CS-GY 6763: Lecture 10
Ellipsoid Method, Linear Programming,
Singular Value Decomposition

NYU Tandon School of Engineering, Prof. Christopher Musco

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

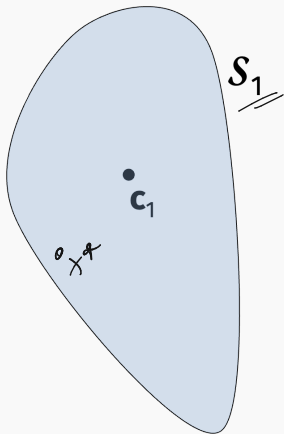### Theorem (Dimension Dependent Convex Optimization)

*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.*
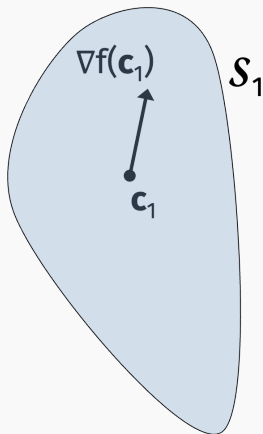
$\gamma_{\epsilon}^2$

$f: \mathbb{B}^d \rightarrow \mathbb{B}$

Natural "cutting plane" method.

- $\mathcal{S}_1 = \mathcal{S}$
- For $t = 1, \ldots, T$:
    - $c_t = $ center of gravity of $\mathcal{S}_t$.
    - Compute $\nabla f(c_t)$.
    - $\mathcal{H} = \{x | \langle \nabla f(c_t), x - c_t \rangle \leq 0\}$.
    - $\mathcal{S}_{t+1} = \mathcal{S}_t \cap H$
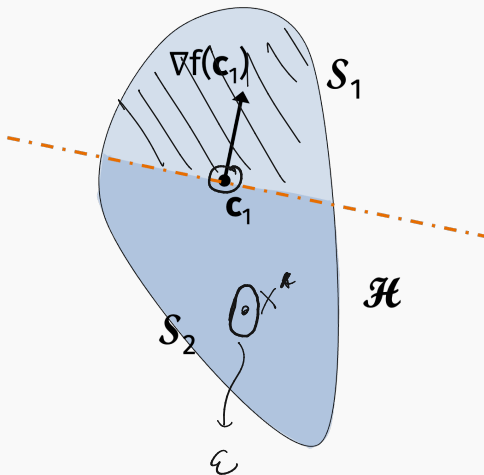- Return $\hat{x} = \arg \min_t f(c_t)$

Natural "cutting plane" method.



- $\mathcal{S}_1 = \mathcal{S}$
- For $t = 1, \ldots, T$:
    - $c_t$ = center of gravity of $\mathcal{S}_t$.
    - Compute $\nabla f(c_t)$.
    - $\mathcal{H} = \{x \,|\, \langle \nabla f(c_t), x - c_t \rangle \leq 0\}$.
    - $\mathcal{S}_{t+1} = \mathcal{S}_t \cap H$
- Return $\hat{x} = \arg\min_t f(c_t)$

4

Natural "cutting plane" method.



- $\mathcal{S}_1 = \mathcal{S}$
- For $t = 1, \ldots, T$:
    - $c_t =$ center of gravity of $\mathcal{S}_t$.
    - Compute $\nabla f(c_t)$.
    - $\mathcal{H} = \{x \mid \langle \nabla f(c_t), x - c_t \rangle \leq 0\}$.
    - $\mathcal{S}_{t+1} = \mathcal{S}_t \cap H$
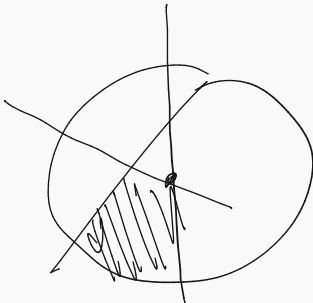- Return $\hat{x} = \arg \min_t f(c_t)$

Proof Reminder:

$V\epsilon$

- By Grünbaum's Theorem, cut the volume of the search space by a constant every step.
- Need to reduce to a convex body whose volume is roughly $\epsilon^d$ smaller than the volume of $\mathcal{S}$.
- Final number of iterations scales with $\log(1/\epsilon^d)$

$$= d\log(1/\epsilon)$$

**In general computing the centroid is hard.** #P-hard even when when $\mathcal{S}$ is an intersection of half-spaces (a polytope).

Even if the problem isn't hard for your starting convex body $\mathcal{S}$, it likely will become hard for $\mathcal{S} \cap \mathcal{H}_1 \cap \mathcal{H}_2 \cap \mathcal{H}_3 \ldots$.

So while the oracle complexity of dimension-dependent optimization was settled in the 60s, a number of basic questions in terms of computational complexity.

We will see how to resolve this issue with an elegant cutting plane methods called the Ellipsoid Method that was introduced by Naum Shor in 1977.

