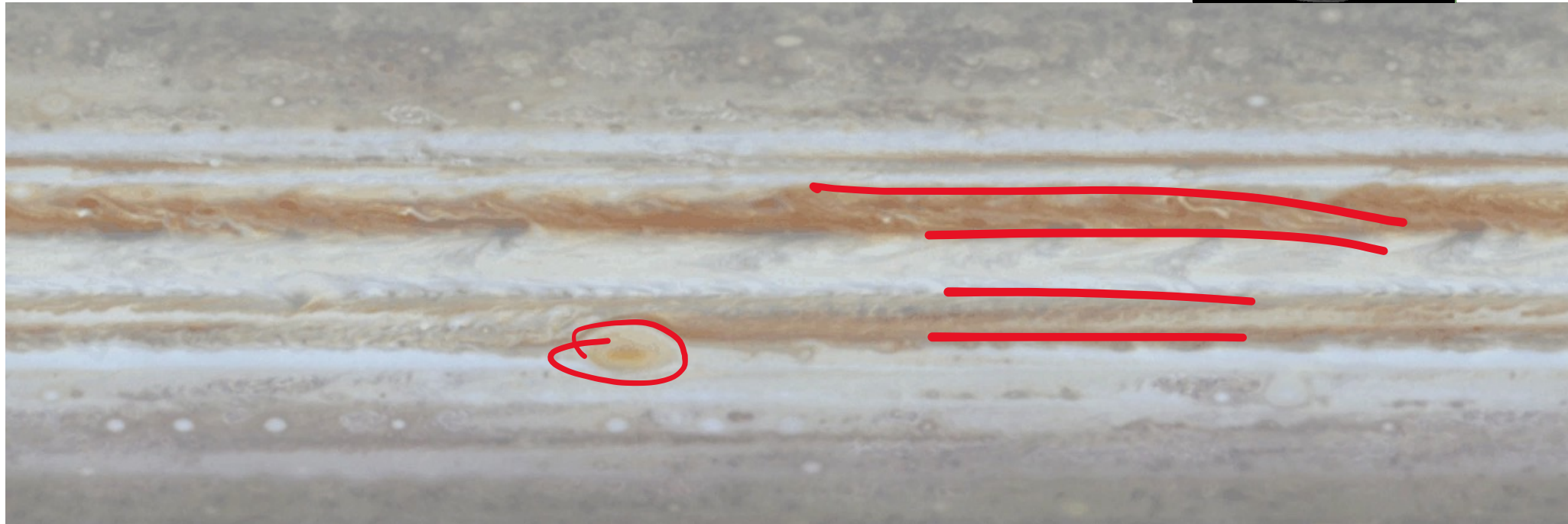
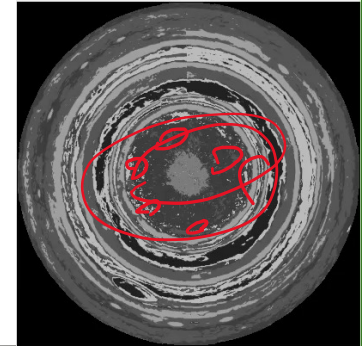
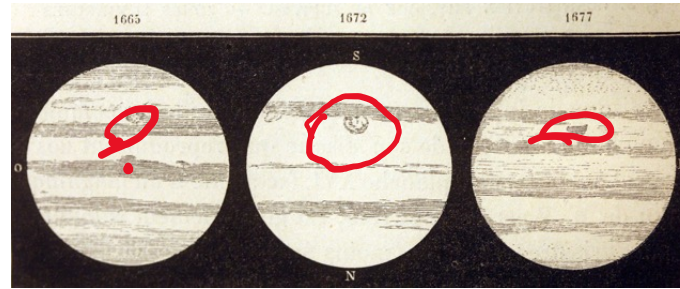
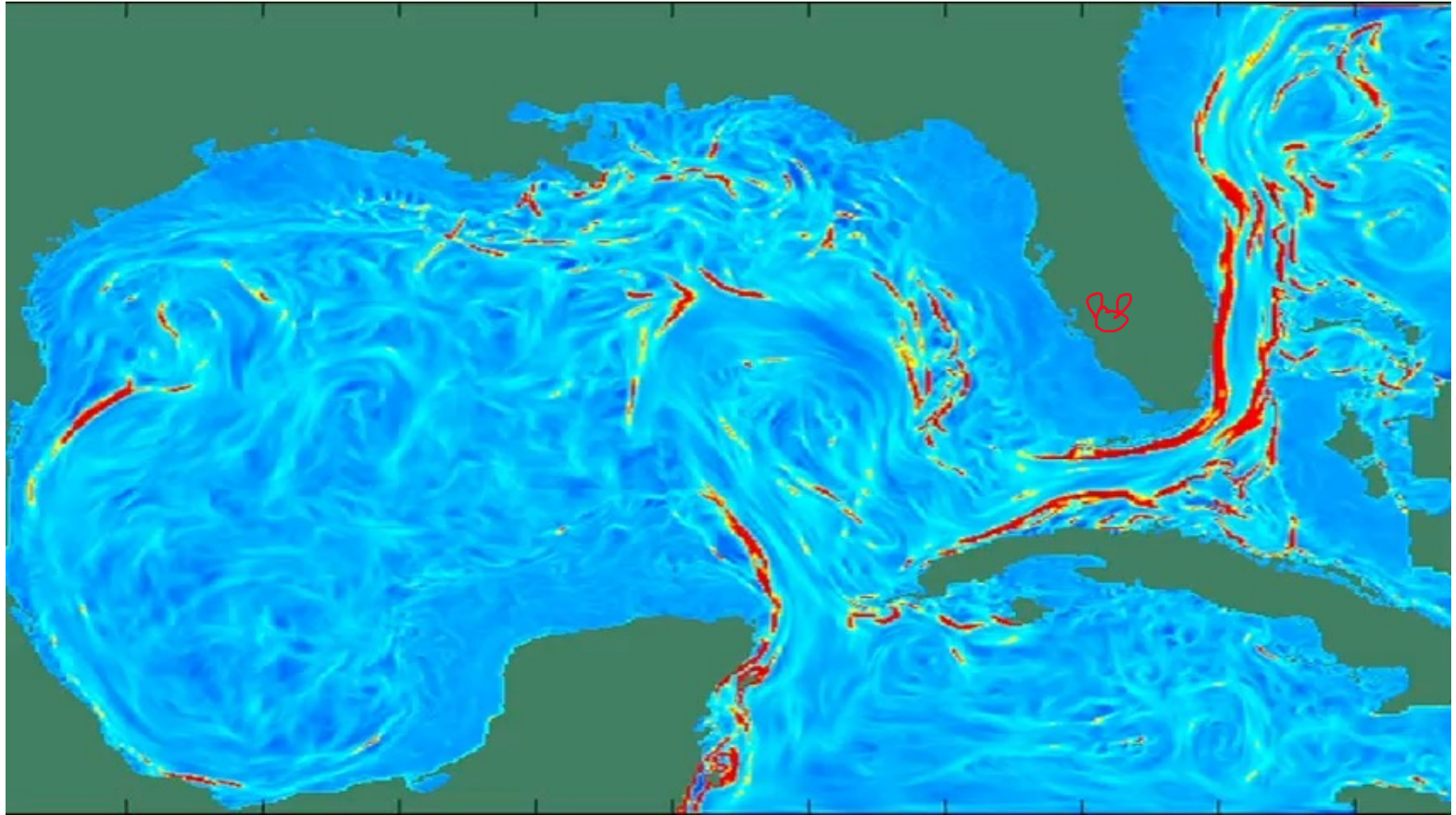


