

CS-GY 6763: Lecture 11

Linear Programming, Singular Value Decomposition

NYU Tandon School of Engineering, Prof. Christopher Musco

DIMENSION DEPENDENT CONVEX OPTIMIZATION

Consider a convex function $f(\mathbf{x})$ be bounded between $[-B, B]$ on a constraint set \mathcal{S} .

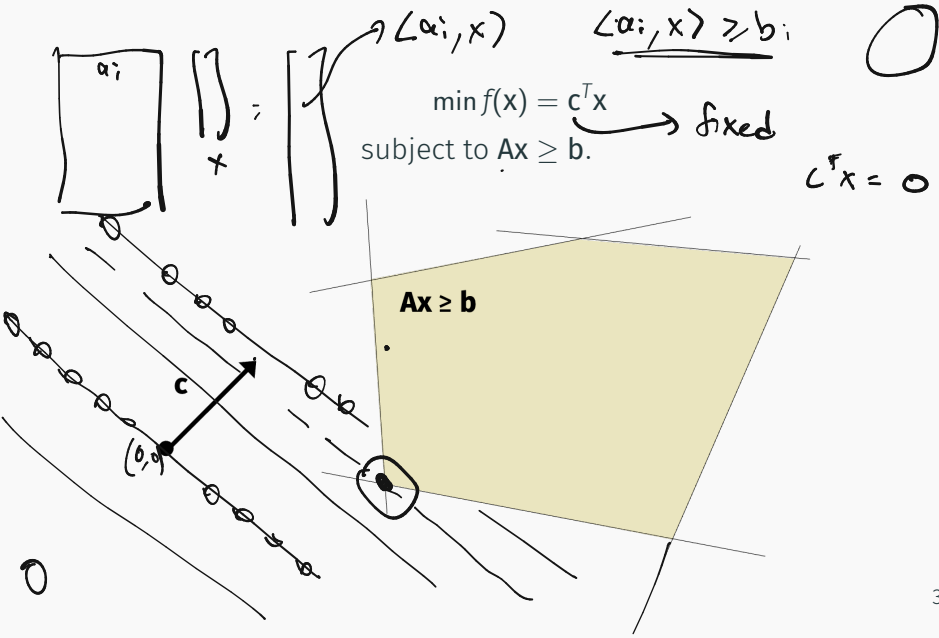
Theorem (Dimension Dependent Convex Optimization)

The Center-of-Gravity Method finds $\hat{\mathbf{x}}$ satisfying $f(\hat{\mathbf{x}}) \leq \min_{\mathbf{x} \in \mathcal{S}} f(\mathbf{x}) + \epsilon$ using $O(d \log(B/\epsilon))$ calls to a function and gradient oracle for convex f .

The center-of-gravity method is not computationally efficient, but inspired the polynomial time ellipsoid method.

$$\text{poly}(d) \quad O(d^7)$$

KILLER APPLICATION: LINEAR PROGRAMMING



Linear programs (LPs) are one of the most basic convex constrained, convex optimization problems:

Let $\mathbf{c} \in \mathbb{R}^d$, $\mathbf{b} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{n \times d}$ be fixed vectors that define the problem, and let \mathbf{x} be our variable parameter.

$$\begin{aligned} \min f(\mathbf{x}) &= \mathbf{c}^T \mathbf{x} \\ \text{subject to } \mathbf{A}\mathbf{x} &\geq \mathbf{b}. \end{aligned}$$

Think about $\mathbf{A}\mathbf{x} \geq \mathbf{b}$ as a union of half-space constraints:

$$\begin{aligned} \{\mathbf{x} : \mathbf{a}_1^T \mathbf{x} &\geq b_1\} \\ \{\mathbf{x} : \mathbf{a}_2^T \mathbf{x} &\geq b_2\} \\ &\vdots \\ \{\mathbf{x} : \mathbf{a}_n^T \mathbf{x} &\geq b_n\} \end{aligned} \quad \bigg\}$$

LINEAR PROGRAMMING APPLICATIONS

- Classic optimization applications: industrial resource optimization problems were important original applications in the 70s. }
- Robust regression: $\min_x \|Ax - b\|_1$.
- $L1$ constrained regression: $\min_x \|x\|_1$ subject to $Ax = b$. Lots of applications in sparse recovery/compressed sensing.
- Solve $\min_x \|Ax - b\|_\infty$.
- Polynomial time algorithms for (Markov Decision Processes) (reinforcement learning).
- Many combinatorial optimization problems can be solved via LP relaxation. }

LINEAR PROGRAMMING

Theorem (Khachiyan, 1979)

Assume $n = d$. The ellipsoid method solves any linear program with L -bit integer valued constraints exactly in $O(n^4 L)$ time.

d : dimension of x
 n : # of constraints

A Soviet Discovery Rocks World of Mathematics

By MALCOLM W. BROWNE

A surprise discovery by an obscure Soviet mathematician has rocked the world of mathematics and computer analysis, and experts have begun exploring its practical applications.

Mathematicians describe the discovery by L.G. Khachian as a method by which computers can find guaranteed solutions to a class of very difficult problems that have hitherto been tackled on a kind of hit-or-miss basis.

Apart from its profound theoretical interest, the discovery may be applicable

in weather prediction, complicated industrial processes, petroleum refining, the scheduling of workers at large factories, secret codes and many other things.

"I have been deluged with calls from virtually every department of government for an interpretation of the significance of this," a leading expert on computer methods, Dr. George B. Dantzig of Stanford University, said in an interview.

The solution of mathematical problems by computer must be broken down into a series of steps. One class of problem sometimes involves so many steps that it

could take billions of years to compute.

The Russian discovery offers a way by which the number of steps in a solution can be dramatically reduced. It also offers the mathematician a way of learning quickly whether a problem has a solution or not, without having to complete the entire immense computation that may be required.

According to the American journal Sci-

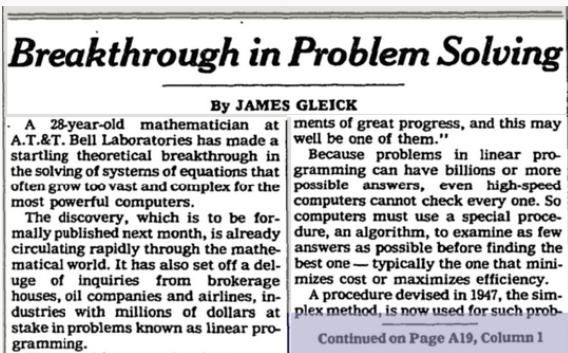
Continued on Page A20, Column 3

ONLY \$10.00 A MONTH!!! 24 Hr. Phone Answering Service. Totally New Concept!!! Incredible!!! 279-3870—ADVT.

Front page of New York Times, November 9, 1979.

Theorem (Karmarkar, 1984)

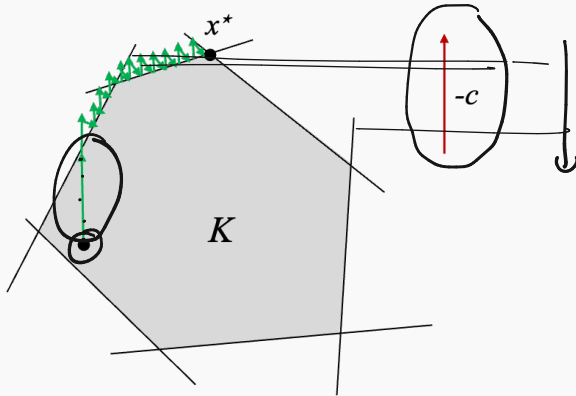
Assume $n = d$. The (interior point method) solves any linear program with L -bit integer valued constraints in $O(n^{3.5}L)$ time.



Front page of New York Times, November 19, 1984.

INTERIOR POINT METHODS

Lecture notes are posted on the website (optional reading).



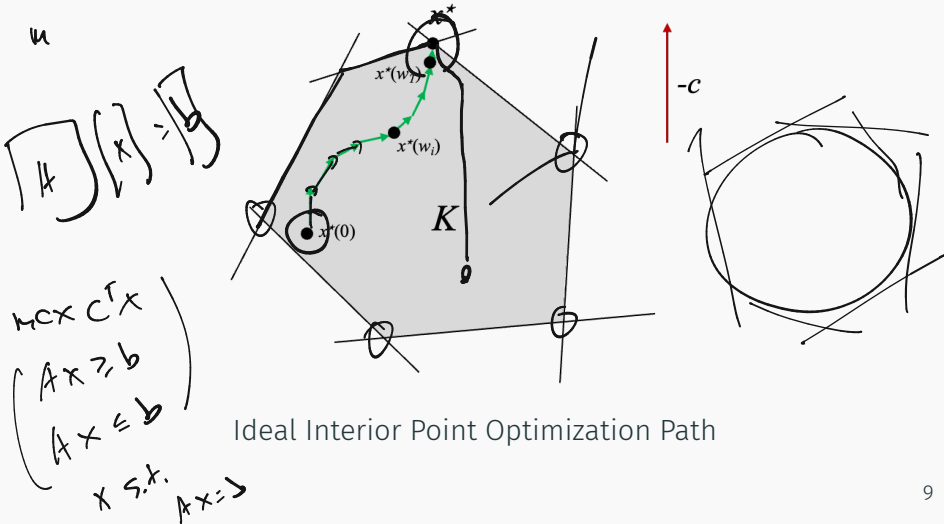
$$f(x) = c^T x$$
$$\nabla f(x) = c$$

Projected Gradient Descent Optimization Path

GD step is $- \mu c$.

