best-fit subspace. If $\boldsymbol{\mu} \neq \mathbf{0}$, then, the best-fit line must pass through $\boldsymbol{\mu}$ by Lemma 3.16. Now, as in the greedy algorithm for finding subsequent singular vectors, we would project perpendicular to the first singular vector. But after the projection, the mean of the Gaussian becomes $\mathbf{0}$ and any vectors will do as subsequent best-fit directions. ■

This leads to the following theorem.

Theorem 3.18 *If p is a mixture of k spherical Gaussians, then the best fit k -dimensional subspace contains the centers. In particular, if the means of the Gaussians are linearly independent, the space spanned by them is the unique best-fit k dimensional subspace.*

Proof: Let p be the mixture $w_1 p_1 + w_2 p_2 + \dots + w_k p_k$. Let V be any subspace of dimension k or less. Then,

$$E_{\mathbf{x} \sim p} (|\operatorname{proj}(\mathbf{x}, V)|^2) = \sum_{i=1}^k w_i E_{\mathbf{x} \sim p_i} (|\operatorname{proj}(\mathbf{x}, V)|^2)$$

If V contains the centers of the densities p_i , by Lemma 3.17, each term in the summation is individually maximized, which implies the entire summation is maximized, proving the theorem. ■

For an infinite set of points drawn according to the mixture, the k -dimensional SVD subspace gives exactly the space of the centers. In reality, we have only a large number of samples drawn according to the mixture. However, it is intuitively clear that as the number of samples increases, the set of sample points will approximate the probability density and so the SVD subspace of the sample will be close to the space spanned by the centers. The details of how close it gets as a function of the number of samples are technical and we do not carry this out here.

3.9.4 Ranking Documents and Web Pages

An important task for a document collection is to rank the documents according to their intrinsic relevance to the collection. A good candidate definition of “intrinsic relevance” is a document’s projection onto the best-fit direction for that collection, namely the top left-singular vector of the term-document matrix. An intuitive reason for this is that this direction has the maximum sum of squared projections of the collection and so can be thought of as a synthetic term-document vector best representing the document collection.

Ranking in order of the projection of each document’s term vector along the best fit direction has a nice interpretation in terms of the power method. For this, we consider a different example, that of the web with hypertext links. The World Wide Web can be represented by a directed graph whose nodes correspond to web pages and directed edges to hypertext links between pages. Some web pages, called *authorities*, are the most prominent sources for information on a given topic. Other pages called *hubs*, are ones that identify the authorities on a topic. Authority pages are pointed to by many hub pages and hub pages point to many authorities. One is led to what seems like a circular definition: a hub is a page that points to many authorities and an authority is a page that is pointed to by many hubs.

One would like to assign hub weights and authority weights to each node of the web. If there are n nodes, the hub weights form an n -dimensional vector \mathbf{u} and the authority weights form an n -dimensional vector \mathbf{v} . Suppose A is the adjacency matrix representing the directed graph. Here a_{ij} is 1 if there is a hypertext link from page i to page j and 0 otherwise. Given hub vector \mathbf{u} , the authority vector \mathbf{v} could be computed by the formula

$$v_j \propto \sum_{i=1}^d u_i a_{ij}$$

since the right hand side is the sum of the hub weights of all the nodes that point to node j . In matrix terms,

$$\mathbf{v} = A^T \mathbf{u} / |A^T \mathbf{u}|.$$

Similarly, given an authority vector \mathbf{v} , the hub vector \mathbf{u} could be computed by $\mathbf{u} = A\mathbf{v} / |A\mathbf{v}|$. Of course, at the start, we have neither vector. But the above discussion suggests a power iteration. Start with any \mathbf{v} . Set $\mathbf{u} = A\mathbf{v}$, then set $\mathbf{v} = A^T \mathbf{u}$, then renormalize and repeat the process. We know from the power method that this converges to the left and right-singular vectors. So after sufficiently many iterations, we may use the left vector \mathbf{u} as the hub weights vector and project each column of A onto this direction and rank columns (authorities) in order of this projection. But the projections just form the vector $A^T \mathbf{u}$ which equals a multiple of \mathbf{v} . So we can just rank by order of the v_j . This is the basis of an algorithm called the HITS algorithm, which was one of the early proposals for ranking web pages.