# EE520 Data Driven Analysis of Complex Systems

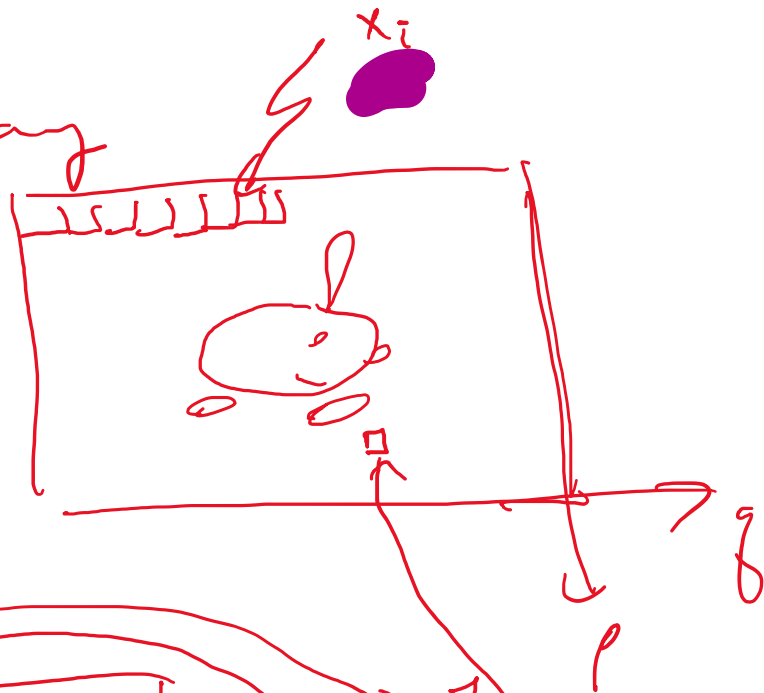
Erik Bollt





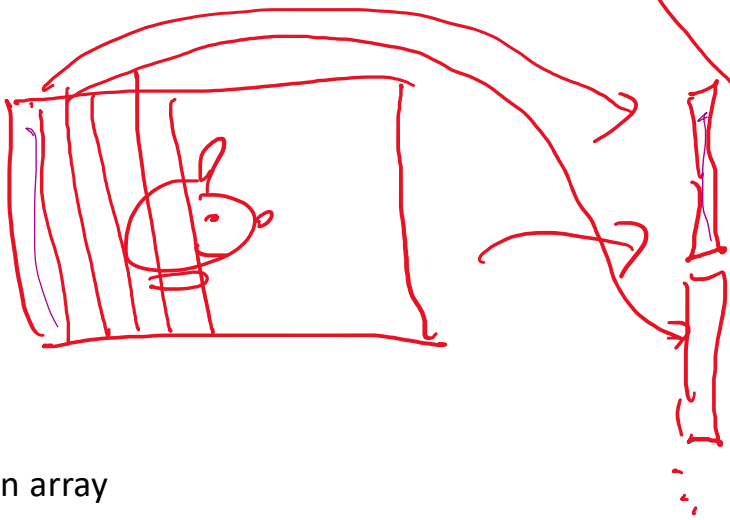
Data as an array

$X =$

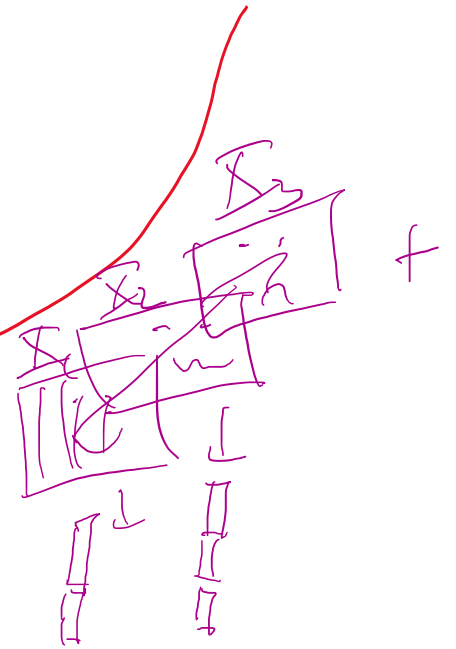


$X_{p \times q} \in \mathbb{R}^{p \times q}$  matrix

$[X]_{i,m} = \text{one pixel}$



slice & stack



Data as an array

# On Matrix Multiplication

$$L(\mathbf{z}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

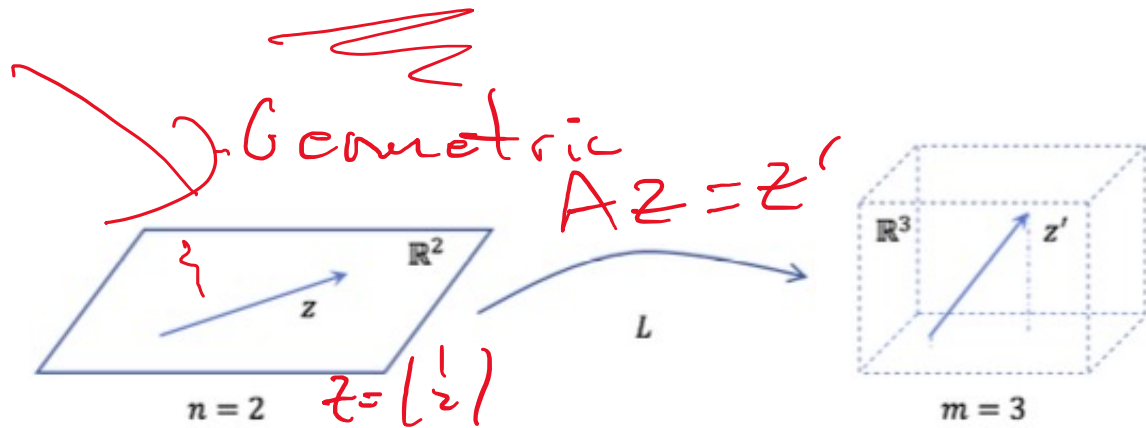
$$\mathbf{z} \mapsto \mathbf{z}' = A\mathbf{z}$$

in terms of the usual matrix times vector multiplication,

$$[\mathbf{z}']_i = \sum_{j=1}^n A_{i,j} [\mathbf{z}]_j, \text{ for each } i = 1, \dots, m,$$

