

# EE445 Mod3-Lec1: Spectral Properties of Matrices

## References:

- [CE-OptMod]: Chapter 3.3, 4, 5

# Outline

1. Review Eigenvalues & Eigenvectors
2. Symmetric Matrices
3. Introduction to Singular values and SVD

# Why are Spectral Properties Important in ML+OPT?

$$\min_{\alpha} \|K\alpha - y\|^2 + \frac{\lambda}{2} \|\alpha\|^2$$

- Computational efficiency

- Analysis ←

- [Dimensionality reduction]

- Numerical stability ←

$$\left( \underbrace{\quad + \lambda I}_{\quad} \right)^{-1}$$

## How will we see it used?

1. Kernel methods
2. Principle component analysis (unsupervised ML)
3. Principle component regression
4. (time permitting) [spectral clustering]

# Reminder: Eigenvalues & Eigenvectors

Some basics: A matrix,  $A \in \mathbb{R}^{m \times n}$

$$\begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$$

- **Def. (Characteristic Polynomial):**  $p(\lambda) = \det(A - \lambda I)$

The roots of  $p(\lambda) = 0$  are the eigenvalues of  $A$ .

- **Def. (Left/Right Eigenvector-value pair):**

A nonzero vector  $x$  s.t.  $Ax = \lambda x$  then  $(\lambda, x)$  is a right eigenvalue-vec pair

—  $\leftarrow$   $y$  s.t.  $y^* A = \lambda y^*$  then  $(\lambda, y)$  is a left eigen-pair

- **Orthogonality:** Let  $\{x_1, \dots, x_n\}$  of  $A$

orthogonal eigenvectors:  $\langle x_i, x_j \rangle = x_i^T x_j = 0 \quad \forall i \neq j$

orthonormal eigvec:  $\langle x_i, x_j \rangle = 0 \quad \forall i \neq j \quad \therefore \|x_i\| = 1 \quad \forall i$

# Reminder: Eigenvalues & Eigenvectors

## Why important?

- Many ML algorithms involve transforming the matrix  $A$  into simpler, or *canonical forms*, from which it is easy to compute its eigenvalues and eigenvectors.
- These transformations are called **similarity transformations**

$S^{-1}$ ,  $\det(S) \neq 0$       Similarity transforms       $S^*$ : complex conjugate transpose

- **Def. [Similarity Transform]:** The matrix  $A$  is similar to  $B$  if there exists a **non-singular (invertible)** matrix  $S$  s.t.  $B = \underbrace{S^{-1}AS}$
- **Proposition.** Similar matrices  $A$  and  $B$  has the same eigenvalues.  
 $x$  is a right eigvec of  $A \iff S^{-1}x$  is a right eigvec of  $B$   
 $y$  is a left eigvec of  $A \iff S^*y$  is a left eigvec of  $B$
- Some special matrices are similar to diagonal matrices—i.e., for some matrices  $A$ , there is a similarity transform  $S$  such that  $\Lambda = S^{-1}AS$  is diagonal, and  $\Lambda$  contains the eigenvalues of  $A$ .
- These matrices are called **diagonalizable**.

## Part 2. Special Matrices [Symmetric and Positive (semi) definite]

# Symmetric Matrices

Symmetric Matrix:  $A \in \mathbb{R}^{n \times n}$  is symmetric if  $A = A^T$

- Symmetric matrices are one of the most important matrices in linear algebra and ML
- **Mod2-L4:** we often use kernel matrices  $K = [K(x^{(i)}, x^{(j)})]$  and these are symmetric—i.e.,  $K = K^T$ —since  $K(x^{(i)}, x^{(j)}) = K(x^{(j)}, x^{(i)})$
- **Mod2-L2:** Gram matrices  $A^T A$  and  $AA^T$  are symmetric;
  - ▶ in fact we can study all kinds of properties of a matrix  $A$  such as the range and null spaces using these gram matrices (cf. Finite Rank Operator Lemma)



# Symmetric Matrices: Examples

↙  
The graph Laplacian is a symmetric matrix

$$L_{ij} = \begin{cases} \# \text{ of edges incident to node } i, & \text{if } i=j \\ -1, & \text{if there is an edge } (i,j) \\ 0 & \text{otherwise} \end{cases}$$

$$L = \begin{bmatrix} L_{11} & L_{12} & \dots \\ & \ddots & \\ & & L \end{bmatrix}$$

↘  
Sample covariance matrix

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \bar{x})(x^{(i)} - \bar{x})^T, \quad \bar{x} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