INTERIOR POINT METHODS

Lecture notes are posted on the website (optional reading).



POLYNOMIAL TIME LINEAR PROGRAMMING

Both results had a huge impact on the theory of optimization, although at the time neither the ellipsoid method or interior point method were faster than a heuristic known at the Simplex Method.

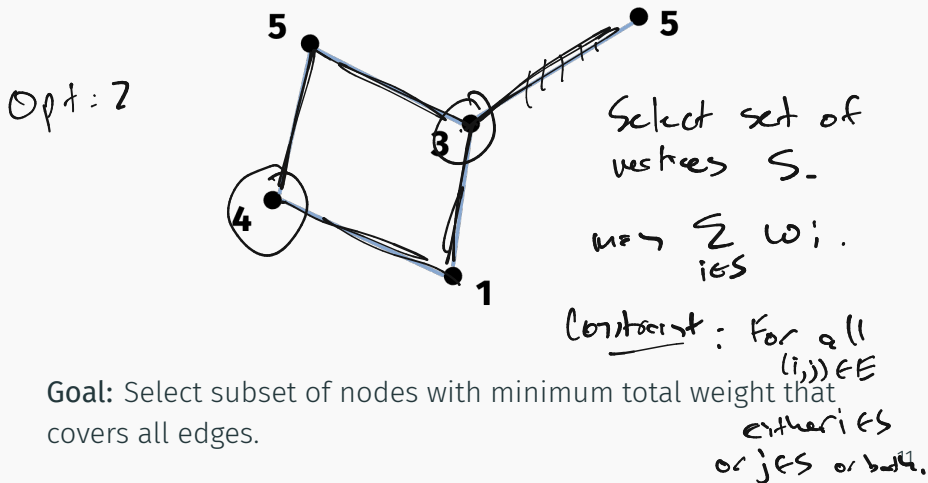
These days, improved interior point methods compete with and often outperform simplex.

Polynomial time linear programming algorithms have also had a huge impact of combinatorial optimization. They are often the work-horse behind approximation algorithms for NP-hard problems.

25

EXAMPLE: VERTEX COVER

Given a graph G with n nodes and edge set E . Each node is assigned a weight w_1, \dots, w_n .



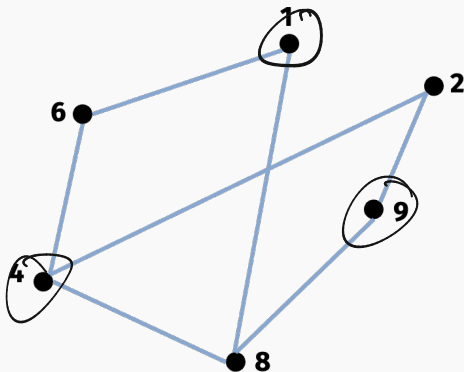
Goal: Select subset of nodes with minimum total weight that covers all edges.

EXAMPLE: VERTEX COVER

NP-hard to solve exactly.

$$\text{OPT} \leq 15$$

$$\text{OPT} = 14$$



EXAMPLE: VERTEX COVER

Given a graph G with n nodes and edge set E . Each node is assigned a weight w_1, \dots, w_n .

Formally: Denote if node i is selected by assigning variable x_i to 0 or 1. Let $x = [x_1, \dots, x_n]$.

$$\left(\min_x \sum_{i=1}^n x_i w_i \right) \text{ subject to } \left(\begin{array}{l} x_i \in \{0, 1\} \text{ for all } i \\ x_i + x_j \geq 1 \text{ for all } (i, j) \in E \end{array} \right)$$

\downarrow
 $= \sum_{i \in S} w_i$

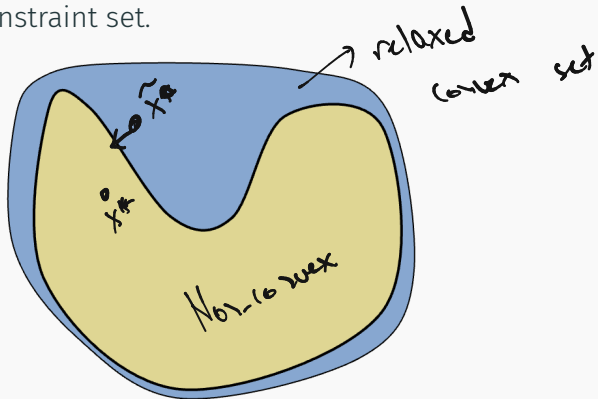
We will use convex optimization give a 2-approximation in polynomial time.

Function to minimize is linear (so convex) but constraint set is not convex. Why?

$$\begin{array}{l} [0, 1, 0] \\ [1, 0, 1] \end{array} = \left(\frac{1}{2}, \frac{1}{2}, \frac{1}{2} \right)$$

High level approach:

- { Relax to a problem with convex constraints.
- { Round optimal solution of convex problem back to original constraint set.



RELAX-AND-ROUND

High level approach:

- Relax to a problem with convex constraints.
- Round optimal solution of convex problem back to original constraint set.

$(0,1)$ $(1,1)$

$(0,0)$

$(1,0)$

$$0 \leq x_1 \leq 1$$

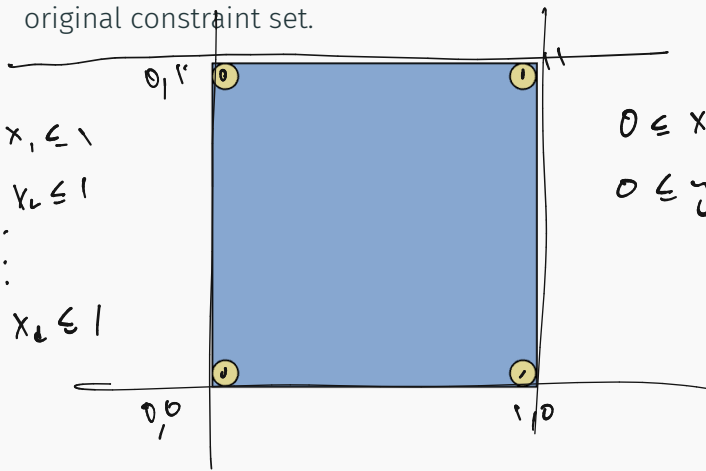
$$0 \leq x_2 \leq 1$$

\vdots

$$0 \leq x_n \leq 1$$

$$0 \leq x \leq 1$$

$$0 \leq y \leq 1$$



RELAX-AND-ROUND

High level approach: *original constraints*

- Relax to a problem with convex constraints.
- Round optimal solution of convex problem back to original constraint set.

Let $\bar{\mathcal{S}} \supseteq \mathcal{S}$ be the relaxed constraint set. Let $\underline{x}^* = \arg \min_{x \in \mathcal{S}} f(x)$ and let $\underline{\bar{x}}^* = \arg \min_{x \in \bar{\mathcal{S}}} f(x)$. We always have that:

$$\underline{f(\bar{x}^*)} \leq f(x^*).$$

So typically the goal is to round \bar{x}^* to \mathcal{S} in such a way that we don't increase the function value too much.