• a vector has length & direction.

$$A = \begin{pmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{pmatrix}, \text{ and each } a_{i,j} \in \mathbb{C}$$

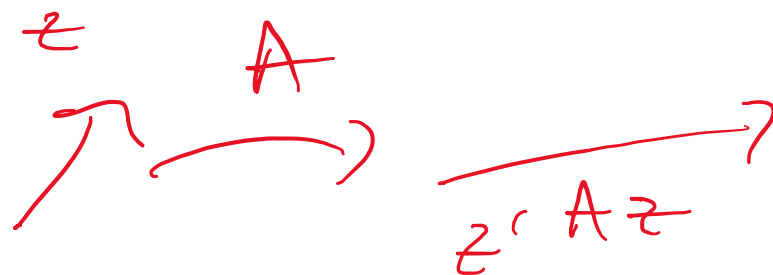


But as linear algebra

$$A_{2 \times 2} = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}; \text{ matrix } \times \text{ vectors}$$

$$A \begin{pmatrix} 3 \\ 4 \end{pmatrix} = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \begin{pmatrix} 3 \\ 4 \end{pmatrix} = \begin{pmatrix} 1 \cdot 3 + 2 \cdot 4 \\ 3 \cdot 3 + 4 \cdot 4 \end{pmatrix} = \begin{pmatrix} 11 \\ 25 \end{pmatrix}$$

$$z' = Az$$

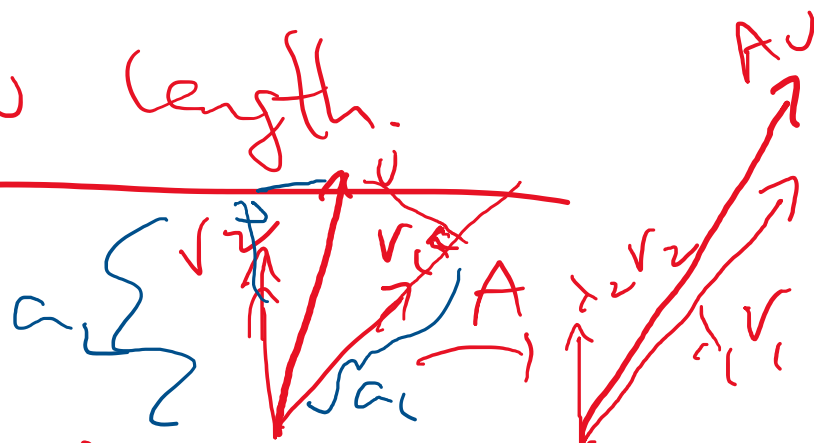


new direction, new length.