A different ranking called *pagerank* is widely used. It is based on a random walk on the graph described above. We will study random walks in detail in Chapter 4.

3.9.5 An Application of SVD to a Discrete Optimization Problem

In clustering a mixture of Gaussians, SVD was used as a dimension reduction technique. It found a k -dimensional subspace (the space of centers) of a d -dimensional space and made the Gaussian clustering problem easier by projecting the data to the subspace. Here, instead of fitting a model to data, we consider an optimization problem where applying dimension reduction makes the problem easier. The use of SVD to solve discrete optimization problems is a relatively new subject with many applications. We start with an important NP-hard problem, the maximum cut problem for a directed graph $G(V, E)$.

The maximum cut problem is to partition the nodes of an n -node directed graph into two subsets S and \bar{S} so that the number of edges from S to \bar{S} is maximized. Let A be the adjacency matrix of the graph. With each vertex i , associate an indicator variable x_i . The variable x_i will be set to 1 for $i \in S$ and 0 for $i \in \bar{S}$. The vector $\mathbf{x} = (x_1, x_2, \dots, x_n)$ is unknown and we are trying to find it or equivalently the cut, so as to maximize the number of edges across the cut. The number of edges across the cut is precisely

$$\sum_{i,j} x_i(1 - x_j)a_{ij}.$$

Thus, the maximum cut problem can be posed as the optimization problem

$$\text{Maximize } \sum_{i,j} x_i(1 - x_j)a_{ij} \quad \text{subject to } x_i \in \{0, 1\}.$$

In matrix notation,

$$\sum_{i,j} x_i(1 - x_j)a_{ij} = \mathbf{x}^T A(\mathbf{1} - \mathbf{x}),$$

where $\mathbf{1}$ denotes the vector of all 1's. So, the problem can be restated as

$$\text{Maximize } \mathbf{x}^T A(\mathbf{1} - \mathbf{x}) \quad \text{subject to } x_i \in \{0, 1\}. \quad (3.1)$$

This problem is NP-hard. However we will see that for dense graphs, that is, graphs with $\Omega(n^2)$ edges and therefore whose optimal solution has size $\Omega(n^2)$,¹³ we can use the SVD to find a near optimal solution in polynomial time. To do so we will begin by computing the SVD of A and replacing A by $A_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ in (3.1) to get

$$\text{Maximize } \mathbf{x}^T A_k(\mathbf{1} - \mathbf{x}) \quad \text{subject to } x_i \in \{0, 1\}. \quad (3.2)$$

Note that the matrix A_k is no longer a 0-1 adjacency matrix.

We will show that:

¹³Any graph of m edges has a cut of size at least $m/2$. This can be seen by noting that the expected size of the cut for a *random* $\mathbf{x} \in \{0, 1\}^n$ is exactly $m/2$.

1. For each 0-1 vector \mathbf{x} , $\mathbf{x}^T A_k (\mathbf{1} - \mathbf{x})$ and $\mathbf{x}^T A (\mathbf{1} - \mathbf{x})$ differ by at most $\frac{n^2}{\sqrt{k+1}}$. Thus, the maxima in (3.1) and (3.2) differ by at most this amount.
2. A near optimal \mathbf{x} for (3.2) can be found in time $n^{O(k)}$ by exploiting the low rank of A_k , which is polynomial time for constant k . By Item 1 this is near optimal for (3.1) where near optimal means with additive error of at most $\frac{n^2}{\sqrt{k+1}}$.