$$f(\text{round}(\bar{x}^*)) \leq 2 f(\bar{x}^*) \leq 2 \cdot f(x^*)$$

RELAXING VERTEX COVER

Vertex Cover:

$$\min_x \sum_{i=1}^n x_i w_i \quad \text{subject to}$$

$$x_i \in \{0, 1\} \text{ for all } i$$

$$x_i + x_j \geq 1 \text{ for all } (i, j) \in E$$

Relaxed Vertex Cover:

$$\min_x \sum_{i=1}^n x_i w_i \quad \text{subject to}$$

$$0 \leq x_i \leq 1 \text{ for all } i$$

$$x_i + x_j \geq 1 \text{ for all } (i, j) \in E$$

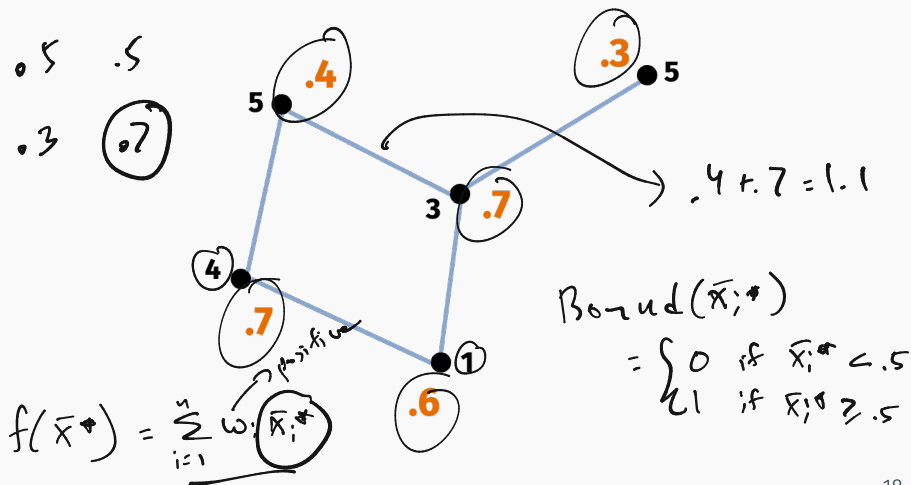
The second problem is a linear program! It can be solved in poly(n) time!

$$\begin{array}{l} \{0 \ 1 \ 0 \ 0 \ 0\} \rightarrow 0 \\ \{0 \ -1 \ 0 \ 0 \ 0\} \rightarrow -1 \\ \{0 \ 1 \ 0 \ 1 \ 0\} \rightarrow 1 \end{array}$$

ROUNDING VERTEX COVER

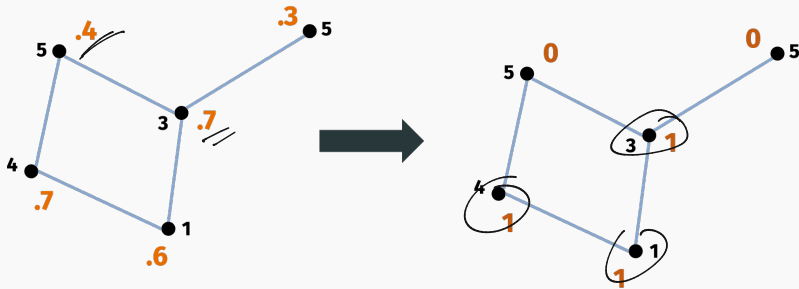
$$5 \cdot .4 + 5 \cdot .3 + 3 \cdot .7 \dots$$

Any ideas on how to round this to a solution to the original problem? I.e., with constraints $x_i \in \{0, 1\}$ for all i .



ROUNDING VERTEX COVER

Simply set all variable $x_i = 1$ if $\bar{x}_i^* \geq 1/2$ and $x_i = 0$ otherwise.



Observation 1: All edges remain covered. I.e., the constraint $x_i + x_j \geq 1$ for all $(i,j) \in E$ is not violated.

ROUNDING VERTEX COVER

Observation 2: Let \underline{x} be the rounded version of \bar{x}^* . We have $f(x) \leq 2 \cdot f(\bar{x}^*)$, and thus $f(x) \leq 2 \cdot f(x^*)$.

Proof:

$$\begin{aligned} f(x) &= \sum_{i=1}^n x_i \cdot w_i = \sum_{i=1}^n \text{round}(\bar{x}_i) \cdot w_i \\ &\leq \sum_{i=1}^n 2 \bar{x}_i \cdot w_i = 2 \sum_{i=1}^n \bar{x}_i \cdot w_i \\ &= 2 \cdot f(\bar{x}^*) \end{aligned}$$

$$f(x) \leq 2 f(\bar{x}^*) \leq 2 \cdot f(x^*) .$$

VERTEX COVER

So, a polynomial time algorithm for solving LPs immediately yields a 2-approximation algorithm for the NP-hard problem of vertex cover.

- Proven that it is NP-hard to do better than a 1.36' approximation in [Dinur, Safra, 2002].
- Recently improved to $\sqrt{2} \approx 1.41$ in [Khot, Minzer, Safra 2018], which proved the 2-to-2 games conjecture.
- Widely believed that doing better than $2 - \epsilon$ is NP-hard for any $\epsilon > 0$, and this is implied by Subhash Khot's (Unique Games Conjecture.)

$$2 - 1/\log(n)$$

There is a simpler greedy 2-approximation algorithm that doesn't use optimization at all!

BREAK

Next section of course: Spectral methods and numerical linear algebra.

Spectral methods generally refer to methods based on the “spectrum” of a matrix. I.e. on its eigenvectors/eigenvalues and singular vectors/singular values. We will look at

- Applications to low-rank approximation and dimensionality reduction.
- Applications to graph problems.
- Fast algorithms for computing spectral information.

Reminder: A vector $\underline{\mathbf{v}} \in \mathbb{R}^d$ is an (eigenvector) of a matrix $\underline{\mathbf{X}} \in \mathbb{R}^{d \times d}$, if there exists a scalar λ such that

$$\underline{\mathbf{X}}\underline{\mathbf{v}} = \underline{\lambda}\underline{\mathbf{v}}$$

The scalar (λ) is called the eigenvalue associated with \mathbf{v} .

Matrices can often be written completely in terms of their eigenvectors and eigenvalues. This is called eigendecomposition.)

We will actually focus on a related tool called (singular value decomposition.)

LINEAR ALGEBRA REMINDER

If a square matrix has orthonormal rows, it also has orthonormal columns:

$$V^T V = I \quad \longleftrightarrow \quad V V^T = I$$

$V^T V = I = V V^T$
 $[V V^T]_{ij} = \langle v_i, v_j \rangle$
 $\langle v_i, v_j \rangle = \|v_i\|^2 = 1$ if $i = j$
 $= 0$ otherwise

$$\begin{bmatrix} -0.62 & 0.78 & -0.11 \\ -0.28 & -0.35 & -0.89 \\ -0.73 & -0.52 & 0.44 \end{bmatrix} \cdot \begin{bmatrix} -0.62 & -0.28 & -0.73 \\ 0.78 & -0.35 & -0.52 \\ -0.11 & -0.89 & 0.44 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

LINEAR ALGEBRA REMINDER

Implies that for any vector x , $\|Vx\|_2^2 = \|x\|_2^2$ and $\|V^T x\|_2^2$.

Same thing goes for Frobenius norm: for any matrix X ,

$$\|VX\|_F^2 = \|X\|_F^2$$

Suppose V has orthonormal columns, then

$$\|Vx\|_2^2 = \|x\|_2^2$$

$$(Vx)^T (Vx) = x^T \underbrace{V^T V}_{=I} x = x^T x = \|x\|_2^2 \rightarrow c_1, \dots, c_d$$

$$\|x\|_2^2 = \sum_{i=1}^d \|c_i\|_2^2$$

$c_i = i$ th column of x

$$\begin{bmatrix} | & & | \\ V & & X \\ | & & | \end{bmatrix} = Vc_1, Vc_2, \dots, Vc_d$$

$$\|VX\|_F^2 = \sum_{i=1}^d \|Vc_i\|_2^2$$

LINEAR ALGEBRA REMINDER

The same is not true for rectangular matrices.

The diagram illustrates the properties of rectangular matrices. On the left, a horizontal teal box labeled V^T is multiplied by a vertical teal box labeled V (which contains four vertical lines representing columns). The result is a white box containing a 4x4 identity matrix I with ones on the diagonal. On the right, a vertical teal box labeled V is multiplied by a horizontal teal box labeled V^T . The result is a white box containing a 4x4 matrix:

.5	-1	.7	-2
1.6	-.44	4.2	-1.5
7.8	.42	-.5	.67
-2	2.0	1.1	8.0
-1.5	.55	3.2	.5
.67	-2.8	-2.4	1.6
9.0	8.7	-7.7	7.8

$$\underline{V^T V} = \underline{I}$$

but

$$V V^T \neq I$$

For any \mathbf{x} , $\|\mathbf{V}\mathbf{x}\|_2^2 = \|\mathbf{x}\|_2^2$ but $\|\mathbf{V}^T\mathbf{x}\|_2^2 \neq \|\mathbf{x}\|_2^2$ in general.

LINEAR ALGEBRA REMINDER

Multiplying a vector by V with orthonormal columns rotates
and/or reflects the vector.