Eig for square -

$$AV = \lambda V$$

$2 \times 2$



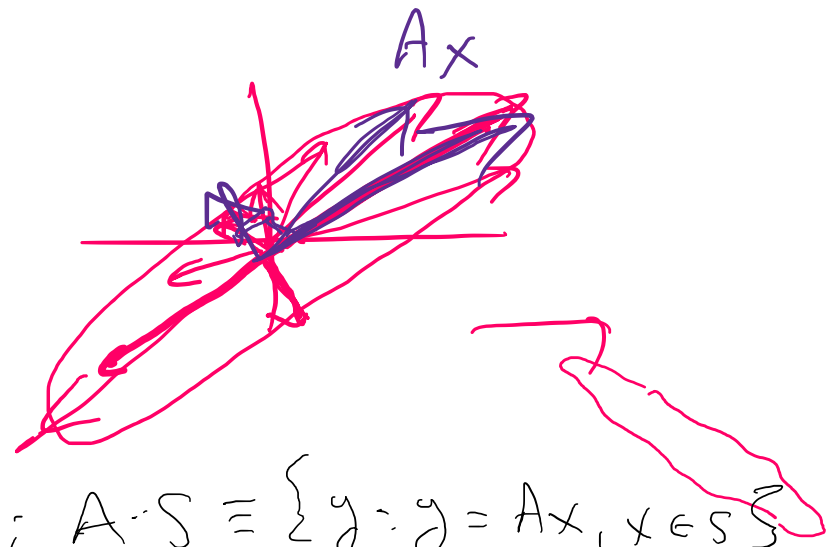
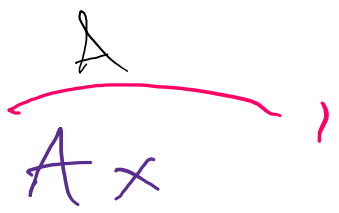
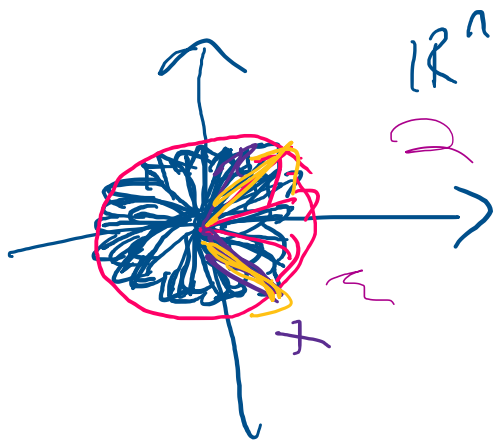
Characterize matrices by knowing just these

- Eig. special directions

$$\begin{aligned}
 \bullet Av &= A(a_1v_1 + a_2v_2) = a_1Av_1 + a_2Av_2 & \det(A - \lambda I) &= 0 \\
 &= a_1\lambda_1v_1 + a_2\lambda_2v_2 & (A - \lambda I)v &= 0
 \end{aligned}$$

Matrix time circle =  
*all vectors of length 1.*

? Matrix  $\times$  circle ?! But matrix  
 times vector.  $\rightarrow =$



$$S = \{x \mid \|x\|_2 = 1, x \in E \cong \mathbb{R}^2\}; \quad A \cdot S = \{y \mid y = Ax, x \in S\}$$

**Theorem 2.1.1 — Singular Value Decomposition.** Let  $A$  be an  $m \times n$  matrix whose entries come from the field  $\mathcal{K}$ , which is either the field of real numbers or the field of complex numbers. Then the singular value decomposition of  $A$  exists, and it takes the form of a product of matrices:

$$A_{m \times n} = U_{m \times m} \Sigma_{m \times n} V_{n \times n}^* \quad (2.5)$$

where

- $U$  is an  $m \times m$  unitary matrix.
- $\Sigma$  is a diagonal  $m \times n$  matrix with non-negative real numbers on the diagonal.
- $V$  is an  $n \times n$  unitary matrix, and  $V^*$  is the conjugate transpose of  $V$ .

The singular values are the nonnegative values:  $\sigma_i \geq 0, i = 1, \dots, n$ ,

The left singular vectors:  $u_i$  are the columns of  $U = [u_1, u_2, \dots, u_m]$ .

The right singular vectors:  $v_i$  are the columns of  $V = [v_1, v_2, \dots, v_n]$ .

**Definition 2.1.1 — Singular values and singular vectors.** The singular values of  $A$  are the scalar values,  $\sigma_i$ , and the columns of  $U$  and  $V$  have columns that are the corresponding  $i^{\text{th}}$  left and right singular vectors,  $u_i$  and  $v_i$ :

The singular values are the nonnegative values:  $\sigma_i \geq 0, i = 1, \dots, n$ ,

The left singular vectors:  $u_i$  are the columns of  $U = [u_1, u_2, \dots, u_m]$ .

The right singular vectors:  $v_i$  are the columns of  $V = [v_1, v_2, \dots, v_n]$ .

Since  $V$  is orthogonal, then right multiplying Eq. (2.5) by  $V$ ,

$$AV = U\Sigma V^*V = U\Sigma, \quad (2.8)$$

$$\Sigma := \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_p), p = \min(m, n),$$

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0.$$

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

$$u_i^T u_j = \delta_{ij}$$

$$u_i^T u_j = \begin{pmatrix} u_{i1} & u_{i2} & \dots & u_{im} \end{pmatrix} \begin{pmatrix} u_{j1} \\ u_{j2} \\ \dots \\ u_{jm} \end{pmatrix} = \delta_{ij} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 \end{pmatrix}$$

$$U^T U = I$$

$$U \text{ unitary} \Leftrightarrow U^* U = U U^* = I$$

$$U^T U = \begin{pmatrix} u_1^T & u_2^T & \dots & u_m^T \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ \dots \\ u_m \end{pmatrix} = I = \begin{pmatrix} 1 & & & \\ & 1 & & \\ & & \dots & \\ & & & 1 \end{pmatrix}$$





Full \*

### The Economy SVD, and Reduced Rank SVD

The general SVD, Eq. (2.5) may be written in terms of submatrices.