↘  
Hessian of a function:

$$H = \nabla^2 f(x), \quad H_{ij} = \frac{\partial^2}{\partial x_i \partial x_j} f(x) \quad \text{mod 2-L2}$$

$$x = (x_1, \dots, x_n), \quad f: \mathbb{R}^n \rightarrow \mathbb{R}$$

# Quadratic Functions

$$f(x) = \underbrace{x^T Q x}$$

- Symmetric matrices play an important role not just in ML but also OPT
- We have seen how to formulate least squares regression as a optimization problem with a quadratic objective:

$$f_{\text{ls}}(x) = \underbrace{\|Ax - b\|_2^2} = \underbrace{(Ax - b)^T (Ax - b)} \quad \leftarrow$$

- A quadratic function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is a second-order multivariate polynomial in  $x$ , that is a function containing a linear combination of all possible monomials of degree at most two—i.e.,

$$\left[ f(x) = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j + \sum_{i=1}^n c_i x_i + d \iff f(x) = x^T A x + \underline{c^T} x + d \right.$$

$$\left[ f_{\text{ls}}(x) = x^T A x - \underbrace{2b^T A x}_c + b^T b \right.$$

$f(x) = x^T Q x$  ←  $Q$  is symmetric for  $n \times n$  **Quadratic Functions**  $a = a^T$  for scalars  $a$

• using properties of **symmetric matrices**, we can express any quadratic function as a quadratic form.

• Fact:  $x^T A x$  is scalar  $\Rightarrow x^T A x = x^T A^T x \Rightarrow x^T A x = \frac{1}{2} x^T \underbrace{(A + A^T)}_{=: H} x$

•  $A = \underbrace{\frac{1}{2} (A + A^T)}_{\text{symmetric}} + \underbrace{\frac{1}{2} (A - A^T)}_{\text{anti-symmetric}}$

More in Mod 4

• Hence,  $f(x) = x^T A x + c^T x + d = \frac{1}{2} x^T H x + c^T x + d$

$$= \frac{1}{2} \begin{bmatrix} x \\ 1 \end{bmatrix}^T \underbrace{\begin{bmatrix} H & c \\ c^T & 2d \end{bmatrix}}_Q \begin{bmatrix} x \\ 1 \end{bmatrix}$$

# Symmetric Matrices: Eigendecomposition (Spectral Theorem)

- Every symmetric matrix  $A$  can be **diagonalized** as  $A = V\Lambda V^T$  with  $V$  formed by the **orthonormal eigenvectors of  $A$**  and  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$  a diagonal matrix of the eigenvalues of  $A$

$v_i^T v_j = 0 \quad \forall i \neq j$   
 $\|v_i\| = 1$   
 $(\lambda_i, v_i)$

$$A = \underbrace{\begin{bmatrix} | & \cdots & | \\ v_1 & \cdots & v_n \\ | & \cdots & | \end{bmatrix}}_V \underbrace{\begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \ddots & & \vdots \\ \vdots & & \lambda_{n-1} & 0 \\ 0 & \cdots & 0 & \lambda_n \end{bmatrix}}_\Lambda \underbrace{\begin{bmatrix} - & v_1^T & - \\ \vdots & \vdots & \vdots \\ - & v_n^T & - \end{bmatrix}}_{V^T} \quad v_i v_i^T$$

$\left[ \begin{array}{c} v_i \\ | \\ 1 \end{array} \right] \left[ -v_i \quad \Rightarrow \right]$

or equivalently,  $A = \lambda_1 v_1 v_1^T + \cdots + \lambda_n v_n v_n^T$  (i.e., weighted sum of dyads)

- Additionally,  $VV^T = V^T V = I$  (i.e.  $V^T = V^{-1}$ )
- This factorization property and the fact that  $A$  has  $n$  orthogonal eigenvectors are two **important properties** for a symmetric matrix.