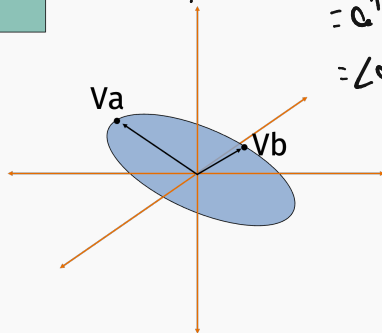
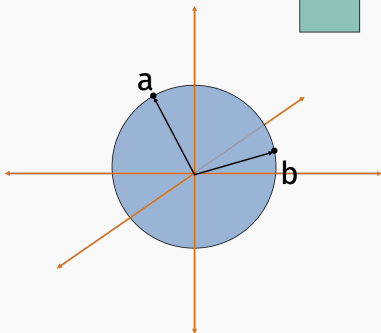
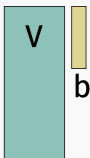
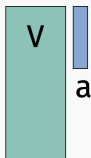
$$\|Va - Vb\|_2^2 = \|V(a-b)\|_2^2$$

$$= \|a-b\|_2^2$$

$$\langle Va, Vb \rangle = a^T V^T V b$$

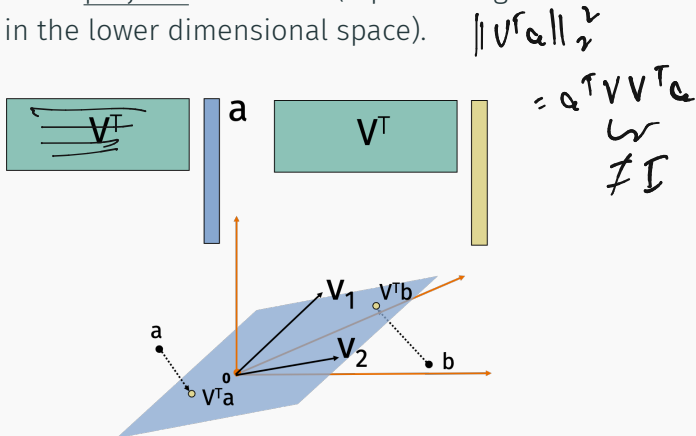
$$= a^T b$$

$$= \langle a, b \rangle$$



LINEAR ALGEBRA REMINDER

Multiplying a vector by a rectangular matrix V^T with orthonormal rows projects the vector (representing it as coordinates in the lower dimensional space).

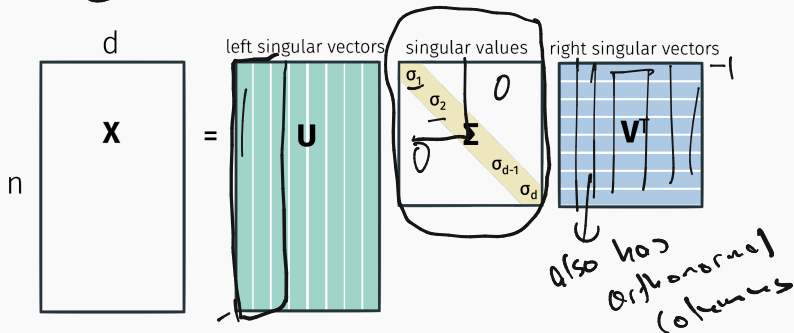


So we always have that $\|V^T x\|_2 \leq \|x\|_2$.

SINGULAR VALUE DECOMPOSITION

One of the most fundamental results in linear algebra.

Any matrix X can be written:



Where $\mathbf{U}^T \mathbf{U} = \mathbf{I}$, $\mathbf{V}^T \mathbf{V} = \mathbf{I}$, and $\underline{\sigma}_1 \geq \underline{\sigma}_2 \geq \dots \underline{\sigma}_d \geq 0$.

Singular values are unique. Factors are not. E.g. would still get a valid SVD by multiplying both i^{th} column of \mathbf{V} and \mathbf{U} by -1 .

SINGULAR VALUE DECOMPOSITION

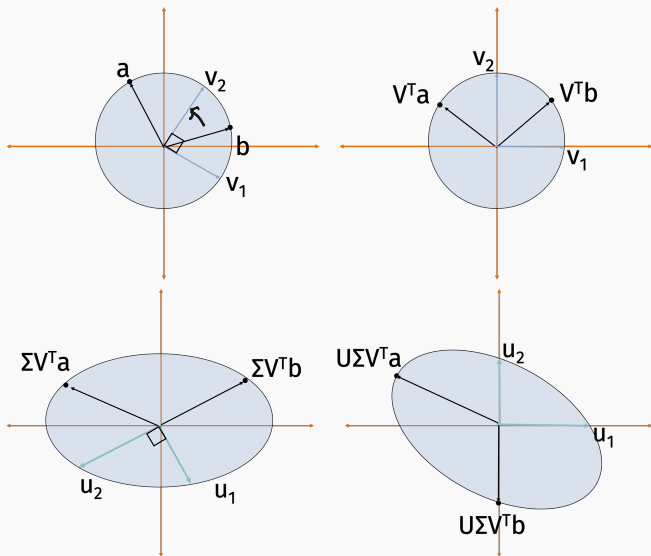
$$X = U \Sigma V^T \quad \textcircled{U} (\Sigma (V^T x))$$

Important take away from singular value decomposition.

Multiplying any vector \mathbf{a} by a matrix \mathbf{X} to form \mathbf{Xa} can be viewed as a composition of 3 operations:

1. Rotate/reflect the vector (multiplication by \mathbf{V}^T).
2. Scale the coordinates (multiplication by $\mathbf{\Sigma}$).
3. Rotate/reflect the vector again (multiplication by \mathbf{U}).

SINGULAR VALUE DECOMPOSITION: ROTATE/REFLECT



COMPARISON TO EIGENDECOMPOSITION

A square matrix has at most d linearly independent eigenvectors. If a matrix has a full set of d eigenvectors $\mathbf{v}_1, \dots, \mathbf{v}_d$ with eigenvalues $\lambda_1, \dots, \lambda_d$ it is called "diagonalizable" and can be written as:

$$\underline{\underline{V \Lambda V^{-1}}}$$

↗ diagonal

V 's columns are $\mathbf{v}_1, \dots, \mathbf{v}_d$.

$$V^T$$

$$V^{-1} = V^T$$

Singular value decomposition

- Exists for all matrices, square or rectangular.
- Singular values are always positive.
- Factors \mathbf{U} and \mathbf{V} are orthogonal.

Eigendecomposition

- Exists for some square matrices.
- Eigenvalues can be positive, negative, or imaginary. Real if \mathbf{X} is symmetric.
- Factor \mathbf{V} is orthogonal if and only if \mathbf{X} is symmetric.

CONNECTION TO EIGENDECOMPOSITION

- U contains the orthogonal eigenvectors of XX^T .
- V contains the orthogonal eigenvectors of X^TX .
- $\sigma_i^2 = \lambda_i(\underline{XX^T}) = \lambda_i(\underline{X^TX})$

$$X = U \Sigma V^T$$

$$\begin{aligned} \underline{XX^T} &= (U \Sigma V^T)(U \Sigma V^T)^T \\ &= U \Sigma \underbrace{V^T V}_{I} \Sigma U^T = \underline{U \Sigma^2 U^T} \end{aligned}$$

SVD APPLICATIONS

Lots of applications.

$$\therefore \boxed{(X^T X)^{-1} X^T} \quad f(w) = \|Xw - b\|_2^2$$

- Compute pseudoinverse $\underline{V \Sigma^{-1} U^T}$.

- Read off condition number of X σ_1^2 / σ_d^2 .

Track with 1
3:30

- Compute matrix norms. E.g. $\|X\|_2 = \sigma_1$, $\|X\|_F = \sqrt{\sum_{i=1}^d \sigma_i^2}$.

(Compute matrix square root – i.e. find a matrix B such that $BB^T = X$. Used e.g. in sampling from Gaussian with covariance X .)

$$V \Sigma V^T \quad B = V \sqrt{\Sigma}$$

- Principal component analysis.

(**Killer app:** Read off optimal low-rank approximation for X .)

$$\max_z \frac{\|Xz\|_2}{\|z\|_2}$$

$$\begin{aligned} \|X\|_F &= \|U \Sigma V^T\|_F \\ &= \|\Sigma V^T\|_F \\ &= \|V \Sigma\|_F = \|\Sigma\|_F \end{aligned}$$

The column span of a matrix \mathbf{X} $\in \mathbb{R}^{n \times d}$ is the set of all vectors that can be written as $\mathbf{X}\mathbf{a}$ for some \mathbf{a} .

The dimension of the column span, D_C , is the maximum number of linear independent vectors in that set.

The row span of a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ is the set of all vectors that can be written as $\mathbf{b}^T \mathbf{X}$ for some \mathbf{b} .

The dimension of the row span, D_r , is the maximum number of linear independent vectors in that set.