**Definition 2.1.3 — The Economy SVD.** For any matrix  $A \in \mathbb{R}^{m \times n}$ , the general SVD Eq. (2.5) can be written in terms of smaller matrices,

$$A_{m \times n} = \hat{U}_{m \times n} \hat{\Sigma}_{n \times n} V_{n \times n}^* \quad (2.21)$$

and  $U = [\hat{U}_{m \times n} | \hat{U}_{(n-m) \times n}]$ , written in terms of an orthogonal "buffer" matrix

**Definition 2.1.4 — Rank Deficient SVD.** For a matrix  $A \in \mathbb{R}^{m \times n}$  such that the SVD results in singular values

$$\sigma_r > \sigma_{r+1} = 0, \text{ for some } r < n. \quad (2.22)$$

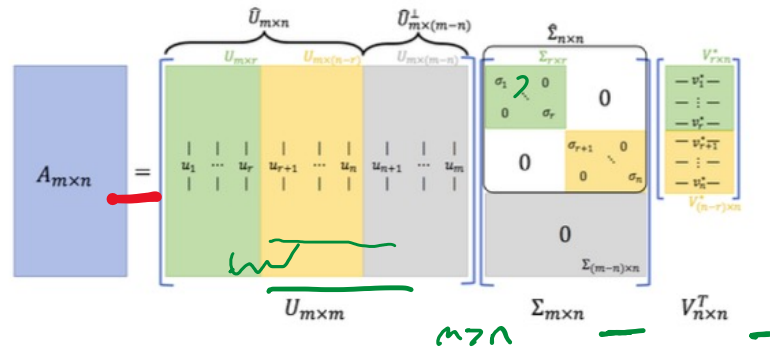
then the SVD can be written in terms of an economy form as smaller matrices,

$$A_{m \times n} = \hat{U}_{m \times r} \hat{\Sigma}_{n \times n} V_{n \times r}^* \quad (2.23)$$

and related to the general SVD Eq. (2.5) by  $U = [\hat{U}_{m \times r} | \hat{U}_{(n-r) \times n}]$ , but  $r < n$ .

$$\sigma_1 > \sigma_2 > \dots > \sigma_r > \sigma_{r+1} = 0 \\ = 0 \dots$$

Rank ill conditioned



Full,  
Economy,  
Truncated  
SVD

Figure 2.3:  $m > n$  tall skinny

Recall that,

$$\begin{aligned}
 A_{m \times n} &= \hat{U}_{m \times n} \hat{\Sigma}_{n \times n} \hat{V}_{n \times n}^T \\
 &= \begin{bmatrix} | & | & | & | \\ u_1 & u_2 & \dots & u_n \\ | & | & | & | \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n \end{bmatrix} \begin{bmatrix} - & v_1^T & - \\ - & v_2^T & - \\ - & \vdots & - \\ - & v_n^T & - \end{bmatrix} \quad (2.24)
 \end{aligned}$$

but  $V^T V = I$ , orthogonality allows:

$$A_{m \times n} \hat{V}_{n \times n} = \hat{U}_{m \times n} \hat{\Sigma}_{n \times n} \quad (2.25)$$

so,

$$A_{m \times n} \begin{bmatrix} | & | & | & | \\ v_1 & v_2 & \dots & v_n \\ | & | & | & | \end{bmatrix} = \begin{bmatrix} | & | & | & | \\ u_1 & u_2 & \dots & u_n \\ | & | & | & | \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n \end{bmatrix} \quad (2.26)$$

but this just states  $n$ -matrix times vector statements:

$$\begin{aligned}
 Av_1 &= \sigma_1 u_1 \\
 Av_2 &= \sigma_2 u_2 \\
 &\vdots \\
 Av_n &= \sigma_n u_n
 \end{aligned} \quad (2.27)$$

Handwritten notes in green ink:

- $\sigma_{r+1} = 0$
- $\sigma_r = 10^{-6}$
- $\sigma_{r+1} < 10^{-6} < \epsilon$

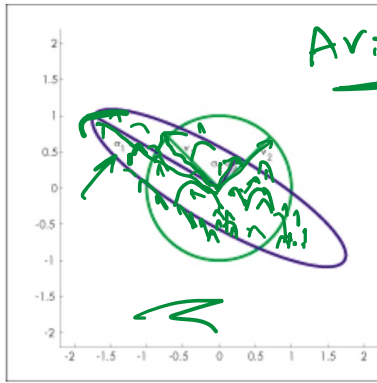
$$A = \underline{U} \underline{\Sigma} \underline{V}^T$$

$$A \underline{V} = \underline{U} \underline{\Sigma} \underline{e}$$

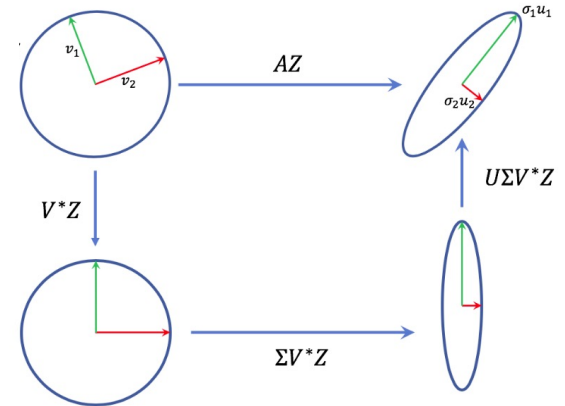
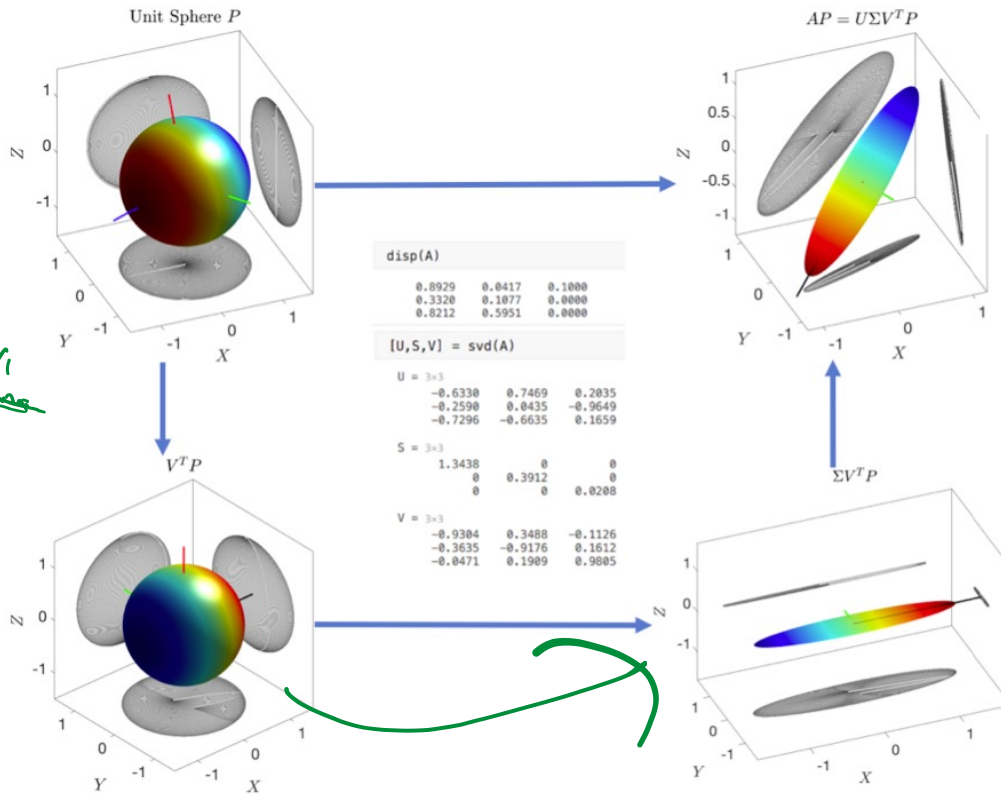
$$A v_i = \sigma_i u_i$$