First, we prove Item 1. Since \mathbf{x} and $\mathbf{1} - \mathbf{x}$ are 0-1 n -vectors, each has length at most \sqrt{n} . By the definition of the 2-norm, $|(A - A_k)(\mathbf{1} - \mathbf{x})| \leq \sqrt{n} \|A - A_k\|_2$. Now since $\mathbf{x}^T (A - A_k)(\mathbf{1} - \mathbf{x})$ is the dot product of the vector \mathbf{x} with the vector $(A - A_k)(\mathbf{1} - \mathbf{x})$,

$$|\mathbf{x}^T (A - A_k)(\mathbf{1} - \mathbf{x})| \leq n \|A - A_k\|_2.$$

By Lemma 3.8, $\|A - A_k\|_2 = \sigma_{k+1}(A)$. The inequalities,

$$(k+1)\sigma_{k+1}^2 \leq \sigma_1^2 + \sigma_2^2 + \dots + \sigma_{k+1}^2 \leq \|A\|_F^2 = \sum_{i,j} a_{ij}^2 \leq n^2$$

imply that $\sigma_{k+1}^2 \leq \frac{n^2}{k+1}$ and hence $\|A - A_k\|_2 \leq \frac{n}{\sqrt{k+1}}$ proving Item 1.

Next we focus on Item 2. It is instructive to look at the special case when $k=1$ and A is approximated by the rank one matrix A_1 . An even more special case when the left and right-singular vectors \mathbf{u} and \mathbf{v} are identical is already NP-hard to solve exactly because it subsumes the problem of whether for a set of n integers, $\{a_1, a_2, \dots, a_n\}$, there is a partition into two subsets whose sums are equal. However, for that problem, there is an efficient dynamic programming algorithm that finds a near-optimal solution. We will build on that idea for the general rank k problem.

For Item 2, we want to maximize $\sum_{i=1}^k \sigma_i (\mathbf{x}^T \mathbf{u}_i) (\mathbf{v}_i^T (\mathbf{1} - \mathbf{x}))$ over 0-1 vectors \mathbf{x} . A piece of notation will be useful. For any $S \subseteq \{1, 2, \dots, n\}$, write $\mathbf{u}_i(S)$ for the sum of coordinates of the vector \mathbf{u}_i corresponding to elements in the set S , that is, $\mathbf{u}_i(S) = \sum_{j \in S} u_{ij}$, and similarly for \mathbf{v}_i . We will find S to maximize $\sum_{i=1}^k \sigma_i \mathbf{u}_i(S) \mathbf{v}_i(\bar{S})$ using dynamic programming.

For a subset S of $\{1, 2, \dots, n\}$, define the $2k$ -dimensional vector

$$\mathbf{w}(S) = (\mathbf{u}_1(S), \mathbf{v}_1(\bar{S}), \mathbf{u}_2(S), \mathbf{v}_2(\bar{S}), \dots, \mathbf{u}_k(S), \mathbf{v}_k(\bar{S})).$$

If we had the list of all such vectors, we could find $\sum_{i=1}^k \sigma_i \mathbf{u}_i(S) \mathbf{v}_i(\bar{S})$ for each of them and take the maximum. There are 2^n subsets S , but several S could have the same $\mathbf{w}(S)$ and in that case it suffices to list just one of them. Round each coordinate of each \mathbf{u}_i to the nearest integer multiple of $\frac{1}{nk^2}$. Call the rounded vector $\tilde{\mathbf{u}}_i$. Similarly obtain $\tilde{\mathbf{v}}_i$. Let $\tilde{\mathbf{w}}(S)$ denote the vector $(\tilde{\mathbf{u}}_1(S), \tilde{\mathbf{v}}_1(\bar{S}), \tilde{\mathbf{u}}_2(S), \tilde{\mathbf{v}}_2(\bar{S}), \dots, \tilde{\mathbf{u}}_k(S), \tilde{\mathbf{v}}_k(\bar{S}))$. We will construct a list of all possible values of the vector $\tilde{\mathbf{w}}(S)$. Again, if several different S 's lead to the