RANK

For a matrix $X \in \mathbb{R}^{n \times d}$ we have:

$$\begin{matrix} n \\ \left[\begin{array}{c} \vdots \\ c \end{array} \right] \\ r \end{matrix} \quad \begin{matrix} w \\ \left[\begin{array}{c} \vdots \\ w \end{array} \right] \end{matrix}$$

$$\begin{aligned} D_c &\leq d \\ D_r &\leq n \\ D_c &= D_r \end{aligned}$$

$$\begin{matrix} n \\ \left[\begin{array}{c} \vdots \\ c \end{array} \right] \\ d \end{matrix} = \begin{matrix} n \\ \left[\begin{array}{c} \vdots \\ c \end{array} \right] \\ r \end{matrix} \quad \begin{matrix} d \\ \left[\begin{array}{c} \vdots \\ d \end{array} \right] \\ r \end{matrix}$$

We call the value of $D_c = D_r$ the rank of X .

$$\text{rank}(X) \leq \min(d, n)$$

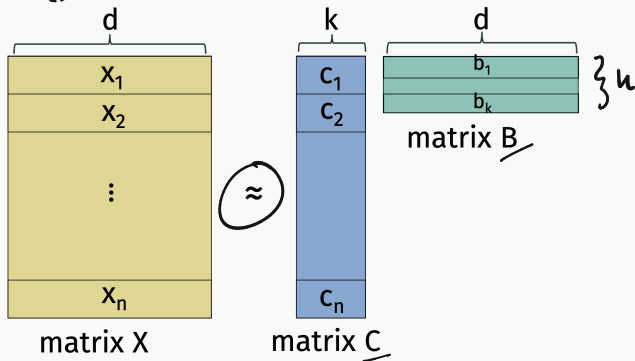
We always have that:

$$\text{rank}(A \cdot B \cdot C \dots) \leq \min(\text{rank}(A), \text{rank}(B), \text{rank}(C), \dots)$$

$$\underbrace{A \cdot B}_{\begin{matrix} r_1 & r_2 \\ \left[\begin{array}{c} \vdots \\ r_1 \end{array} \right] & \left[\begin{array}{c} \vdots \\ r_2 \end{array} \right] \end{matrix}} \rightarrow \begin{matrix} r_1 & r_2 \\ \left[\begin{array}{c} \vdots \\ r_1 \end{array} \right] & \left[\begin{array}{c} \vdots \\ r_2 \end{array} \right] \end{matrix}$$

LOW-RANK APPROXIMATION

Approximate X as a rank k matrix:

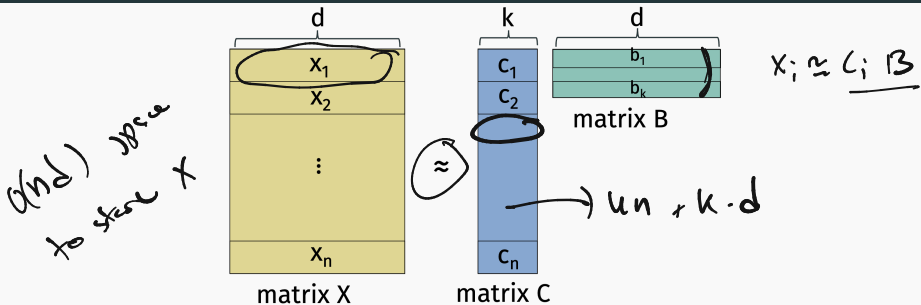


Choose C and B to minimize:

$$\min_{B, C} \|X - CB\|_2$$

for some matrix norm. Common choice is $\|X - CB\|_F^2$.

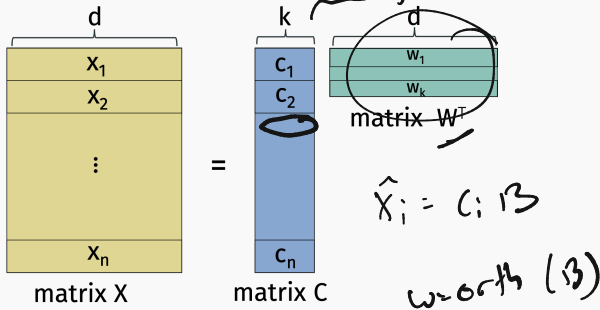
APPLICATIONS OF LOW-RANK APPROXIMATION



- CB takes $O(k(n + d))$ space to store instead of $O(nd)$.
 - Important in many applications, including e.g. [LoRA: Low-Rank Adaptation of Large Language Models](#) (Can be used to compress vector databases.)
 - Many more applications.
- Many linear algebraic problems involving CB can be solved in $O(nk^2)$ instead of $O(nd^2)$ time.

LOW-RANK APPROXIMATION

Without loss of generality can assume that the right matrix is orthogonal. I.e. W^T with $W^T W = I$



Then we should choose left matrix C to minimize:

$$\min_C \|X - \underline{C}W^T\|_F^2$$

(This is just n least squares regression problems!)

LOW-RANK APPROXIMATION

$$\begin{aligned}
 & \|X - CW^T\|_F^2 \\
 &= \|WC^T - X^T\|_F^2 \\
 &= \left\| \begin{bmatrix} \downarrow \\ C^T \end{bmatrix} - \begin{bmatrix} \downarrow \\ X^T \end{bmatrix} \right\|_F^2 = \sum_{i=1}^n \| \underbrace{WC_i - X_i}_{} \|_2^2 \\
 & \quad \downarrow \quad \downarrow \\
 & \quad C_i \quad X_i \\
 & \quad C_i = W^T X_i \\
 & \quad \underline{C} = \underline{XW}
 \end{aligned}$$

$C_i = (W^T W) W^T X_i = W^T X_i$
 $C_i = \arg \min_c \|WC - X_i\|_2^2$
 $= I$ (least squares)

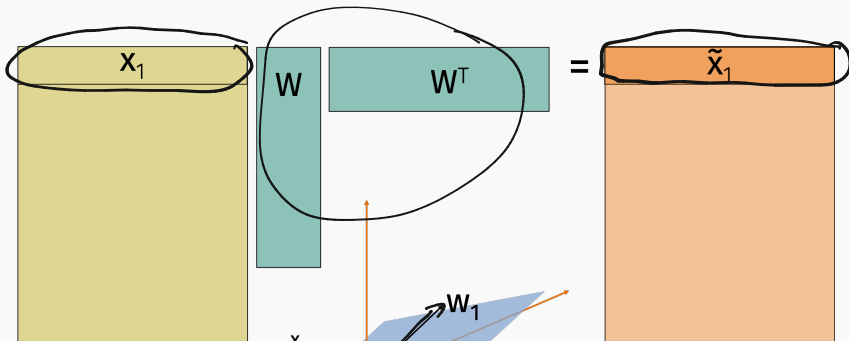
So our optimal low-rank approximation always has the form:

$$X \approx XWW^T$$

PROJECTION MATRICES

WW^T is a symmetric projection matrix.

$$x W W^T$$



$$x, W W^T$$

$$\underline{W W^T} x_1$$

$$\| \tilde{x}_1 - x_1 \|^2 = (1 - \alpha) \| x_1 - y_1 \|^2$$

DATA COMPRESSION

$$X W W^T$$

$C = XW$ can be used as a meaningful compressed version of data matrix X . We have that:

$$\tilde{X}_i \approx X_i$$

$$\|x_i - x_j\|_2 \approx \|WW^T x_i - WW^T x_j\|_2 = \|c_i - c_j\|_2$$

So we expect that:

- $\|x_i\|_2 \approx \|c_i\|_2$
- $\langle x_i, x_j \rangle \approx \langle c_i, c_j \rangle$
- etc.

$$= \|\tilde{x}_i - \tilde{x}_j\|_2$$

$$= \|\omega^T x_i - \omega^T x_j\|_2$$

length k

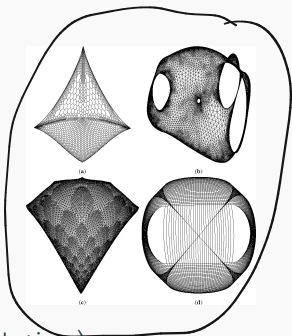
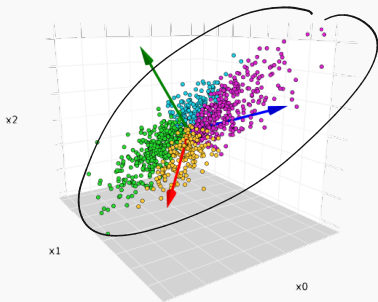
How does this compare to Johnson-Lindenstrauss projection?



APPLICATIONS OF LOW-RANK APPROXIMATION

Also useful in:

(Data visualization when $k = 2$ or 3 .)

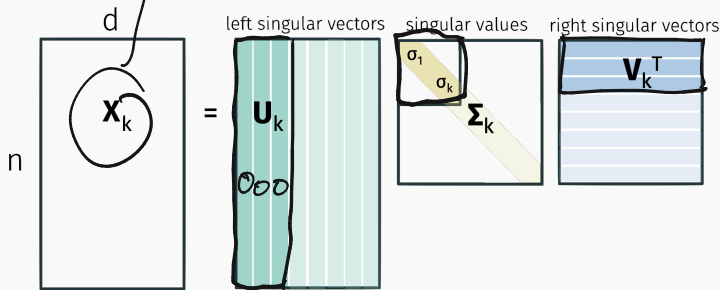


- Data denoising (e.g. distance triangulation).
- Feature selection.

