

Applications

- Compression, matrix completion
- Recommender systems
- Sports team rankings

Topics

- Low-rank approximations
- Matrix completion problem
- Frobenius norm
- SVD for optimal low-rank approximations

These lecture notes are mostly based off lecture 9 from Stanford CS168:

<https://web.stanford.edu/class/cs168/l/l9.pdf>

What are the Missing Entries?

Suppose that I run a web streaming service for movies for three of my friends, Amy, Bob, and Carol. It's a very specialized movie service, with only five movie options: The Matrix, Inception, Star Wars: Episode I, Moana, and Inside Out. After 1 month, we ask our friends Amy, Bob, and Carol to rate the movies they've watched from one to five. We collect their ratings into a table below (we mark unwatched movies with ?):

	The Matrix	Inception	Star Wars: Ep I	Moana	Inside Out
Amy	2	?	?	?	5
Bob	?	3	4	?	2
Carol	?	?	2	1	?

and are asked to provide recommendations to Amy, Bob, and Carol as to which movie they should watch next. Said another way, we are asked to fill in the unknown ? entries in the table above.

This seems a bit unfair! Each of the unknown entries could be any value in 1-5 after all! But what if I told you an additional hint: Amy, Bob, and Carol have the same relative preferences for each movie: for example, Amy likes Inside Out $\frac{5}{2}$ more than Bob likes Inside Out, and this ratio is the same across all movies. Mathematically, we are making the assumption that **all columns of the table above are multiples of each other**.

Thus we can conclude that Bob likes The Matrix $\frac{2}{5}$. (Amy's rating) = $\frac{4}{5}$. Similarly, Carol's rating of Inception is $\frac{1}{2}$ (Bob's rating) = $\frac{1.5}{5}$, Carol's rating of Inside Out is $\frac{1}{2} \times (\text{Bob's rating}) = 1$, and so on. Here's the completed matrix:

$$M = \begin{bmatrix} 2 & 7.5 & 10 & 5 & 5 \\ 0.8 & 3 & 4 & 2 & 2 \\ 0.4 & 1.5 & 2 & 1 & 1 \end{bmatrix}$$

The point of this example is that when you know something about the **structure** of a partially known matrix, then sometimes it is possible to intelligently fill in missing entries. In this previous example, the assumption that every column is a multiple of each other means that $\text{rank } M = 1$ ($\dim \text{Col}(M) = 1$), which is pretty extreme! One natural and useful definition is that assuming a matrix M has **low-rank**. What rank counts as "low" is application dependent, but it typically means that for a matrix $M \in \mathbb{R}^{m \times n}$, that $\text{rank } M = r \ll \min\{m, n\}$.

This lecture will explore how we can use this idea of structure to solve the matrix completion problem by finding the best low-rank approximation to a partially known matrix. The SVD will of course be our main tool.

Low-Rank Matrix Approximations: Motivation

Before diving into the math, let's highlight some applications of low-rank matrix approximations:

1. Compression: we saw this idea last class, but it's worth revisiting through the lens of low-rank approximations. If the original matrix $M \in \mathbb{R}^{mn}$ is described by mn numbers, then a rank k approximation only requires $k(mn)$ numbers. To see this, recall that if \hat{M} has rank k , then we can write its SVD as

$$\begin{aligned}\hat{M} &= \begin{bmatrix} U \\ \vdots \\ U \end{bmatrix}_{m \times k} \begin{bmatrix} Z \end{bmatrix}_{k \times k} \begin{bmatrix} V^T \end{bmatrix}_{k \times n} \\ &= \begin{bmatrix} U \Sigma^{1/2} \\ \vdots \\ U \Sigma^{1/2} = Y \end{bmatrix}_{m \times k} \begin{bmatrix} \Sigma^{1/2} V^T \\ \vdots \\ \Sigma^{1/2} V^T = Z^T \end{bmatrix}_{k \times n} \quad (\Sigma^{1/2} = \text{diag}(\sigma_1^{1/2}, \dots, \sigma_k^{1/2}))\end{aligned}$$

a product $\hat{M} = YZ^T$, where $Y \in \mathbb{R}^{m \times k}$ and $Z \in \mathbb{R}^{n \times k}$. For example, if M represents a Gray Scale image (with entries = pixel intensities), m and n are typically in the 100s (or 1000s for HD images), and a modest value of k ($\sim 100-150$) is usually enough to give a good approximation of the original image.