same vector $\tilde{\mathbf{w}}(S)$, we will keep only one copy on the list. The list will be constructed by dynamic programming. For the recursive step, assume we already have a list of all such vectors for $S \subseteq \{1, 2, \dots, i\}$ and wish to construct the list for $S \subseteq \{1, 2, \dots, i+1\}$. Each $S \subseteq \{1, 2, \dots, i\}$ leads to two possible $S' \subseteq \{1, 2, \dots, i+1\}$, namely, S and $S \cup \{i+1\}$. In the first case, the vector $\tilde{\mathbf{w}}(S') = (\tilde{\mathbf{u}}_1(S), \tilde{\mathbf{v}}_1(\bar{S}) + \tilde{v}_{1,i+1}, \tilde{\mathbf{u}}_2(S), \tilde{\mathbf{v}}_2(\bar{S}) + \tilde{v}_{2,i+1}, \dots, \dots)$. In the second case, it is $\tilde{\mathbf{w}}(S') = (\tilde{\mathbf{u}}_1(S) + \tilde{u}_{1,i+1}, \tilde{\mathbf{v}}_1(\bar{S}), \tilde{\mathbf{u}}_2(S) + \tilde{u}_{2,i+1}, \tilde{\mathbf{v}}_2(\bar{S}), \dots, \dots)$. We put in these two vectors for each vector in the previous list. Then, crucially, we prune - i.e., eliminate duplicates.

Assume that k is constant. Now, we show that the error is at most $\frac{n^2}{\sqrt{k+1}}$ as claimed. Since \mathbf{u}_i and \mathbf{v}_i are unit length vectors, $|\mathbf{u}_i(S)|, |\mathbf{v}_i(\bar{S})| \leq \sqrt{n}$. Also $|\tilde{\mathbf{u}}_i(S) - \mathbf{u}_i(S)| \leq \frac{n}{nk^2} = \frac{1}{k^2}$ and similarly for \mathbf{v}_i . To bound the error, we use an elementary fact: if a and b are reals with $|a|, |b| \leq M$ and we estimate a by a' and b by b' so that $|a - a'|, |b - b'| \leq \delta \leq M$, then $a'b'$ is an estimate of ab in the sense

$$|ab - a'b'| = |a(b - b') + b'(a - a')| \leq |a||b - b'| + (|b| + |b - b'|)|a - a'| \leq 3M\delta.$$

Using this,

$$\left| \sum_{i=1}^k \sigma_i \tilde{\mathbf{u}}_i(S) \tilde{\mathbf{v}}_i(\bar{S}) - \sum_{i=1}^k \sigma_i \mathbf{u}_i(S) \mathbf{v}_i(\bar{S}) \right| \leq 3k\sigma_1 \sqrt{n}/k^2 \leq 3n^{3/2}/k \leq n^2/k,$$

and this meets the claimed error bound.

Next, we show that the running time is polynomially bounded. First, $|\tilde{\mathbf{u}}_i(S)|, |\tilde{\mathbf{v}}_i(S)| \leq 2\sqrt{n}$. Since $\tilde{\mathbf{u}}_i(S)$ and $\tilde{\mathbf{v}}_i(S)$ are all integer multiples of $1/(nk^2)$, there are at most $2n^{3/2}k^2$ possible values of $\tilde{\mathbf{u}}_i(S)$ and $\tilde{\mathbf{v}}_i(S)$ from which it follows that the list of $\tilde{\mathbf{w}}(S)$ never gets larger than $(2n^{3/2}k^2)^{2k}$ which for fixed k is polynomially bounded.

We summarize what we have accomplished.

Theorem 3.19 *Given a directed graph $G(V, E)$, a cut of size at least the maximum cut minus $O\left(\frac{n^2}{\sqrt{k}}\right)$ can be computed in time polynomial in n for any fixed k .*

Note that achieving the same accuracy in time polynomial in n and k would give an exact max cut in polynomial time.

3.10 Bibliographic Notes

Singular value decomposition is fundamental to numerical analysis and linear algebra. There are many texts on these subjects and the interested reader may want to study these. A good reference is [GvL96]. The material on clustering a mixture of Gaussians in Section 3.9.3 is from [VW02]. Modeling data with a mixture of Gaussians is a standard tool in statistics. Several well-known heuristics like the expectation-minimization