PARTIAL SVD

→ Optimal rank k approx to X .

Key result: Can find the best projection from the singular value decomposition. **Note:** $X_k = \underline{U_k \Sigma_k V_k^T} = \underline{U_k U_k^T X} = \underline{X V_k V_k^T}$ ω



$$U_k = \arg \min_{\text{orthogonal } Z \in \mathbb{R}^{d \times k}} \|X - ZZ^T X\|_F^2$$

$$V_k = \arg \min_{\text{orthogonal } W \in \mathbb{R}^{d \times k}} \|X - XWW^T\|_F^2$$

OPTIMAL LOW-RANK APPROXIMATION

Goal: Minimize $\|X - B\|_F$.

$$X = \underline{U \Sigma U^T}$$

Claim 1: Without loss of generality, can assume $B = \underline{UZV^T}$ for some other rank k matrix Z .

$$\|X - B\|_F = \underbrace{\|XV - BV\|_F} \leq \|U^T XV - U^T BV\|_F$$

$$= \|U(U^T XV - U^T BV)U^T\|_F$$



$$= \|UU^T U \Sigma V^T U^T - UU^T B V U^T\|_F$$

$$= \|U \Sigma V^T - \underbrace{UU^T B V U^T}_{Z}\|_F^2$$

$$Z = U^T B U$$

OPTIMAL LOW-RANK APPROXIMATION

Goal: Minimize $\|X - B\|_F$.

Claim 2: Should choose Z to be the best rank k approximation to Σ . (We will then show this equals Σ_k .)

$$\begin{aligned} \min_{\text{rank } k \ Z} \|X - UZV^T\|_F &= \|U\Sigma V^T - UZV^T\|_F \\ &= \|\Sigma V^T - ZV^T\|_F \\ &= \|\Sigma - Z\|_F \end{aligned}$$

$$UZV^T$$

(Choose Z to be the optimal rank k approx. to Σ .)

OPTIMAL LOW-RANK APPROXIMATION

$U \geq U^T$
 Prove that
 $\arg \min_{rank \leq 2} \|X - Z\|_F$
 $= \Sigma_k$

$$\begin{array}{c} d \\ X_k \end{array} = \begin{array}{c} \text{left singular vectors} \\ U_k \end{array} \begin{array}{c} \text{singular values} \\ \begin{array}{c} \sigma_1 \\ \sigma_k \\ \sigma_{k+1} \end{array} \end{array} \begin{array}{c} \text{right singular vectors} \\ V_k^T \end{array}$$

$$U \Sigma_k U^T = U_k \Sigma_k V_k^T$$

Claim 3:

orthogonal \rightarrow

$$\arg \min_{W \in \mathbb{R}^{d \times k}} \|X - XWW^T\|_F^2 = \arg \max_{W \in \mathbb{R}^{d \times k}} \|XWW^T\|_F^2$$

Follows from fact that for all orthogonal W :

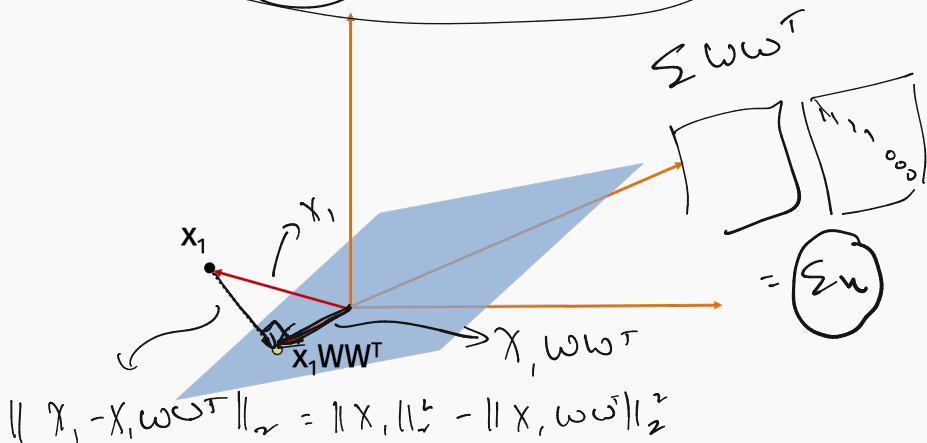
$$\|X - XWW^T\|_F^2 = \|X\|_F^2 - \|XWW^T\|_F^2$$

$$\min_W \|X - XWW^T\|_F^2 = \min_W (\|X\|_F^2 - \|XWW^T\|_F^2) = \max_W \|XWW^T\|_F^2$$

OPTIMAL LOW-RANK APPROXIMATION

Claim 3:

$$\|X - XWW^T\|_F^2 = \|X\|_F^2 - \|XWW^T\|_F^2$$



OPTIMAL LOW-RANK APPROXIMATION

Final Step: Let $W^* \in \mathbb{R}^{d \times k}$ contain the first k standard basis vectors. Then we claim that $W^* = \arg \max_W \|\Sigma W W^T\|_F^2$.

$$\Sigma = \begin{bmatrix} 6_1 & & \\ & \ddots & \\ & & 6_k \\ & & & 0 \end{bmatrix} \quad \left\| \Sigma W W^T \right\|_F^2 = \left\| W^T \Sigma \right\|_F^2$$

$$= \begin{matrix} \downarrow \\ \boxed{\omega} \end{matrix} \quad \boxed{\| \omega \|^2} \quad \begin{matrix} \downarrow \\ \omega_i \end{matrix} \quad \begin{matrix} \downarrow \\ \omega_i \end{matrix}$$

$$\max_W \left\| W^T \Sigma \right\|_F^2 = \sum_{i=1}^n \left\| \omega_i \right\|_2^2 (6_i^2)$$

$$\left\| \omega \frac{\omega_i}{\left\| \omega_i \right\|_2} \right\|_2^2 = \left\| \frac{\omega_i}{\left\| \omega_i \right\|_2} \right\|_2^2 = 1$$

if the entry is $\left\| \omega_i \right\|_2^2 / \left\| \omega_i \right\|_2 = \left\| \omega_i \right\|_2$

$$\left(\sum_{i=1}^n \left\| \omega_i \right\|_2^2 = k \right)$$

(For all i , $\left\| \omega_i \right\|_2^2 \leq 1$)

USEFUL OBSERVATIONS

$$\|X - X_k\|_F^2 = \|X\|_F^2 - \|X_k\|_F^2$$

Diagram illustrating the SVD decomposition of a matrix X_k (size $n \times d$):

- X_k is decomposed into U_k (left singular vectors, size $n \times d$), Σ_k (singular values, size $d \times d$), and V_k^T (right singular vectors, size $d \times d$).
- The singular values σ_1, σ_k are highlighted in the Σ_k matrix.

Observation: The optimal low-rank approximation error

$E_k = \|X - X_k\|_F^2 = \|X\|_F^2 - \|X_k\|_F^2$ can be written:

$$E_k = \sum_{i=k+1}^d \sigma_i^2$$

$$\|X\|_F^2 = \sum_{i=1}^d \sigma_i^2$$

$$\|X_k\|_F^2 = \sum_{i=1}^k \sigma_i^2$$

SPECTRAL PLOTS

Observation: The optimal low-rank approximation error

$E_k = \|\mathbf{X} - \mathbf{X}_k\|_F^2 = \|\mathbf{X}\|_F^2 - \|\mathbf{X}_k\|_F^2$ can be written:

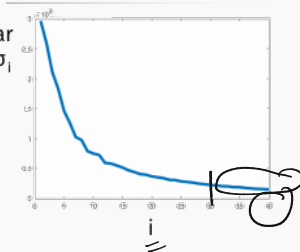
$$E_k = \sum_{i=k+1}^d \sigma_i^2.$$

Can immediately get a sense of “how low-rank” a matrix is from it’s spectrum:

784 dimensional vectors



singular
value σ_i



SPECTRAL PLOTS

Observation: The optimal low-rank approximation error

$E_k = \|\mathbf{X} - \mathbf{X}_k\|_F^2 = \|\mathbf{X}\|_F^2 - \|\mathbf{X}_k\|_F^2$ can be written:

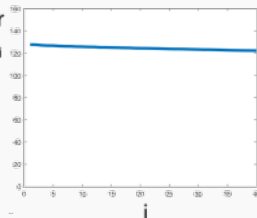
$$E_k = \sum_{i=k+1}^d \sigma_i^2.$$

Can immediately get a sense of “how low-rank” a matrix is from it’s spectrum:

784 dimensional vectors



singular
value σ_i



COMPUTING THE SVD

Suffices to compute right singular vectors V :

- Compute $X^T X$.
- Find eigendecomposition $V \Lambda V^T = X^T X$ using e.g. QR algorithm.
- Compute $L = \underline{XV}$. Set $\sigma_i = \|L_i\|_2$ and $U_i = L_i / \|L_i\|_2$.

$$\begin{array}{c} \textcircled{X^T X} \\ O(\underline{n d^2}) \end{array}$$

$$O(d^3)$$

$$\begin{array}{c} XV = U \Sigma V^T V \\ U \Sigma \end{array}$$

$$\text{Total runtime} \approx O(\underline{n d^2} + d^3)$$

COMPUTING THE SVD (FASTER)

How to go faster?