# Example Problem: Eigenvalues are Real

**Problem:** Consider a symmetric matrix  $A$ . Show that the eigenvalues of  $A$  are real.

**Solution.**

• Consider  $Ax = \lambda x$ ,  $x \neq 0$ ,  $x \in \mathbb{C}^n$ ,  $\lambda \in \mathbb{C}$

• Recall:  $\langle x, x \rangle = x^* x = (\bar{x})^T x$

• WTS  $\lambda = \bar{\lambda}$   $\begin{cases} \lambda = a + ib \text{ then } b = 0 \\ \bar{\lambda} = a - ib \end{cases}$

In general for a matrix  $A$   
eigenvalues; vectors can be  
Complex.

$$\lambda \langle x, x \rangle = (\bar{x})^T \underbrace{(\lambda x)}_{Ax} = (\bar{x})^T Ax = (A^T \bar{x})^T x = (A \bar{x})^T x = (\bar{\lambda} \bar{x})^T x = \bar{\lambda} \langle x, x \rangle$$

$$\Leftrightarrow \lambda = \bar{\lambda} \quad \Leftrightarrow \lambda \in \mathbb{R}, \quad x \in \mathbb{R}^n$$

# Example Problem: Orthogonality of Eigenvectors

**Problem:** Consider a symmetric matrix  $A$ . Show that eigenvectors corresponding to distinct eigenvalues are orthogonal.

**Solution.**

$$\left. \begin{array}{l} (\lambda, x) \text{ eigenpair for } A \\ (\mu, z) \text{ — " — } A \end{array} \right\} \text{ where } \mu \neq \lambda$$

$$\lambda \langle x, z \rangle = \langle Ax, z \rangle = (Ax)^T z = x^T A^T z = x^T A z = \mu \langle x, z \rangle$$

$$\text{Since } \mu \neq \lambda \text{ we have that } (\lambda - \mu) \langle x, z \rangle = 0 \Rightarrow \underline{x^T z = 0}$$

# Matrix powers with eigendecomposition

- Recall from **Mod1** we saw many applications with matrix powers such as computing the number of paths of length  $k$  in a graph
- For symmetric matrices, computing matrix powers is easy

$$A = V \underline{\Lambda} V^T, \quad V^T V = I$$

$$A^k = \underbrace{(V \underline{\Lambda} V^T)(V \underline{\Lambda} V^T) \dots (V \underline{\Lambda} V^T)}_{k \text{ times}} = V \underline{\Lambda}^k V^T$$

# Positive Definite Matrices

$$\langle Ax, x \rangle = x^T A^T x > 0$$

- Another important class of matrices are positive definite matrices
- The matrix  $A$  is **positive definite** if  $\langle Ax, x \rangle > 0$ ; sometimes we write  $A \succ 0$
- And,  $A$  is **positive semidefinite (PSD)** if  $\langle Ax, x \rangle \geq 0$ ; sometimes we write  $A \succeq 0$
- Positive definite matrices need not be symmetric, but often we are interested in positive definite symmetric matrices
- Eigenvalues: let  $\lambda_1(A) \geq \dots \geq \lambda_n(A)$  be the order set of eigenvalues of  $A = A^T$

$$\underbrace{f(x)} \quad \underbrace{[PSD]} \quad \underbrace{[A \succeq 0 \iff \lambda_i(A) \geq 0]} \quad \forall i \in \{1, \dots, n\}$$

$$\frac{d}{dx} f(x^*) = 0, \quad \frac{d^2}{dx^2} f(x^*) > 0 \quad \underbrace{[PD]} \quad \underbrace{[A \succ 0 \iff \lambda_i(A) > 0]} \quad \forall i \in \{1, \dots, n\}$$

$$\underbrace{H = \nabla^2 f(x)}, \quad \underbrace{\nabla f(x^*) = 0} \quad ; \quad \underbrace{\text{eigenvalues of } H \text{ positive}} \iff H \succ 0$$



# Example Problem

$$A = A^T$$

**Problem:** Show that  $A \succeq 0 \iff \lambda_i(A) \geq 0, \forall i \in \{1, \dots, n\}$

**Solution.**

$$A = A^T \iff A = V \Lambda V^T \iff x^T A x = \underbrace{x^T V}_{=: z^T} \Lambda \underbrace{V^T x}_z = z^T \Lambda z = \sum_{i=1}^n \lambda_i(A) z_i^2 \quad (1)$$

$$\iff \left[ x^T A x \geq 0 \quad \forall x \in \mathbb{R}^n \iff z^T \Lambda z \geq 0 \quad \forall z \in \mathbb{R}^n \right]$$



$$\lambda_i(A) \geq 0 \quad \forall i \in \{1, \dots, n\}$$

due to (1)



# Examples of PSD Matrices from ML and OPT

- **Mod2-L4:** we often use kernel matrices  $K = [K(x^{(i)}, x^{(j)})]$  and these are symmetric and in general PSD
- **Mod2-L2:** Gram matrices  $A^T A$  and  $AA^T$  are PSD