Generalization

**Slightly more general problem:** Given a convex set $\mathcal{K}$ via access to separation oracle $S_{\mathcal{K}}$ for the set, determine if $\mathcal{K}$ is empty, or otherwise return any point $x \in \mathcal{K}$

$$S_{\mathcal{K}}(y) = \begin{cases} \emptyset & \text{if } y \in \mathcal{K}. \\ \text{separating hyperplane } (a, c) & \text{if } y \notin \mathcal{K}. \end{cases}$$
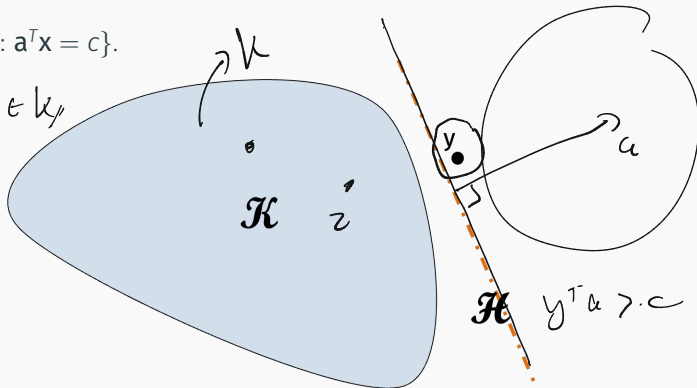
Let $\mathcal{H} = \{x : a^T x = c\}$.
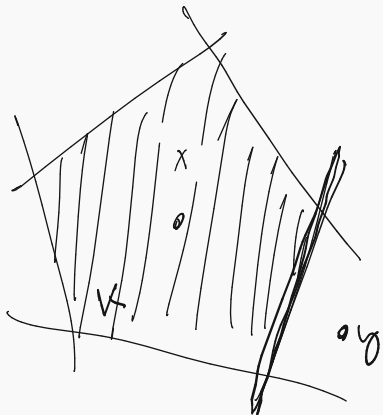
For any $z \in \mathcal{K}$,
$z^T a < c$

Also
$y^T a \geq c$



8

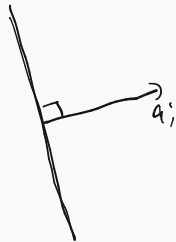**Example:** How would you implement a separation oracle for a polytope $\{x : \underline{Ax \geq b}\}$.



$$x^{\top} a_1 \geq b_1$$
$$\boxed{x^{\top} a_2 \geq b_2}$$
$$\vdots$$
$$x^{\top} a_n \geq b_n$$

$$y^{\top} a_i < b_i$$

but

$$\underline{z^{\top} a_i \geq b_1} \quad \text{for all } z \in k$$

Original problem: $\min_{x \in \mathcal{S}} f(x)$.
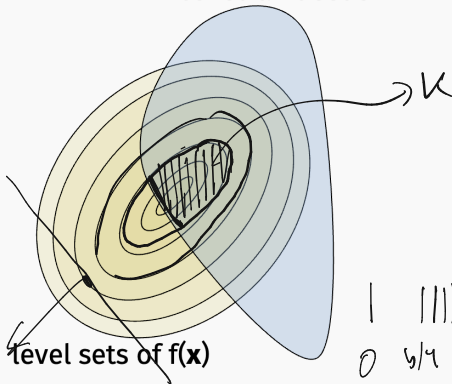
$K = \mathcal{S} \cap \{f(x) \leq c\}$

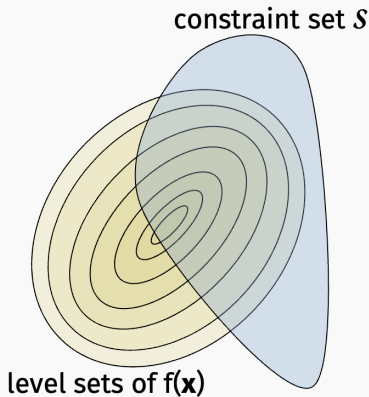How to reduce to determining if a convex set $\mathcal{K}$ is empty or not?



$\log(1/\omega)$

constraint set $S$

$\log(b/\omega)$

$\rightarrow K$

level sets of f(x)

$| \quad |||| \quad | \qquad |$
$0 \quad b/4 \quad b/v \qquad b$

10

Original problem: $\min_{x \in \mathcal{S}} f(x)$.

How to reduce to determining if a convex set $\mathcal{K}$ is empty or not?



constraint set $\mathcal{S}$
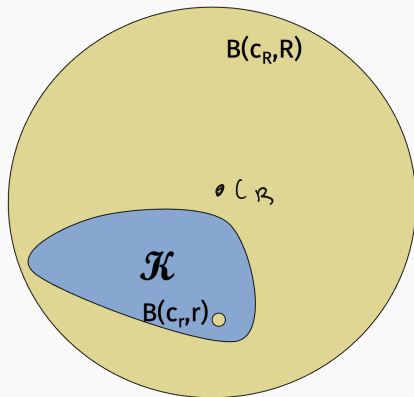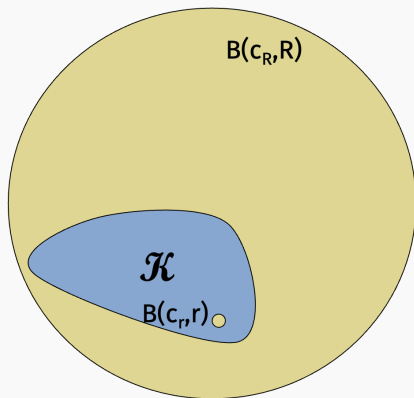
level sets of f($x$)

**Binary search!** For a convex function $f(x)$, $\{x : f(x) \le c\}$ is convex, and you can get a separation oracle via the gradient.

11

**Goal of ellipsoid algorithm:** Solve "Is $\mathcal{K}$ empty or not?" given a separation oracle for $\mathcal{K}$ under the assumptions that:

1. $\mathcal{K} \subset B(c_R, R)$.
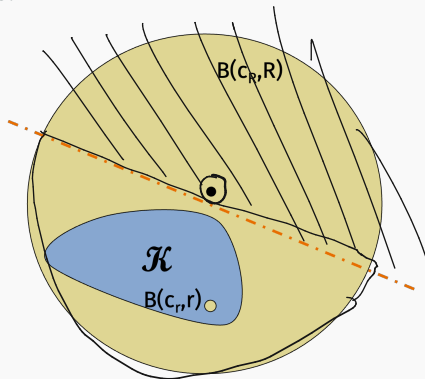2. If non-empty, $\mathcal{K}$ contains $B(c_r, r)$ for some $r < R$.