### Geometry:

1.  $V^*$  rotates to a standard configuration.
2.  $\Sigma$  stretches each orthogonal axis to the major covariance axis of the corresponding ellipsoid, and
3.  $U$  rotates results back to the configuration that associates with  $A$ .

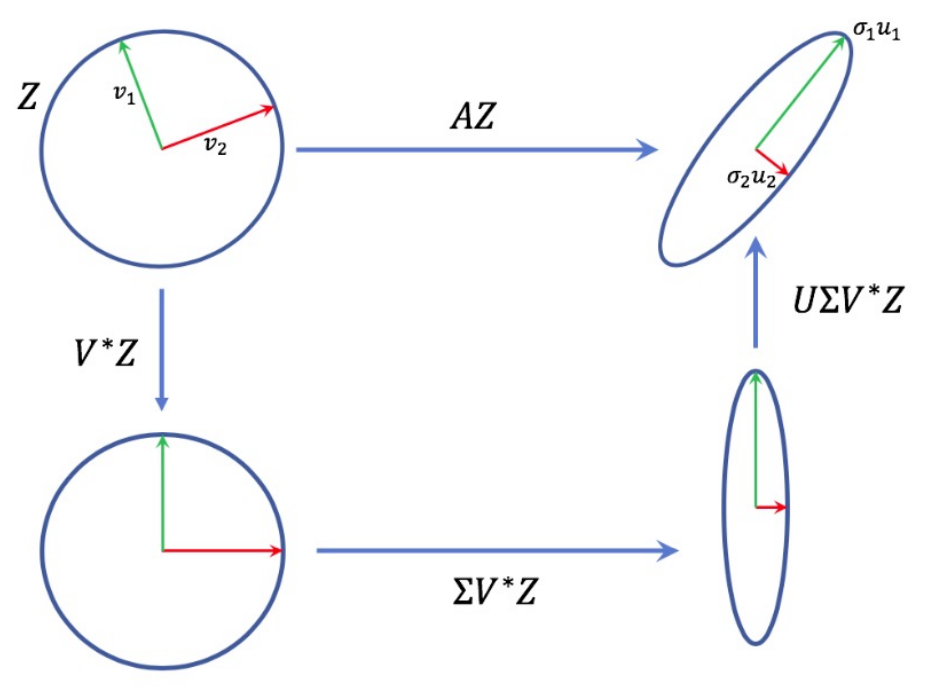
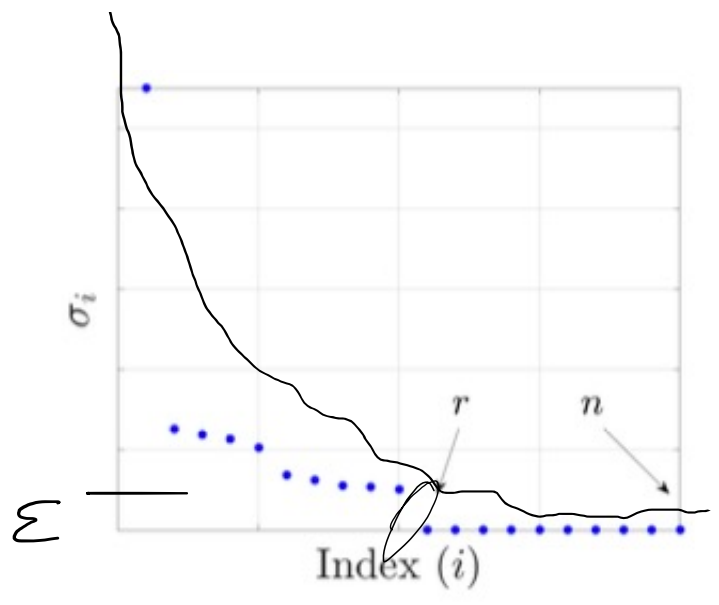