2. Updating Huge AI Models: A modern application of low-rank matrix approximations is for "fine-tuning" huge AI models. In the setting of Large Language Models (LLMs), like (but GPT), we are typically given some huge off-the-shelf model with billions (or more) parameters. Given this large model that has been trained on an enormous but generic corpus of text from the web, one often performs "fine-tuning". This fine-tuning is a second round of training, typically using a much smaller domain specific dataset (for example, the lecture notes for this class could be used to fine-tune a "LinearAlgebra(GPT)"). The challenge of fine-tuning is that because these models are so big, making these updates is extremely challenging. The 2021 paper "LOLA: Low-Rank Adaptation of Large Language Models" argued that fine-tuning updates are generally approx. low-rank and that one can learn these updates in their factorized YZ^T form, allowing model fine-tuning with $1000 \times 1000 \times$ fewer parameters!

3. Denoising: If M is a noisy version of some "true" matrix that is approx. low-rank, then finding a low-rank approximation to M will typically remove a lot of noise (and maybe some signal), resulting in a matrix that is actually more informative than the original.

4. Matrix Completion: Low-rank approximations offer a way of solving the matrix completion problem we introduced above. Given a matrix M with missing entries, the first step is to obtain a full matrix \hat{M} by filling in the missing entries with "default" values.

What these default values should be often requires trial and error, but natural things to try include \bar{O} , the average of known entries in the same column, row, or the entire matrix. The second step is then to find a rank k approximation to \tilde{M} . This approach works well when the unknown matrix is close to a rank k matrix and there aren't too many missing entries.

With this motivation in mind, let's see how the SVD can help us in finding a good rank r approximation of a matrix M . Once we've described our procedure, and seen some examples of it in action, we'll make precise how our method is actually producing the "best" rank r approximation possible.

Low-Rank Approximations from the SVD

Given an $m \times n$ matrix $M \in \mathbb{R}^{m \times n}$, which we'll assume has rank r . Then the SVD of M is given by

$$M = U \Sigma V^T = \sum_{i=1}^r \sigma_i u_i v_i^T, \quad (\text{SVD})$$

for $U = [u_1 \dots u_r] \in \mathbb{R}^{m \times r}$, $V = [v_1 \dots v_r] \in \mathbb{R}^{n \times r}$, and $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r)$ the matrices of left singular vectors, right singular vectors, and singular values, respectively.

The right-most expression of (SVD) is a particularly convenient expression for our purposes, which expresses M as a sum of rank 1 matrices $\sigma_i u_i v_i^T$ with mutually orthogonal column and row spaces.

This sum expression suggests a very natural way of forming a rank k approximation to M : simply truncate the sum to the top k terms, as measured by the singular values σ_i :

$$\tilde{M}_k = \sum_{i=1}^k \sigma_i u_i v_i^T = U_k \Sigma_k V_k^T, \quad (\text{SVD-}k)$$

where the right-most expression is defined in terms of the truncated matrices:

$$U_k = [u_1 \dots u_k] \in \mathbb{R}^{m \times k}, \quad V_k = [v_1 \dots v_k] \in \mathbb{R}^{n \times k}, \quad \Sigma_k = \text{diag}(\sigma_1, \dots, \sigma_k) \in \mathbb{R}^{k \times k}$$

Before analyzing the properties of $\tilde{M}_k = U_k \Sigma_k V_k^T$, let's examine if \tilde{M}_k could plausibly address our motivating application. Storing the matrices U_k , V_k , and Σ_k requires storing $km + kn + k^2 \approx k(m+n)$ numbers if $k < \min\{m, n\}$, which is much less than the mn numbers needed to store $M \in \mathbb{R}^{m \times n}$ when m and n are relatively large.

It is also natural to interpret (SVD- k) as approximating the raw data M in terms of k "concepts" (e.g., "sci-fi", "romcom", "drama", "classic"), where the singular values $\sigma_1, \dots, \sigma_k$ express the "prominence" of the concepts, the rows of V^T and columns of U express the "typical row(column)" associated with each concept (e.g., a viewer likes only sci-fi movies, or a movie liked only by romcom viewers), and the rows of U (or columns of V^T) approximately express each row

(or column) of M as a linear combination (scaled by g_1, \dots, g_r) of the "typical rows" (or "typical columns"),

this method of producing a low-rank approximation is beautiful: we interpret the SVD of a matrix M as a list of "ingredients" ordered by "importance", and we retain only the k most important ingredients. But is this elegant procedure only "good"?

A Matrix Norm

For an $m \times n$ matrix $M \in \mathbb{R}^{m \times n}$, let \tilde{M} be a low-rank approximation of M , and define the approximation error as $E = M - \tilde{M}$. Intuitively, a "good" approximation will lead to "small" error E . But we need to quantify the "size" of $E \in \mathbb{R}^{m \times n}$: we saw that for vectors $x \in \mathbb{R}^n$, the right way to quantify the size of x was through its norm $\|x\|$, where $\|\cdot\|$ is a function that needs to satisfy the axioms of a norm.