**Application to original problem:** For a convex function $f$ such that $\|\nabla f(\mathbf{x})\|_2 \leq G$, it can be checked that the convex set $\{\mathbf{x} : f(\mathbf{x}) \leq \epsilon\}$ contains a ball of radius $\epsilon/G$.
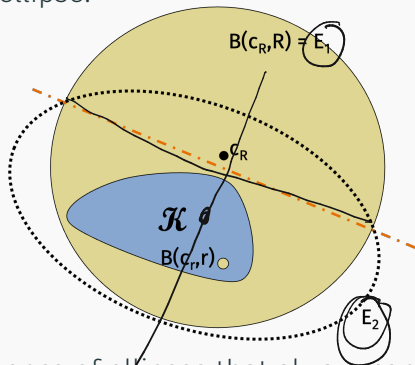
Iterative method similar to center-of-gravity:

1. Check if center $c_R$ of $B(c_R, R)$ is in $\mathcal{K}$.
2. If it is, we are done.
3. If not, cut search space in half, using separating hyperplane.

**Key insight:** Before moving on, approximate new search region by something that we can easily compute the centroid of. Specifically an ellipse!



Produce a sequence of ellipses that <u>always contain</u> $\mathcal{K}$ and decrease in volume: $B(c_R, R) = E_1, E_2, \ldots$. Once we get to an ellipse with volume $\leq B(c_r, r)$, we know that $\mathcal{K}$ must be empty.

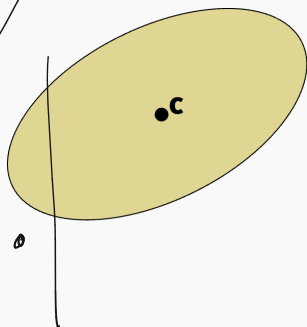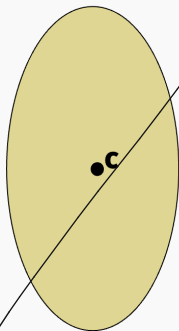An ellipse is a convex set of the form: $\{x : \|A(x - c)\|_2^2 \leq \alpha\}$ for some constant $c$ and matrix $A$. The center-of-mass is $c$.

) center

**{x: $\|I(x-c)\|^2 < \alpha$}**     **{x: $\|D(x-c)\|^2 < \alpha$}**     **{x: $\|A(x-c)\|^2 < \alpha$}**



$\bullet$ **c**

$\bullet$ **c**

$\bullet$ **c**

$Q = \left(A^\top A\right)^{-1} \alpha$

$\varnothing$

Often re-parameterized to say that the ellipse is all $x$ with
$\{x : (x - c)^T Q^{-1} (x - c) \leq 1\}$

There is a closed form solution for the equation of the smallest ellipse containing a given half-ellipse. I.e. let $E_i$ have parameters $Q_i, c_i$ and consider the half-ellipse:

$$E_i \cap \{x : a_i^T x \leq a_i^T c_i\}.$$

Then $E_{i+1}$ is the ellipse with parameters:

$$Q_{i+1} = \frac{d^2}{d^2 - 1}\left(Q_i - \frac{2}{d+1}hh^T\right) \qquad c_{i+1} = c_i - \frac{1}{n+1}h,$$

where $h = \sqrt{a_i^T Q_i a_i} \cdot a_i$.

Computing the update takes $O(d^2)$ time.

$a_i^\dagger Q_i a_i :$

$d \times d$

$\approx .64$

**Claim:** $\mathrm{vol}(E_{i+1}) \leq (1 - \frac{1}{2d}) \, \mathrm{vol}(E_i)$.

$\left(1 - \frac{1}{2d}\right)^{2d} \approx \frac{1}{e}$

**Proof:** Via reduction to the "isotropic case". I will post a proof on the course website if you are interested.



Not as good as the $(1 - \frac{1}{e})$ constant-factor volume reduction we got from center-of-gravity, but still very good!

18

$d \log(1/\epsilon)$

**Claim:** $\text{vol}(E_{i+1}) \leq (1 - \frac{1}{2d}) \text{vol}(E_i)$

If run for

$\tilde{d}^2 \log(B/r)$ steps)

then

$\text{Vol}(E_T) \leq \left(\frac{r}{B}\right)^d$

$\left(\frac{B}{r}\right)^d$

$\text{Vol}(E_i)$

$\text{Vol}(E_T) \leq \text{Vol}(B(c_r, r))$



$B(c_R, R) = E_1$

$c_R$

$\mathcal{K}$

$B(c_r, r)$

$E_2$

After $O(d)$ iterations, we reduce the volume by a constant.

In total require $O(d^2 \log(R/r))$ iterations to solve the problem.

19

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

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

$$\min \underbrace{f(x)} = \underbrace{c^T x}_{\text{fixed}}$$
$$\text{subject to } \underbrace{Ax \geq b.}$$

Think about $Ax \geq b$ as a union of half-space constraints:

$$\{x : a_1^T x \geq b_1\}$$
$$\{x : a_2^T x \geq b_2\}$$
$$\vdots$$
$$\{x : a_n^T x \geq b_n\}$$

20

$$\min f(\mathbf{x}) = \mathbf{c}^T \mathbf{x}$$
$$\text{subject to } \mathbf{Ax} \geq \mathbf{b}$$

$\mathbf{Ax} \geq \mathbf{b}$

$\mathbf{c}$

- Classic optimization applications: industrial resource optimization problems were killer app 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.