$b_1 \sim 1.5$   
 $b_2 = 0.0003$   
 $b_3 = 2$   
 $b_4 = 0.0003$



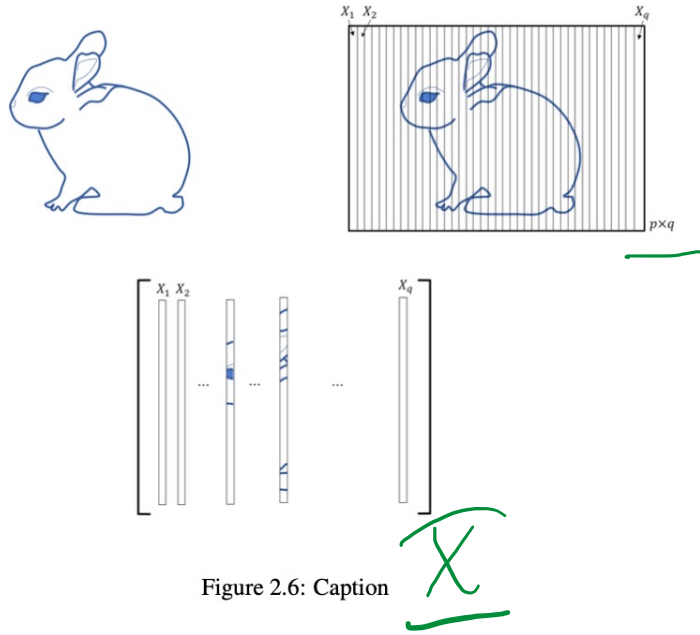
$A w = U \underline{\Sigma} V^T w$   
 $v_i^T w = v_i \cdot w$   
 $\rightarrow \begin{pmatrix} X \\ Y \\ Z \end{pmatrix} = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 4 & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \end{pmatrix}$   
~~redundant~~  
 if  $b_i \approx 0$

And ROM



$\mathcal{B} \subset \mathcal{E} \rightarrow \mathcal{B} \subset \mathcal{H}$   
 $\sim$   
 $0$

# Bunny Compression



Covariance – notice the demean step

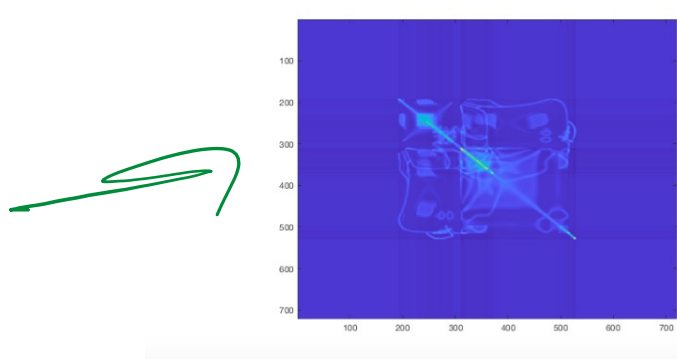
$$C_I = \frac{1}{n-1} (X - \bar{X})^T (X - \bar{X})$$

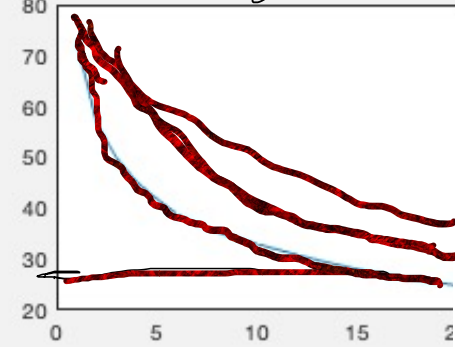
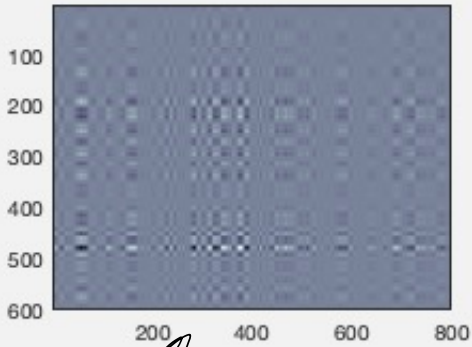
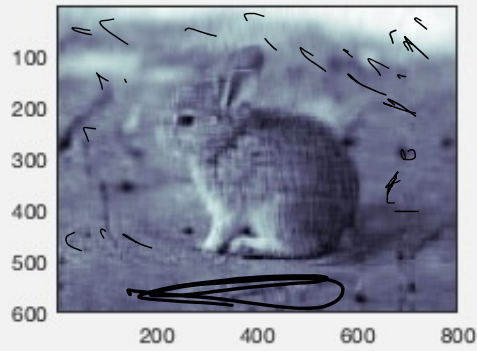
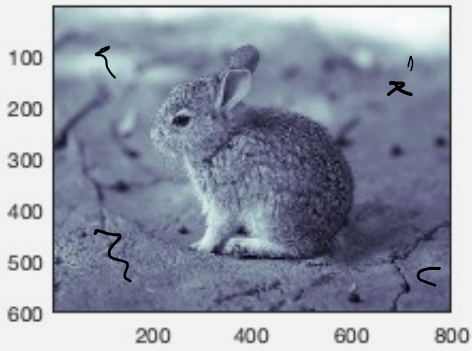
A A

X

X<sup>T</sup>X

X -  $\bar{X}$





800x150

$$U(x,t) = \sum a_i(t) \phi_i(x)$$

$|a_i(t)| \xrightarrow{i \rightarrow \infty} 0$   
 as fast as possible.

$$I = U \Sigma^T$$

$\approx \sum \lambda_i u_i v_i^T$

```

1 I = imread('Bunny.jpg');
2
3 figure
4 subplot(1,2,1)
5 imshow(I)
6 xticks({}); yticks({});
7 pbaspect([1 1 1])
8 title('RGB Image')
9
10 I = rgb2gray(I); %Convert the 3D RGB color to 1D grayscale
11 I = im2double(I); %Convert integer value to double (scaled ...
    from 0 to 1)
12
13 subplot(1,2,2)
14 imshow(I)
15 xticks({}); yticks({});
16 pbaspect([1 1 1])
17 title('Grayscale Image')
  