- { Compute approximate solution. }
- { Only compute top k singular vectors/values. }
- Iterative algorithms achieve runtime $\approx O(ndk)$ vs. $O(nd^2)$ time.
 - { **Krylov subspace methods** like the Lanczos method are most commonly used in practice.
 - { **Power method** is the simplest Krylov subspace method, and still works very well.

Handwritten diagram illustrating the sparsity of the matrix x . A vertical rectangle represents the matrix, with a small d above it and a small n to its left. Inside the rectangle are several small circles representing non-zero entries. An arrow points from the matrix to the equation $nnz(x) \ll n \cdot d$. Below this equation is the equation $nnz(x) = k$.

POWER METHOD

Today: Consider ^{simplest} ~~simplest~~ case when $k=1$.

$$X = U \Sigma U^T$$

Goal: Find some $\underline{z} \approx \underline{v}_1$

Input: $\underline{X} \in \mathbb{R}^{n \times d}$ with SVD $\underline{U} \underline{\Sigma} \underline{V}^T$.

$$O(nd) \quad O(n \log_2(X))$$

$$X^T X z^{(i-1)}$$



$$\sqrt{\frac{1}{d} \sum_{i=1}^d x_i^2}$$

Power method:

- Choose $\underline{z}^{(0)}$ randomly. $\underline{z}_0 \sim \mathcal{N}(0, 1)$.

- $\underline{z}^{(0)} = \underline{z}^{(0)} / \|\underline{z}^{(0)}\|_2$

- For $i = 1, \dots, T$

- $\underline{z}^{(i)} = X^T \cdot (X \underline{z}^{(i-1)})$

- $n_i = \|\underline{z}^{(i)}\|_2$

- $\underline{z}^{(i)} = \underline{z}^{(i)} / n_i$

Return $\underline{z}^{(T)}$

$$X \underline{z}^{(i-1)} \in \mathbb{R}^d$$

$$\boxed{X^T}$$

$$X^T X \underline{z}^{(i-1)} \in \mathbb{R}^d$$

$$\underline{z}^T = C \cdot \underline{(X^T X)}^T \underline{z}^{(0)}$$

Theorem (Basic Power Method Convergence)

Let $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$ be parameter capturing the “gap” between the first and second largest singular values. If Power Method is initialized with a random Gaussian vector then, with high probability, after $T = O\left(\frac{\log d/\epsilon}{\gamma}\right)$ steps, we have either:

$$\|\mathbf{v}_1 - \mathbf{z}^{(T)}\|_2 \leq \epsilon \quad \text{or} \quad \|\mathbf{v}_1 - (-\mathbf{z}^{(T)})\|_2 \leq \epsilon.$$

Total runtime: $O\left(nd \cdot \frac{\log d/\epsilon}{\gamma}\right)$

ONE STEP ANALYSIS OF POWER METHOD

Write $\mathbf{z}^{(i)}$ in the right singular vector basis:

$$\mathbf{z}^{(0)} = c_1^{(0)}\mathbf{v}_1 + c_2^{(0)}\mathbf{v}_2 + \dots + c_d^{(0)}\mathbf{v}_d$$

$$\mathbf{z}^{(1)} = c_1^{(1)}\mathbf{v}_1 + c_2^{(1)}\mathbf{v}_2 + \dots + c_d^{(1)}\mathbf{v}_d$$

$$\vdots$$

$$\mathbf{z}^{(i)} = c_1^{(i)}\mathbf{v}_1 + c_2^{(i)}\mathbf{v}_2 + \dots + c_d^{(i)}\mathbf{v}_d$$

Note: $[c_1^{(i)}, \dots, c_d^{(i)}] = \mathbf{c}^{(i)} = \mathbf{V}^T \mathbf{z}^{(i)}$.

Also: Since \mathbf{V} is orthogonal and $\|\mathbf{z}^{(i)}\|_2 = 1$, $\|\mathbf{c}^{(i)}\|_2^2 = 1$.

ONE STEP ANALYSIS OF POWER METHOD

Claim: After update $\mathbf{z}^{(i)} = \frac{1}{n_i} \mathbf{X}^T \mathbf{X} \mathbf{z}^{(i-1)}$,

$$c_j^{(i)} = \frac{1}{n_i} \sigma_j^2 c_j^{(i-1)}$$

$$\mathbf{z}^{(i)} = \frac{1}{n_i} \left[c_1^{(i-1)} \sigma_1^2 \cdot \mathbf{v}_1 + c_2^{(i-1)} \sigma_2^2 \cdot \mathbf{v}_2 + \dots + c_d^{(i-1)} \sigma_d^2 \cdot \mathbf{v}_d \right]$$

Equivalently: $\mathbf{c}^{(i)} = \frac{1}{n_i} \mathbf{\Sigma}^2 \mathbf{c}^{(i-1)}$.

Claim: After T updates:

$$\mathbf{z}^{(T)} = \frac{1}{\prod_{i=1}^T n_i} \left[c_1^{(0)} \sigma_1^{2T} \cdot \mathbf{v}_1 + c_2^{(0)} \sigma_2^{2T} \cdot \mathbf{v}_2 + \dots + c_d^{(0)} \sigma_d^{2T} \cdot \mathbf{v}_d \right]$$

Let $\alpha_j = \frac{1}{\prod_{i=1}^T n_i} c_j^{(0)} \sigma_j^{2T}$. **Goal:** Show that $\alpha_j \ll \alpha_1$ for all $j \neq 1$.

POWER METHOD FORMAL CONVERGENCE

Since $\mathbf{z}^{(T)}$ is a unit vector, $\sum_{i=1}^d \alpha_i^2 = 1$. So $|\alpha_1| \leq 1$.

If we can prove that $\left| \frac{\alpha_j}{\alpha_1} \right| \leq \sqrt{\frac{\epsilon}{2d}}$ then we will have that $\|\mathbf{v}_1 - \mathbf{z}^{(T)}\|_2^2 \leq \epsilon$.

$$\alpha_j^2 \leq \alpha_1^2 \cdot \frac{\epsilon}{2d}$$

$$1 = \alpha_1^2 + \sum_{j=2}^d \alpha_j^2 \leq \alpha_1^2 + \frac{\epsilon}{2}$$

$$\alpha_1^2 \geq 1 - \frac{\epsilon}{2}$$

$$|\alpha_1| \geq 1 - \frac{\epsilon}{2}$$

$$\|\mathbf{v}_1 - \mathbf{z}^{(T)}\|_2^2 = 2 - 2\langle \mathbf{v}_1, \mathbf{z}^{(T)} \rangle \leq \epsilon$$

POWER METHOD FORMAL CONVERGENCE

Let's see how many steps T it takes to ensure that $\left| \frac{\alpha_j}{\alpha_1} \right| \leq \sqrt{\frac{\epsilon}{2d}}$ where

$\alpha_j = \frac{1}{\prod_{i=1}^T n_i} c_j^{(0)} \sigma_j^{2T}$. Answer will depend on $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$. **Assumption:**

Starting coefficient on first eigenvector is not too small:

$$|c_1^{(0)}| \geq o\left(\frac{1}{\sqrt{d}}\right).$$

We will prove shortly that it holds with probability 99/100.

$$\frac{|\alpha_j|}{|\alpha_1|} = \frac{\sigma_j^{2T}}{\sigma_1^{2T}} \cdot \frac{|c_j^{(0)}|}{|c_1^{(0)}|} \leq$$

Need to set $T =$

STARTING COEFFICIENT ANALYSIS

Need to prove: Starting coefficient on first eigenvector is not too small. I.e., with probability 99/100,

$$|c_1^{(0)}| \geq O\left(\frac{1}{\sqrt{d}}\right).$$

Prove using Gaussian anti-concentration. First use rotational invariance of Gaussian:

$$\mathbf{c}^{(0)} = \frac{\mathbf{V}^T \mathbf{z}^{(0)}}{\|\mathbf{z}^{(0)}\|_2} = \frac{\mathbf{V}^T \mathbf{z}^{(0)}}{\|\mathbf{V}^T \mathbf{z}^{(0)}\|_2} \sim \frac{\mathbf{g}}{\|\mathbf{g}\|_2},$$

where $\mathbf{g} \sim \mathcal{N}(0, 1)^d$.