Many combinatorial optimization problems can be solved via LP relaxations.

### 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.*

## 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-

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.

23

## 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.*



# Breakthrough in Problem Solving

## By JAMES GLEICK

A 28-year-old mathematician at A.T.&T. Bell Laboratories has made a startling theoretical breakthrough in the solving of systems of equations that often grow too vast and complex for the most powerful computers.

The discovery, which is to be formally published next month, is already circulating rapidly through the mathematical world. It has also set off a deluge of inquiries from brokerage houses, oil companies and airlines, industries with millions of dollars at stake in problems known as linear programming.

ments of great progress, and this may well be one of them.''

Because problems in linear programming can have billions or more possible answers, even high-speed computers cannot check every one. So computers must use a special procedure, an algorithm, to examine as few answers as possible before finding the best one — typically the one that minimizes cost or maximizes efficiency.

A procedure devised in 1947, the simplex method, is now used for such prob-
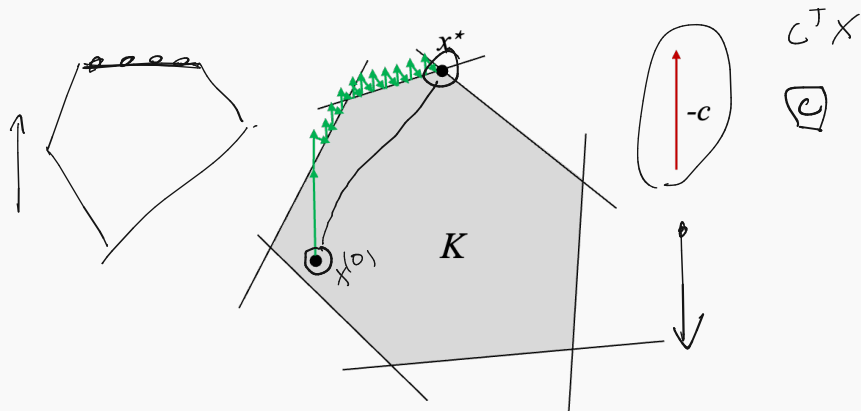
Front page of New York Times, November 19, 1984.
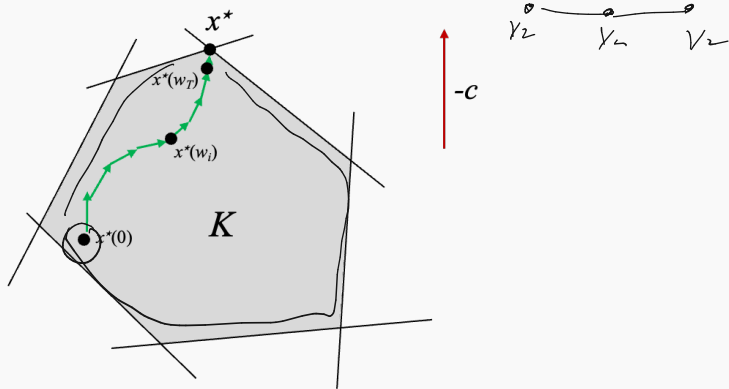
Lecture notes are posted on the website (optional reading).



Projected Gradient Descent Optimization Path

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


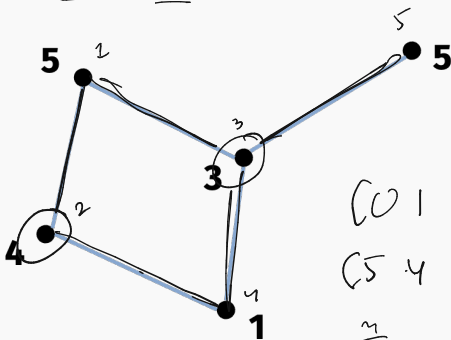
Ideal Interior Point Optimization Path

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.

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



$$[0\ 1\ 1\ 0\ 0] = x$$
$$[5\ 4\ 3\ 1\ 5] = w$$
$$\sum_{i=1}^{n} x_i w_i = 4+3 = 7$$

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

28

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

$0, 1, 2$

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

$$\min_{\mathbf{x}} \sum_{i=1}^{n} x_i w_i \quad \text{subject to} \quad x_i \in \{0, 1\} \text{ for all } i$$

→ constraints

$f(x)$

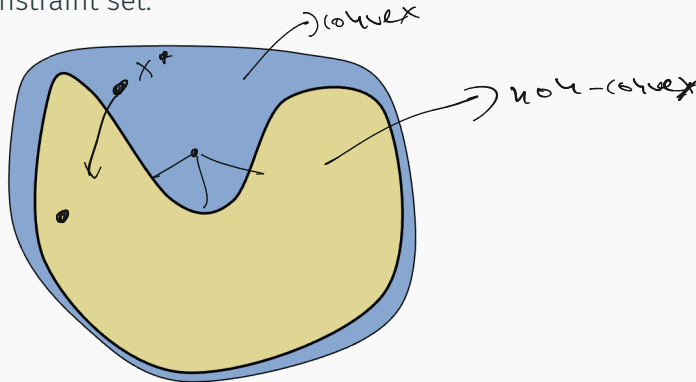$$x_i + x_j \geq 1 \text{ for all } (i, j) \in E$$