1. $\|ax\| = |a|\|x\|$ for all $x \in \mathbb{R}^n$ and $a \in \mathbb{R}$
2. $\|x\| \geq 0$ for all $x \in \mathbb{R}^n$, and $\|x\| = 0$ if and only if $x = 0$
3. $\|x+y\| \leq \|x\| + \|y\|$ for all $x, y \in \mathbb{R}^n$.

It turns out we can define functions on the vector space of $m \times n$ matrices that satisfy these same properties: these are called **matrix norms**. We'll introduce one of them here that is particularly relevant to low-rank matrix approximations, but be aware that just as for vectors, there are many kinds of matrix norms.

The Frobenius Norm

The **Frobenius norm** of a $m \times n$ matrix $M \in \mathbb{R}^{m \times n}$ simply computes the Euclidean norm of M as if it were a mn vector:

$$\|M\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n m_{ij}^2}. \quad (F)$$

Example $M = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$ has $\|M\|_F^2 = 1^2 + 2^2 + 3^2 + 4^2 = 30$.

which is the same as $\|\text{Vec}(M)\|_2^2 = \left\| \begin{bmatrix} 1 \\ 3 \\ 2 \\ 4 \end{bmatrix} \right\|_2^2$.

We need a couple of properties of the Frobenius norm before we can connect the SVD to low-rank matrix approximation.

Property 1: For $A \in \mathbb{R}^{m \times n}$ a square matrix, $\|A\|_F = \|A^T\|_F$.

This isn't too hard to check from the definition of $\|\cdot\|_F$: taking the transpose just swaps the role (i,j) in the sum, but you still end up adding together the square of all entries in A , which are the same as the square of all of the entries in A^T .

Property 2: If $Q \in \mathbb{R}^{m \times m}$ is an orthogonal matrix and $A \in \mathbb{R}^{m \times n}$ is a square matrix, then $\|QA\|_F = \|A^TQ^T\|_F = \|A\|_F$, i.e., the Frobenius norm of a matrix A is unchanged by left or right multiplication by an orthogonal matrix.

To see why this is true, recall that if $A = [\underline{a}_1 \cdots \underline{a}_n]$ are the columns of A , then $QA = [Q\underline{a}_1 \cdots Q\underline{a}_n]$. Then, since we can write the Frobenius norm squared of a matrix as the sum of the Euclidean norm squared of its columns, we have:

$$\begin{aligned}\|QA\|_F^2 &= \|Q\underline{a}_1\|^2 + \cdots + \|Q\underline{a}_n\|^2 \\ &= \|\underline{a}_1\|^2 + \cdots + \|\underline{a}_n\|^2 = \|A\|_F^2.\end{aligned}$$

Here, the second equality holds because multiplying a vector by an orthogonal matrix does not change its Euclidean norm. Finally we use this and Property 1 to conclude

$$\|A^TQ^T\|_F = \|Q^TA\|_F = \|A^T\|_F = \|A\|_F.$$

since Q^T is
also orthogonal

We will measure the quality of our rank k approximation $(\text{SVD-}k) \hat{M}$ to M in terms of the Frobenius norm of their difference.

The following theorem tells us that the SVD-based approximation $(\text{SVD-}k)$ is optimal with respect to the Frobenius norm of the approximation error!

Theorem: For every $m \times n$ matrix $M \in \mathbb{R}^{m \times n}$, every rank target $k \geq 1$, and every rank k $m \times n$ matrix $B \in \mathbb{R}^{m \times n}$,

$$\|M - \hat{M}_k\|_F \leq \|M - B\|_F$$

where \hat{M}_k is the rank k approximation derived from the SVD $M = U \Sigma V^T$ as in $(\text{SVD-}k)$.

We won't formally prove this theorem, but let's get some intuition as to why this is true. To keep things simple, we'll assume M is square and full rank, i.e., $M \in \mathbb{R}^{n \times n}$ with rank $M=n$. Nearly the exact same argument works for general M , but we have to use

the non-compact SVD of M (which keeps zero singular values around).