**Problem:** Show that Gram and Kernel matrices are PSD.

**solution.**

$$G_{ij} = K(x^{(i)}, x^{(j)}) = \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$$

$$[v^T G v \geq 0 \quad \forall v]$$

$$v^T G v = \sum_{i=1}^n \sum_{j=1}^n v_i v_j G_{ij} = \sum_{i=1}^n \sum_{j=1}^n v_i v_j \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$$

$$= \left\langle \sum_{i=1}^n v_i \phi(x^{(i)}), \sum_{j=1}^n v_j \phi(x^{(j)}) \right\rangle = \left\| \sum_{i=1}^n v_i \phi(x^{(i)}) \right\|^2 \geq 0$$

# Example Problem

$$A = A^T$$

**Problem.** Show that a matrix  $A$  is PSD if and only if  $A = B^T B$  for some real matrix  $B$ .

**Solution.**

(a)

eigendecomp

(b)

(a  $\Rightarrow$  b) : Suppose  $A = A^T$  is PSD.  $A = V \Lambda V^T \Leftrightarrow AV = V \Lambda$

Set  $B := \sqrt{\Lambda} V^T$  where  $\sqrt{\Lambda} = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n})$

which exists because  $\lambda_i(A) \geq 0$

$$\text{Hence } B^T B = V \sqrt{\Lambda} \sqrt{\Lambda} V^T = V \Lambda V^T = A \underbrace{V V^T}_{=I} = A$$

(b  $\Rightarrow$  a)  $A = B^T B$ . for any vector  $v$  w/  $v^T A v \geq 0$

$$v^T A v = v^T \underbrace{B^T B}_{=I} v = \langle Bv, Bv \rangle = \|Bv\|^2 \geq 0 \Rightarrow A \text{ is PSD.}$$

# Part 3. SVD

# Overview

- We just talked about special classes of matrices that have a nice decomposition in terms of their eigenvalues—namely, symmetric PSD matrices.
- Now, we will talk about a matrix decomposition that every matrix has—i.e., **SVD**
- And, it is fundamentally related to a key ML analysis tool: **PCA**

# Matrix Decomposition

- Matrix decomposition, also known as matrix factorization, involves describing a given matrix using its constituent elements.
- Recall that you saw QR decomposition in **Module 1** and then its use in **Module 2** (e.g., solving least squares, in particular sparse problems)
- Perhaps the most known and widely used matrix decomposition method is the **Singular-Value Decomposition**, or **SVD**.
- All matrices have an **SVD**, which makes it more stable than other methods, such as the eigen-decomposition.
- We will see the **SVD** is useful for computing the pseudoinverse efficiently and for dimensionality reduction

# Singular Value Decomposition

# What is SVD?

- One can generalize eigenvalues/vectors to non-square matrices, in which case they are called **singular vectors** and **singular values**.
- The **SVD** is a unique matrix decomposition that exists for every matrix  $A \in \mathbb{R}^{m \times n}$ :

$$A = U\Sigma V^{\top}$$

where  $U \in \mathbb{R}^{m \times m}$  and  $V \in \mathbb{R}^{n \times n}$  are *unitary* matrices, and  $\Sigma \in \mathbb{R}^{m \times n}$  is a matrix with non-negative entries on the diagonal and zeros on the off diagonal.

- Unitary:  $UU^{\top} = I$  and  $VV^{\top} = I$



# SVD as a Dyadic Expansion

An equivalent way to express the **SVD**  $A = U\Sigma V^T$  is as a dyadic expansion:

- That is, the **SVD** expresses  $A$  as a nonnegative linear combination of  $\min\{m, n\}$  rank-1 matrices
- the singular values provide the multipliers
- the outer products of the left and right singular vectors provide the rank-1 matrices.