**NP-hard to solve exactly.** 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?

$\min_{x} f(x)$     where     $f(x) = \sum_{i=1}^{n} x_i w_i$

29
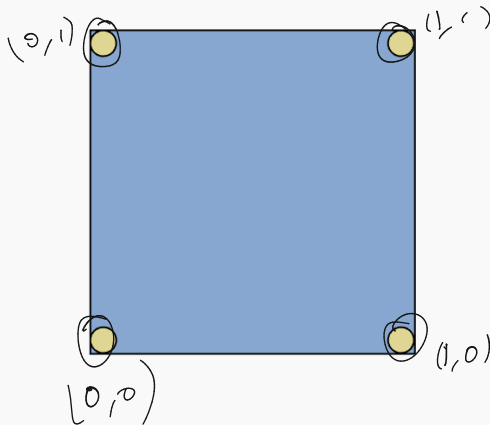
High level approach:

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

High level approach:

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

**High level approach**:

→ *original constraints*

*relaxed 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 $x^* = \arg\min_{x \in \mathcal{S}} f(x)$ and let $\bar{x}^* = \arg\min_{x \in \bar{\mathcal{S}}} f(x)$. We always have that:

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

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.

$$\bar{x}^* \longrightarrow y \in \mathcal{S} \qquad f(y) \leq f(\bar{x}^*) + \epsilon$$

32

Vertex Cover:

$$\min_{\mathbf{x}} \sum_{i=1}^{n} x_i w_i \qquad \text{subject to} \qquad x_i \in \{0, 1\} \text{ for all } i$$

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

Relaxed Vertex Cover:

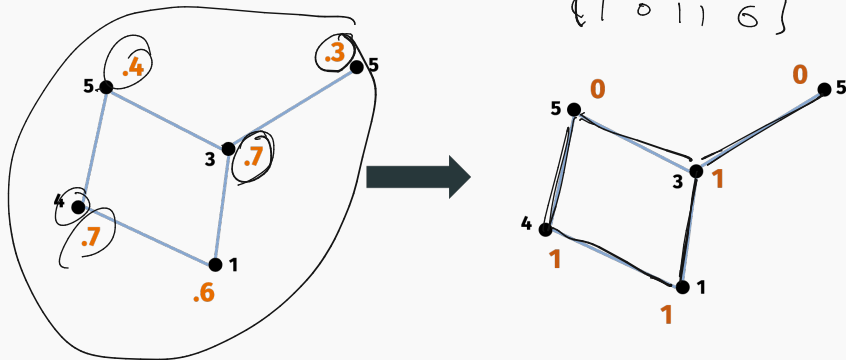$$\min_{\mathbf{x}} \sum_{i=1}^{n} x_i w_i \qquad \text{subject to} \qquad 0 \leq x_i \leq 1 \text{ for all } i$$

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

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

$$f(\overline{x}^*) \leq \min_{x \in S} f(x)$$

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

$$\{ 1 \quad 0 \quad 1 \; 1 \quad 6 \}$$



**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.

**Observation 2:** Let 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:**

$$f(x) = \sum_{i=1}^{n} w_i \; \mathbb{1}[\bar{x}_i^* \geq 1/2] \leq \sum_{i=1}^{n} w_i \cdot 2 \cdot \bar{x}_i^*$$

$$= 2 f(\bar{x}^*)$$

$$f(\bar{x}^*) = \sum_{i=1}^{n} w_i \; \bar{x}_i^*$$

$$. \, 7 \rightarrow 1$$

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

35

## 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.

There is a simpler greedy 2-approximation algorithm that doesn't use optimization at all. Try coming up with it on your own!

BREAK

Next section of course: <u>Spectral methods</u> and <u>numerical linear algebra</u>.

Spectral methods generally refer to methods based on the "spectrum" of a matrix. I.e. on it's eigenvectors/eigenvalues and singular vectors/singular values. We will look at applications in:

- Low-rank approximation and dimensionality reduction.
- Data clustering and related problems.
- Constructing data embeddings (e.g. Word2Vec).

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

$$Xv = \lambda v$$

The scalar $\lambda$ is called the eigenvalue associated with $v$.

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.

If a <u>square</u> matrix has orthonormal rows, it also has orthonormal columns:

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

$$V^T V = I = VV^T$$

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

39

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

Same thing goes for Frobenius norm: for any matrix $\mathbf{X}$,
$\|\mathbf{VX}\|_F^2 = \|\mathbf{X}\|_F^2$ and $\|\mathbf{V}^T\mathbf{X}\|_F^2 = \|\mathbf{X}\|_F^2$.

$$\|Vx\|_2^2 = (Vx)^T Vx = x^T V^T V x = x^T x = \|x\|_2^2$$
$$\underbrace{}_{\mathcal{I}}$$

$$\|V^T x\|_2^2 = x V V^T x = \|x\|_2$$
$$\underbrace{}_{\mathcal{I}}$$

40

The same is <u>not true</u> for rectangular matrices:



$$V^T V = I \qquad \text{but} \qquad VV^T \neq I$$

For any $x$, $\|Vx\|_2^2 = \|x\|_2^2$ <u>but</u> $\|V^T x\|_2^2 \neq \|x\|_2^2$ in general.

Multiplying a vector by **V** with orthonormal columns <u>rotates</u> <u>and/or reflects</u> the vector.

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

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

Multiplying a vector by a rectangular matrix $\mathbf{V}^T$ with orthonormal rows <u>projects</u> the vector (representing it as coordinates in the lower dimensional space).



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

One of the most fundamental results in linear algebra.

Underline(Any) matrix $X$ can be written:



Where $U^T U = I$, $V^T V = I$, and $\sigma_1 \geq \sigma_2 \geq \ldots \sigma_d \geq 0$.

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

Important <u>take away</u> from singular value decomposition.