Need to show that with high probability, first entry of

$$\frac{g}{\|g\|_2} \geq c \cdot \frac{1}{\sqrt{d}}.$$

Part 1: With super high probability (e.g. 99/100),

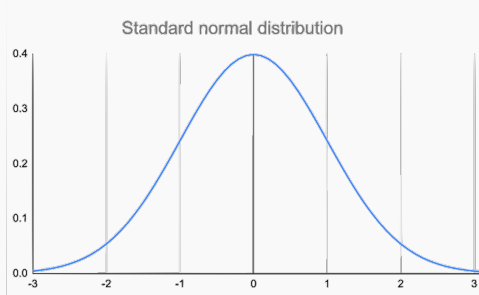
$$\|g\|_2^2 \leq$$

STARTING COEFFICIENT ANALYSIS

Need to show that with high probability, the magnitude of the first entry of $\mathbf{g} \geq c$ for a constant c . Think e.g. $c = 1/100$.

Part 2: With probability $1 - O(\alpha)$,

$$|g_1| \geq \alpha.$$



Theorem (Basic Power Method Convergence)

Let $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$ be parameter capturing the “gap” between the first and second largest singular values. If Power Method is initialized with a random Gaussian vector then, with high probability, after $T = O\left(\frac{\log d/\epsilon}{\gamma}\right)$ steps, we have either:

$$\|\mathbf{v}_1 - \mathbf{z}^{(T)}\|_2 \leq \epsilon \quad \text{or} \quad \|\mathbf{v}_1 - (-\mathbf{z}^{(T)})\|_2 \leq \epsilon.$$

The method truly won't converge if γ is very small. Consider extreme case when $\gamma = 0$.

$$\mathbf{z}^{(T)} = \frac{1}{\prod_{i=1}^T n_i} \left[c_1^{(0)} \sigma_1^{2T} \cdot \mathbf{v}_1 + c_2^{(0)} \sigma_2^{2T} \cdot \mathbf{v}_2 + \dots + c_d^{(0)} \sigma_d^{2T} \cdot \mathbf{v}_d \right]$$

Theorem (Gapless Power Method Convergence)

If Power Method is initialized with a random Gaussian vector then, with high probability, after $T = O\left(\frac{\log d/\epsilon}{\epsilon}\right)$ steps, we obtain a \mathbf{z} satisfying:

$$\|\mathbf{X} - \mathbf{X}\mathbf{z}\mathbf{z}^T\|_F^2 \leq (1 + \epsilon)\|\mathbf{X} - \mathbf{X}\mathbf{v}_1\mathbf{v}_1^T\|_F^2$$

Intuition: For a good low-rank approximation, we don't actually need to converge to \mathbf{v}_1 if σ_1 and σ_2 are the same or very close. Would suffice to return either \mathbf{v}_1 or \mathbf{v}_2 , or some linear combination of the two.

GENERALIZATIONS TO LARGER k

- Block Power Method aka Simultaneous Iteration aka Subspace Iteration aka Orthogonal Iteration

Power method:

- Choose $\mathbf{G} \in \mathbb{R}^{d \times k}$ be a random Gaussian matrix.
- $\mathbf{Z}_0 = \text{orth}(\mathbf{G})$.
- For $i = 1, \dots, T$
 - $\mathbf{Z}^{(i)} = \mathbf{X}^T \cdot (\mathbf{X}\mathbf{Z}^{(i-1)})$
 - $\mathbf{Z}^{(i)} = \text{orth}(\mathbf{Z}^{(i)})$

Return $\mathbf{Z}^{(T)}$

Guarantee: After $O\left(\frac{\log d/\epsilon}{\epsilon}\right)$ iterations:

$$\|\mathbf{X} - \mathbf{X}\mathbf{Z}\mathbf{Z}^T\|_F^2 \leq (1 + \epsilon)\|\mathbf{X} - \mathbf{X}\mathbf{V}_k\mathbf{V}_k^T\|_F^2.$$

Runtime: $O(\text{nnz}(\mathbf{X}) \cdot k \cdot T) \leq O(ndk \cdot T)$.

Possible to “accelerate” these methods.

Convergence Guarantee: $T = O\left(\frac{\log d/\epsilon}{\sqrt{\epsilon}}\right)$ iterations to obtain a nearly optimal low-rank approximation:

$$\|\mathbf{X} - \mathbf{X}\mathbf{Z}\mathbf{Z}^T\|_F^2 \leq (1 + \epsilon)\|\mathbf{X} - \mathbf{X}\mathbf{V}_k\mathbf{V}_k^T\|_F^2.$$

For a normalizing constant c , power method returns:

$$\mathbf{z}^{(q)} = c \cdot (\mathbf{X}^T \mathbf{X})^q \cdot \mathbf{g}$$

Along the way we computed:

$$\mathcal{K}_q = [\mathbf{g}, (\mathbf{X}^T \mathbf{X}) \cdot \mathbf{g}, (\mathbf{X}^T \mathbf{X})^2 \cdot \mathbf{g}, \dots, (\mathbf{X}^T \mathbf{X})^q \cdot \mathbf{g}]$$

\mathcal{K} is called the Krylov subspace of degree q .

Idea behind Krylov methods: Don't throw away everything before $(\mathbf{X}^T \mathbf{X})^q \cdot \mathbf{g}$.

Want to find \mathbf{v} , which minimizes $\|\mathbf{X} - \mathbf{X}\mathbf{v}\mathbf{v}^T\|_F^2$.

Lanczos method:

- Let $\mathbf{Q} \in \mathbb{R}^{d \times k}$ be an orthonormal span for the vectors in \mathcal{K} .
- Solve $\min_{\mathbf{v}=\mathbf{Q}\mathbf{w}} \|\mathbf{X} - \mathbf{X}\mathbf{v}\mathbf{v}^T\|_F^2$.
 - Find best vector in the Krylov subspace, instead of just using last vector.
 - Can be done in $O(ndk + dk^2)$ time.
 - What you're using when you run `svds` or `eigs` in MATLAB or Python.

For a degree t polynomial p , let $\mathbf{v}_p = \frac{p(\mathbf{X}^T \mathbf{X}) \mathbf{g}}{\|p(\mathbf{X}^T \mathbf{X}) \mathbf{g}\|_2}$. We always have that $\mathbf{v}_p \in \mathcal{K}_t$, the Krylov subspace constructed with t iterations.

Power method returns:

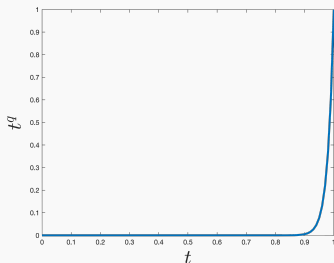
$$\mathbf{v}_p \text{ where } p = x^q \text{ for } q = 2T.$$

Lanczos method returns \mathbf{v}_{p^*} where:

$$p^* = \arg \min_{\text{degree } t \text{ } p} \|\mathbf{X} - \mathbf{X} \mathbf{v}_p \mathbf{v}_p^T\|_F^2.$$

LANCZOS METHOD ANALYSIS

Claim: There is a $t = O\left(\sqrt{q \log \frac{1}{\Delta}}\right)$ degree polynomial \hat{p} approximating \mathbf{x}^q up to error Δ on $[0, \sigma_1^2]$.



$$\|\mathbf{X} - \mathbf{X}\mathbf{v}_{p^*}\mathbf{v}_{p^*}^T\|_F^2 \leq \|\mathbf{X} - \mathbf{X}\mathbf{v}_{\hat{p}}\mathbf{v}_{\hat{p}}^T\|_F^2 \approx \|\mathbf{X} - \mathbf{X}\mathbf{v}_{x^q}\mathbf{v}_{x^q}^T\|_F^2 \approx \|\mathbf{X} - \mathbf{X}\mathbf{v}_1\mathbf{v}_1^T\|_F^2$$

Runtime: $O\left(\frac{\log(d/\epsilon)}{\sqrt{\epsilon}} \cdot \text{nnz}(\mathbf{X})\right)$ vs. $O\left(\frac{\log(d/\epsilon)}{\epsilon} \cdot \text{nnz}(\mathbf{X})\right)$

- Block Krylov methods
- Let $\mathbf{G} \in \mathbb{R}^{d \times k}$ be a random Gaussian matrix.
- $\mathcal{K}_q = \left[\mathbf{G}, (\mathbf{X}^T \mathbf{X}) \cdot \mathbf{G}, (\mathbf{X}^T \mathbf{X})^2 \cdot \mathbf{G}, \dots, (\mathbf{X}^T \mathbf{X})^q \cdot \mathbf{G} \right]$

Runtime: $O\left(\text{nnz}(\mathbf{X}) \cdot k \cdot \frac{\log d/\epsilon}{\sqrt{\epsilon}}\right)$ to obtain a nearly optimal low-rank approximation.