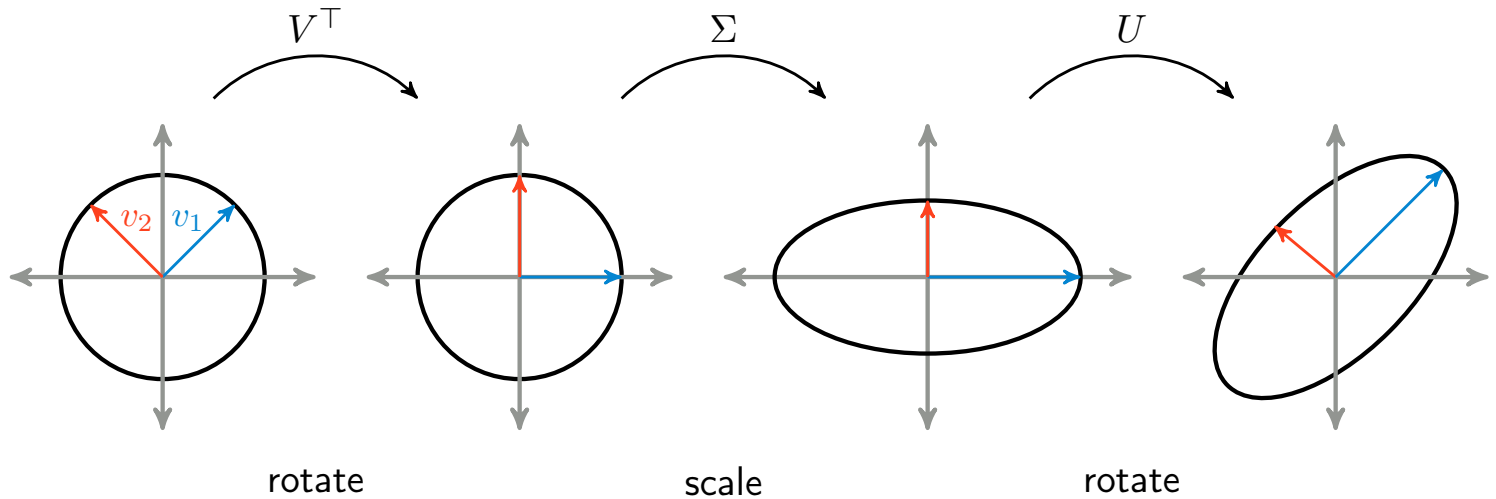
# SVD

- The diagonal entries of  $\Sigma$  are called the **singular values** of  $A$
- The column vectors of  $V$  are called the **right singular vectors** of  $A$
- The column vectors of  $U$  are called the **left singular vectors** of  $A$ .
- The number of nonzero singular values is equal to the rank of the matrix  $A$ .

A diagram illustrating the SVD decomposition of a matrix  $A$ . On the left, a brown rectangle labeled  $A$  is shown with a bracket underneath indicating its dimensions are  $m \times n$ . This is followed by an equals sign. To the right of the equals sign are four components: a red square labeled  $U$  with a bracket underneath indicating its dimensions are  $m \times m$ ; a yellow rectangle labeled  $\Sigma$  with a bracket underneath indicating its dimensions are  $m \times n$ ; and a blue square labeled  $V^T$  with a bracket underneath indicating its dimensions are  $n \times n$ .

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \ddots & 0 \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \sigma_n \\ 0 & \cdots & 0 & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & \cdots & 0 \end{bmatrix}$$

# Geometric View of SVD



# Unpacking the SVD

- Let  $A \in \mathbb{R}^{m \times n}$
- **Fact 1.** Both  $A^T A \in \mathbb{R}^{n \times n}$  and  $AA^T \in \mathbb{R}^{m \times m}$  are symmetric square matrices:
  
- **Fact 2.** Both  $A^T A$  and  $AA^T$  share the same non-zero eigenvalues:

# Unpacking the SVD

- According to the orthogonally diagonalizable property of symmetric matrices, the matrices  $A^T A$  and  $AA^T$  can be decomposed as following:
  
  
  
  
  
  
  
  
  
  
  
  
  
  
  
- **How to obtain the SVD?**: Compute by diagonalizing the PSD symmetric matrices  $A^T A$  and  $AA^T$

# Using the SVD to Compute Pseudo Inverses

- It turns out that using the **SVD** we have a very easy way to compute the pseudo-inverse of  $A$ —i.e.,  $A^\dagger = (A^\top A)^{-1} A^\top$  which we saw in **Mod1 & Mod2**

# Matrix Norms and Connections to Singular values

- Matrix norms and singular values have special relationships.
- **Forbenius Norm:**

$$\|A\|_F = \left( \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{1/2} = (\text{Tr}(A^\top A))^{1/2}$$

- **Matrix  $p$ -norm:** matrix  $p$ -Norm is defined as the largest scalar that you can get for a unit vector

$$\|A\|_p = \sup_{x \neq 0} \frac{\|Ax\|_p}{\|x\|_p} = \max_{\|x\|_p=1} \|Ax\|_p$$