Multiplying any vector **a** by a matrix **X** to form **Xa** can be viewed as a composition of 3 operations:

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

$$X\,a$$

$$U\left(\Sigma\left(V^T a\right)\right)$$

Recall that an eigenvalue of a <u>square</u> matrix $X \in \mathbb{R}^{d \times d}$ is any vector $v$ such that $Xv = \lambda v$. A matrix has at most $d$ linearly independent eigenvectors. If a matrix has a full set of $d$ eigenvectors $v_1, \ldots, v_d$ with eigenvalues $\lambda_1, \ldots, \lambda_d$ it is called "diagonalizable" and can be written as:

$$X = V \Lambda V^{-1}.$$

$V$'s columns are $v_1, \ldots, v_d$.

$$\begin{bmatrix} v_1 & \cdots & v_d \end{bmatrix} \qquad \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_d \end{bmatrix}$$

### Singluar value decomposition

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

### Eigendecomposition

- Exists for <u>some</u> square matrices.
- Eigenvalues can be positive or negative.
- Factor $V$ is orthogonal if and only if $X$ is symmetric.

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

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

$$= U\Sigma^2 U^T$$

Lots of applications.

- Compute pseudoinverse $V\Sigma^{-1}U^T$.
- Read off condition number of $X$, $\sigma_1^2/\sigma_d^2$.
- 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 $X$.

$$B = U\Sigma^{1/2}U^T$$

- Principal component analysis.

**Killer app:** Read off optimal low-rank approximations for $X$.

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

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

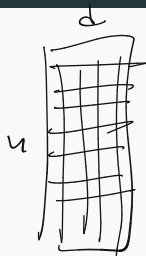The row span of a matrix $X \in \mathbb{R}^{n \times d}$ is the set of all vectors that can be written as $b^T X$ for some $b$.

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

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

$$D_c \leq d$$
$$D_r \leq n$$
$$D_c = D_r.$$

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

$$\text{rank} \leq \min(d, n)$$
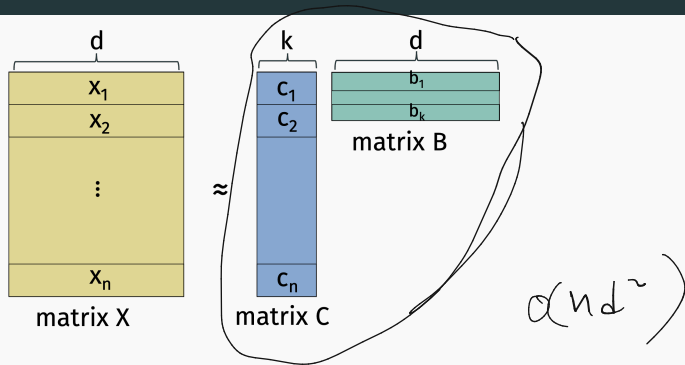
Approximate X as the product of two rank $k$ matrices: $k \ll d$



Typically choose C and B to minimize:

$$\min_{B,C} \|X - CB\|$$
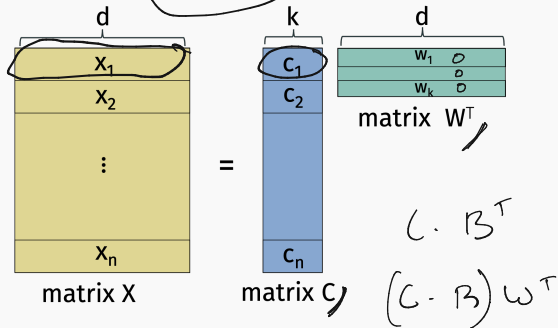
for some matrix norm. Common choice is $\|X - CB\|_F^2$. 53

- **CB** takes $O(k(n + d))$ space to store instead of $O(nd)$.
- Regression problems involving **CB** can be solved in $O(nk^2)$ instead of $O(nd^2)$ time.
- Will see a bunch more in a minute.

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



matrix X          matrix C

$C \cdot B^T$

$(C \cdot B)W^T$

Then we should choose $C$ to minimize:

$$\min_{C,W}\|X - CW^T\|_F^2$$

This is just $n$ least squares regression problems!

$\min \|x_1 - C_1^T W\|_2^2$

$= \arg\min_{C_1} \|Wc_1 - x_1\|_2^2$

55

$x_i \approx x_i \omega \omega^T$

$x_j \approx x_j \omega \omega^T$

$X - X G G^T$

$$c_i = \arg\min_c \|Wc - x_i\|_2^2$$

$c_i = (\omega^T \omega)^{-1} \omega^T x_i = \omega^T x_i$

$\underbrace{}_{I}$

$\min \|X - X\omega\omega^T\|_F^2$

$= \sum_{r=1}^{ } \|x_i - x_i \omega\omega^T\|_F^2$

$$c_i = W^T x_i$$

$$C = XW$$

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

$$X \approx XWW^T$$

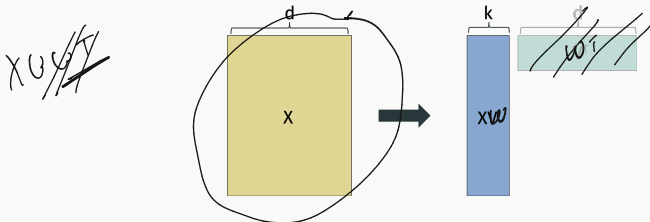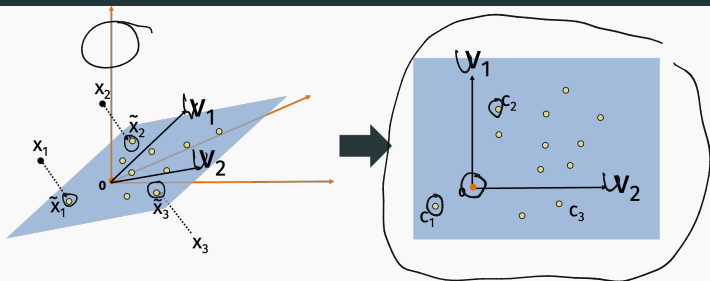$\min_{C, B} \|X - CB\|_F^2$ $\quad$ $\min_\omega \|X - X\omega\omega^T\|_F^2$ $\quad$ where $\omega^T\omega = I$

$WW^T$ is a symmetric <u>projection matrix</u>.