Our goal is to find a rank k matrix \hat{M} which minimizes $\|\hat{M} - M\|_F^2$. Let $M = U \Sigma V^\top$ be the SVD of M , where $U, \Sigma, V \in \mathbb{R}^{n \times n}$ since $\text{rank } M = n$. By Property 2 of the Frobenius norm, we then have the following sequence of equalities:

$$\begin{aligned}\|\hat{M} - M\|_F^2 &= \|\hat{M} - U \Sigma V^\top\|_F^2 \\ &= \|U^\top (\hat{M} - U \Sigma V^\top)\|_F^2 \quad (\|A\|_F = \|QA\|_F) \\ &= \|U^\top \hat{M} - \Sigma V^\top\|_F^2 \quad (U^\top U = I) \\ &= \|(U^\top \hat{M} - \Sigma V^\top) V\|_F^2 \quad (\|A\|_F = \|AQ\|_F) \\ &= \|U^\top \hat{M} V - \Sigma\|_F^2 \quad (V^\top V = I)\end{aligned}$$

Now notice that since Σ is a diagonal matrix, any non-diagonal entry in $U^\top \hat{M} V$ adds to our approx error, so $U^\top \hat{M} V$ should be diagonal. Let's posit $\hat{M} = U D V^\top$ for some diagonal matrix D . Then

$$\|\hat{M} - M\|_F^2 = \|U^\top (U D V^\top) V - \Sigma\|_F^2 = \|D - \Sigma\|_F^2 = \sum_{i=1}^n (d_{ii} - \sigma_i)^2. \quad (*)$$

Therefore, we want to pick the diagonal entries d_{ii} of D to minimize the right-most expression in $(*)$. If there were no rank restriction on \hat{M} , we simply would set $d_{ii} = \sigma_i$. However, notice $\hat{M} = U D V^\top$ is an SVD of \hat{M} ! Therefore, for \hat{M} to be rank k , only k of the d_{ii} can be nonzero: if we can only knock off k of the $(d_{ii} - \sigma_i)^2$ terms in $(*)$, we should pick the top k , i.e., $d_{ii} = \sigma_i$ for $i=1, \dots, k$, and $d_{ii} = 0$ for $i=k+1, \dots, n$.

$$\text{Then, } \hat{M} = [U_1 \dots U_k \underbrace{U_{k+1} \dots U_n}_G] \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_k & \\ & & & 0 \end{bmatrix} \begin{bmatrix} V_1^\top \\ \vdots \\ V_k^\top \\ \vdots \\ V_{k+1}^\top \\ \vdots \\ V_n^\top \end{bmatrix} = \sum_{i=1}^k \sigma_i U_i V_i^\top = U_k \Sigma_k V_k^\top$$

is exactly the expression in (SVD-1c), and the approximation error it means is

$$\|\hat{M} - M\|_F^2 = \|\hat{M} - M\|_F^2 = \sum_{i=k+1}^n \sigma_i^2,$$

i.e. the sum of the squares of the "tail" singular values of M .

Example: Recall the matrix $A = \begin{bmatrix} 4 & 11 & 14 \\ 8 & 7 & -2 \end{bmatrix}$ from Lecture 18; we computed its SVD as:

$$A = \begin{bmatrix} 3/\sqrt{10} & 1/\sqrt{10} \\ 1/\sqrt{10} & -3/\sqrt{10} \end{bmatrix} \begin{bmatrix} 6\sqrt{10} & 0 \\ 0 & 3\sqrt{10} \end{bmatrix} \begin{bmatrix} 1/3 & 2/3 & 2/3 \\ -2/3 & -1/3 & 2/3 \end{bmatrix} = U \Sigma V^\top.$$

A is rank 2, and its rank 1 approximation is, according to (SVD-k), given by

$$\hat{A}_1 = \begin{bmatrix} 3\sqrt{10} \\ 1/\sqrt{10} \end{bmatrix} 6\sqrt{10} \begin{bmatrix} 1/3 & 2/3 & 2/3 \end{bmatrix} = \begin{bmatrix} 6 & 12 & 12 \\ 2 & 4 & 4 \end{bmatrix}$$

If we compute $\|\hat{A}_1 - A\|_F^2$ we get:

$$\left\| \begin{bmatrix} 2 & -1 & 2 \\ -6 & -3 & 6 \end{bmatrix} \right\|_F^2 = 2^2 + (-1)^2 + 2^2 + (-6)^2 + (-3)^2 + 6^2 = 90$$

which is exactly $6^2 = (3\sqrt{10})^2 = 90$.

Finally, we address an obvious question when applying these ideas in practice: how should we pick the rank k of our approximation?

In a perfect world, the singular values of the original data matrix will give strong guidance: if the top few singular values are much larger than the rest, then the obvious solution is to take $k = \# \text{ of } b_1 \text{ values}$. This was the case in the Lander example last class!

The 1st singular value was significantly larger than others, suggesting a rank 1 approximation would be a good choice (this was the image (d)).

In less clear settings, the rule of thumb is to take k as small as possible while still providing a "useful" approximation of the original data. For example, it is common to choose k so that the sum of the top k singular values is at least c times larger than the sum of the other singular values. The ratio c is typically a domain-dependent constant picked based on the application.