- **Aside:** supremum  $\sup(\cdot)$  is the “least upper bound” of its argument

# Matrix Norms: Spectral Norm

- **Spectral Norm (Matrix 2-norm):** Largest singular value of the matrix  $\sigma_1(A)$

- **Fact:** show that  $\|A\|_F = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2(A)}$  using the fact that  $\|A\|_F = \sqrt{\text{Tr}(A^\top A)}$



# Reduced SVD & Low Rank Approximation

- Rank of  $\Lambda$  is  $r \implies$  there are  $r$  non-zero eigenvalues of the matrices  $A^T A$  and  $AA^T$
- Reduced **SVD**:

# Low Rank Structure

The diagram illustrates the low-rank structure of matrix  $A$ . It shows the equation  $A = Y Z^T$ . Matrix  $A$  is represented by a large orange rectangle with the label  $A$  in the center. Below it is a bracket indicating its dimensions are  $m \times n$ . Matrix  $Y$  is a smaller, vertical red rectangle with the label  $Y$  in the center. Below it is a bracket indicating its dimensions are  $m \times r$ . Matrix  $Z^T$  is a yellow rectangle with the label  $Z^T$  in the center. Below it is a bracket indicating its dimensions are  $r \times n$ . An equals sign is placed between  $A$  and  $Y Z^T$ .

$$A = Y Z^T$$

$m \times n$        $m \times r$        $r \times n$

# Low Rank Structure

$$A = uv^{\top} =$$

$$A = uv^{\top} + wz^{\top} =$$

# Low Rank Approximation

- Low Rank Approximation: take only top  $k$ -singular values and corresponding dyads in the dyadic expansion
  
- Low Rank Approximation is an important tool for many applications including
  - ▶ Linear system identification: approximating matrix is Hankel structured. (You saw this in M2-N2.ipynb)
  - ▶ ML: feature space dimensionality reduction
  - ▶ Recommender systems: matrix completion
  - ▶ Distance matrix completion where there is a positive definiteness constraint.
  - ▶ Natural language processing where the approximation is non-negative.
  - ▶ Image or video compression

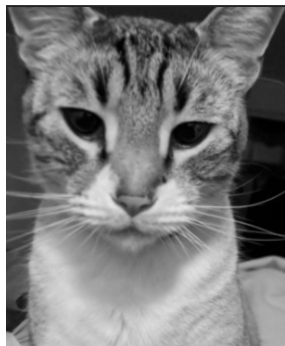
# Example: Compression

- **Compression.** A low-rank approximation provides a (lossy) compressed version of the data matrix.
  - ▶ The original matrix  $A$  is described by  $mn$  numbers, while describing  $Y$  and  $Z^T$  requires only  $k(m+n)$  numbers.
  - ▶ When  $k$  is small relative to  $m$  and  $n$ , replacing the product of  $m$  and  $n$  by their sum is a big win.
  - ▶ With images, a modest value of  $k$  (say 100 or 150) is usually enough to achieve approximations that look a lot like the original image.

bender in color



gray scale



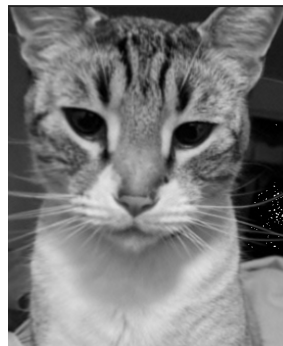
rank 2 bender



rank 20 bender



rank 100 bender



# Optimality of Low Rank Approximation

- The low rank approximation obtained via **SVD** is optimal in the following sense.
- Recall the Forbenius norm:
  - i.e.,  $\ell_2$ -norm (i.e., usual Euclidean norm) applied to the matrix as if it were a vector
  - **Theorem [Eckat-Young-Mirsky].**

# How to choose $k$ ?

- When producing a low-rank matrix approximation, we have been taking as a parameter the target rank  $k$ .
- **Ideal Setting:** the singular values of  $A$  give strong guidance
  - ▶ if the top few singular values are big and the rest are small, then the obvious solution is to take  $k$  equal to the number of "big values".
- **Less Ideal Setting:** take  $k$  as small as possible subject to obtaining a useful approximation, where what "useful" means depends on the application.
  - ▶ e.g., a common rule of thumb is to choose  $k$  such that the sum of the top  $k$  singular values is at least  $c$  times as big as the sum of the other singular values, where  $c$  is a domain-dependent constant (like 10, say).

# Next Up

- Next lecture we will talk about PCA, and show that PCA reduces to SVD and is fundamentally connected to low rank approximations.