$C = XW$ can be used as a compressed version of data matrix $X$.

$$c_i = x_{\cdot i}\, w \qquad c_j = x_j\, w$$

Let $C = XW$. We have that:

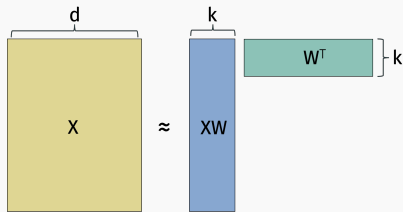$$\|x_i - x_j\|_2 \approx \|x_i^T WW^T - x_j^T WW^T\|_2 = \|c_i - c_j\|_2$$

Similarly, we expect that:

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

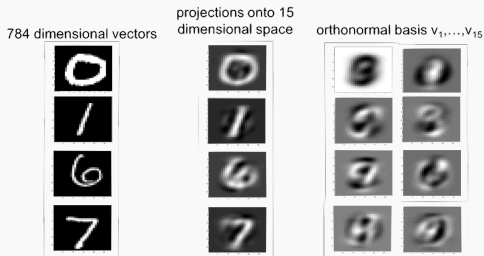How does this compare to Johnson-Lindenstrauss projection?

Rows of X (data points) are approximately spanned by *k* vectors. Columns of X (data features) are approximately spanned by *k* vectors.

If a data set only had $k$ unique data points, it would be exactly rank $k$. If it has $k$ "clusters" of data points (e.g. the 10 digits) it's often very close to rank $k$.
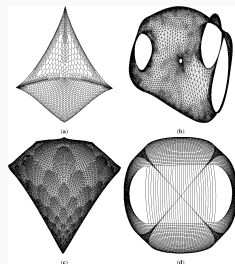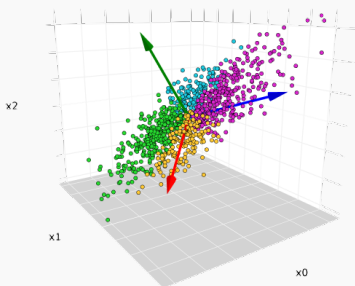


784 dimensional vectors

projections onto 15 dimensional space

orthonormal basis $v_1,...,v_{15}$

Colinearity/correlation of data features leads to a low-rank data matrix.

| | bedrooms | bathrooms | sq.ft. | floors | list price | sale price |
|---|---|---|---|---|---|---|
| home 1 | 2 | 2 | 1800 | 2 | 200,000 | 195,000 |
| home 2 | 4 | 2.5 | 2700 | 1 | 300,000 | 310,000 |
| . | . | . | . | . | . | . |
| . | . | . | . | . | . | . |
| . | . | . | . | . | . | . |
| home n | 5 | 3.5 | 3600 | 3 | 450,000 | 450,000 |

Fact that $\|x_i - x_j\|_2 \approx \|x_i^T WW^T - x_j^T WW^T\|_2 = \|c_i - c_j\|_2$ leads to lots of applications.

- Data compression. E.g. used in state-of-the-art data dependent methods for nearest neighbor search.
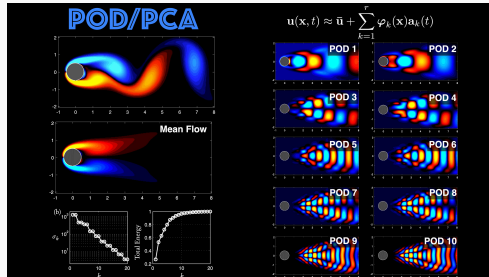- Data visualization when $k = 2$ or $3$.

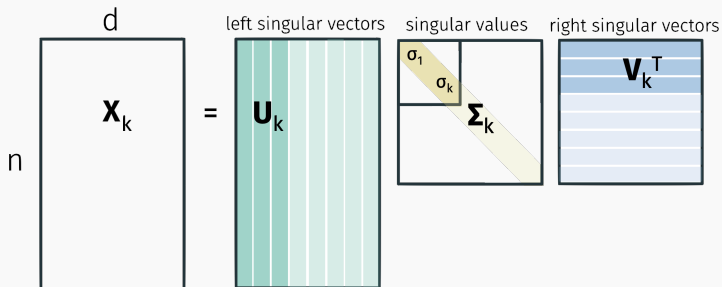

- Data embeddings (next lecture).

- Reduced order modeling for solving physical equations.



- Constructing preconditioners in optimization.
- Noisy triangulation (on problem set).

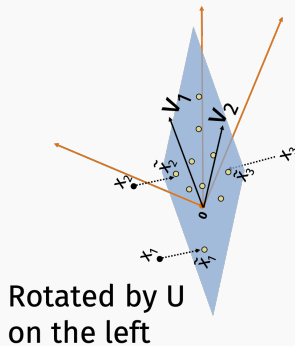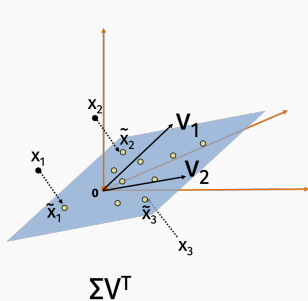Can find the best projection from the singular value decomposition.



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

Claim: $X_k = U_k \Sigma_k V_k^T = X V_k V_k^T$.

Claim 1:

$$\underset{\text{rank } k \ B}{\arg \min} \|X - B\|_F^2 = U \cdot \underset{\text{rank } k \ B}{\arg \min} \|\Sigma V^T - B\|_F^2$$



$\Sigma V^T$

Rotated by U
on the left

Claim 2:

$$\operatorname*{arg\,min}_{\text{rank } k \text{ B}} \|\mathbf{\Sigma}\mathsf{V}^T - \mathsf{B}\|_F^2 = \operatorname*{arg\,min}_{\text{rank } k \text{ B}} \|\mathsf{V}\mathbf{\Sigma} - \mathsf{B}^T\|_F^2$$

Claim 3:

$$\operatorname*{arg\,min}_{\text{rank } k \text{ B}} \|\mathsf{V}\mathbf{\Sigma} - \mathsf{B}^T\|_F^2 = \operatorname*{arg\,min}_{\text{rank } k \text{ B}} \|\mathbf{\Sigma} - \mathsf{V}^T\mathsf{B}^T\|_F^2$$

Chose $\mathsf{B}^T$ so that $\mathsf{V}^T\mathsf{B}^T = \mathbf{\Sigma}_k$.

left singular vectors   singular values   right singular vectors

Observation 1:

$$\underset{\mathbf{W}\in\mathbb{R}^{d\times k}}{\arg\min} \|\mathbf{X} - \mathbf{X}\mathbf{W}\mathbf{W}^T\|_F^2 = \underset{\mathbf{W}\in\mathbb{R}^{d\times k}}{\arg\max} \|\mathbf{X}\mathbf{W}\mathbf{W}^T\|_F^2$$
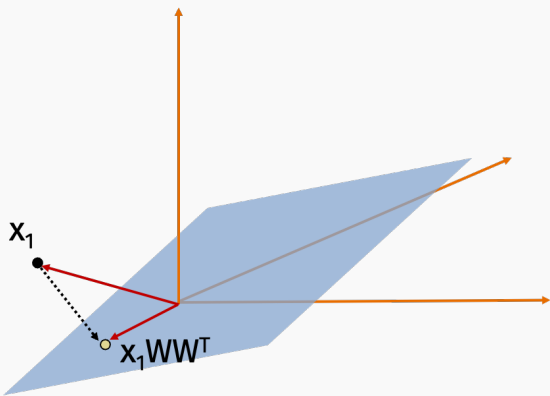
Follows from fact that for all orthogonal $\mathbf{W}$:

$$\|\mathbf{X} - \mathbf{X}\mathbf{W}\mathbf{W}^T\|_F^2 = \|\mathbf{X}\|_F^2 - \|\mathbf{X}\mathbf{W}\mathbf{W}^T\|_F^2$$
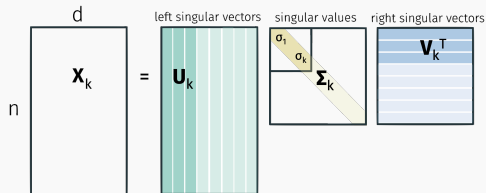
69

Claim:

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

Observation 2: The optimal low-rank approximation error
$E_k = \|X - XV_kV_k^T\|_F^2 = \|X\|_F^2 - \|XV_kV_k^T\|_F^2$ can be written:

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

**Observation 2:** The optimal low-rank approximation error
$E_k = \|X - XV_kV_k^T\|_F^2 = \|X\|_F^2 - \|XV_kV_k^T\|_F^2$ can be written:
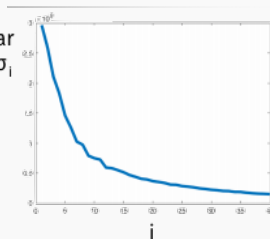
$$E_k = \sum_{i=k+1}^{d} \sigma_i^2.$$

Can immediately get a sense of "how low-rank" a matrix is from it's spectrum:

784 dimensional vectors



singular value $\sigma_i$

i

72

Observation 2: The optimal low-rank approximation error
$E_k = \|X - XV_kV_k^T\|_F^2 = \|X\|_F^2 - \|XV_kV_k^T\|_F^2$ can be written:
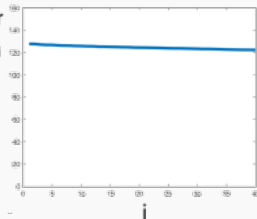
$$E_k = \sum_{i=k+1}^{d} \sigma_i^2.$$

Can immediately get a sense of "how low-rank" a matrix is from it's spectrum:
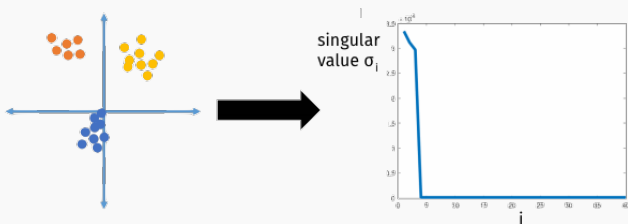


784 dimensional vectors

singular value $\sigma_i$

i

73

**Observation 2:** The optimal low-rank approximation error
$E_k = \|X - XV_kV_k^T\|_F^2 = \|X\|_F^2 - \|XV_kV_k^T\|_F^2$ 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:

## 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 $L = XV$. Set $\sigma_i = \|L_i\|_2$ and $U_i = L_i / \|L_i\|_2$.

Total runtime $\approx$

## COMPUTING THE SVD (FASTER)

Next class:

- Compute <u>approximate</u> solution.
- Only compute <u>top $k$ singular vectors/values</u>. Runtime will depend on $k$. When $k = d$ we can't do any better than classical algorithms based on eigendecomposition.
- <u>Iterative algorithms</u> 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.