```

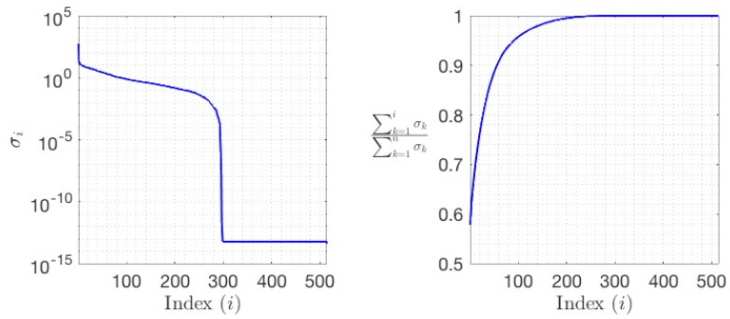
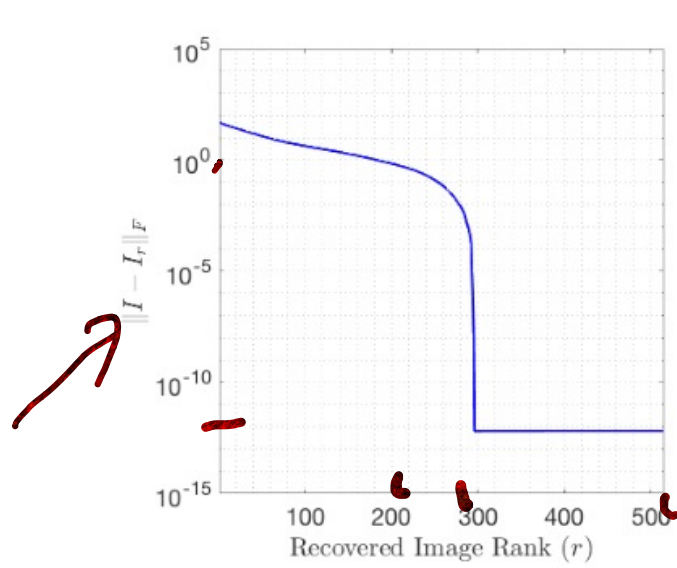
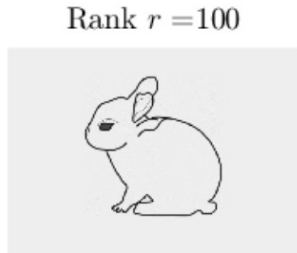


Figure 2.8: (Left) Singular Values. (Right) Energy



$$e_{\text{ref}} = \frac{\|I - I_r\|_F}{\|I\|_F}$$



: Distance  $\|I - I_r\|_F$ , where  $I_r$  is the recovered image using the reduced

Code 2.1: Read, convert, and display images.

```
1 I = imread('Bunny.jpg');
2
3 figure
4 subplot(1,2,1)
5 imshow(I)
6 xticks({}); yticks({});
7 paspect([1 1 1])
8 title('RGB Image')
9
10 I = rgb2gray(I); %Convert the 3D RGB color to 1D grayscale
11 I = im2double(I); %Convert integer value to double (scaled ...
    from 0 to 1)
12
13 subplot(1,2,2)
14 imshow(I)
15 xticks({}); yticks({});
16 paspect([1 1 1])
17 title('Grayscale Image')
```



## History



Gene Golub's license plate, photographed by Professor P. M. Kroonenberg of Leiden University. Gene Howard Golub (February 29, 1932 – November 16, 2007), Fletcher Jones Professor of Computer Science at Stanford University. His work made fundamental contributions that have made the singular value decomposition practical as one of the most powerful and widely used tools in modern matrix computation.

Lots of Machine Learning

& Data Analysis

is solving an ill-posed

- optimize a cost function.

$$AA^T = U \Sigma V^T (V \Sigma^T U^T)$$
$$= U \Sigma \Sigma^T U^T$$

$$(AA^T) \underline{U} = U (\Sigma \Sigma^T) = (\Sigma \Sigma^T) \underline{U}$$

$$\underline{U} = \begin{pmatrix} \underline{u}_1 & \underline{u}_2 & \dots & \underline{u}_m \end{pmatrix}$$

**Definition 2.1.2 — Induced Norm.** Suppose a vector norm  $\|\cdot\|$  on  $\mathcal{K}^m$  is given. Any matrix  $A_{m \times n}$  induces a linear operator from  $\mathcal{K}^n$  to  $\mathcal{K}^m$  with respect to the standard basis, and one defines the corresponding induced norm or operator norm on the space  $\mathcal{K}^{m \times n}$  of all  $m \times n$  matrices as follows:

$$\|A\|_p = \sup_{x \neq 0} \frac{\|Ax\|_p}{\|x\|_p} \quad (2.14)$$

or, taking a vector  $x$  such that  $\|x\|_p = 1$ , then we have

$$\|A\|_p = \sup_{\|x\|_p=1} \|Ax\|_p \quad (2.15)$$

### Some Special (Simple) Matrix Norms

The first 3 of these are induced norms, but the 4th is not.

- For  $p = 1$ :

$$\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}| \quad (2.16)$$

- For  $p = \infty$ :

$$\|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}| \quad (2.17)$$

- A special case is the spectral norm when  $p = 2$ , in which we have:

$$\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)} = \sigma_{\max} \quad (2.18)$$

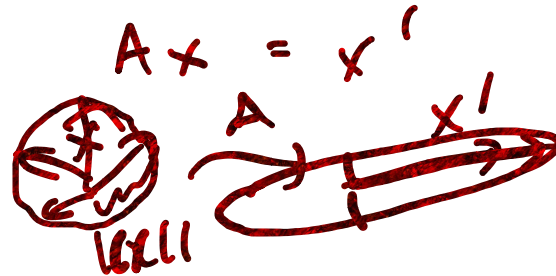
where  $\sigma_{\max}$  is the maximum singular value of the matrix  $A$ .

- The Frobenius norm is given by:

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2} \quad (2.19)$$

**Theorem 2.1.2** For a matrix  $A$ , the product of the singular values of  $A$ , equals the absolute value of its determinant:

$$|\det(A)| = \prod_{i=1}^n \sigma_i \quad (2.20)$$



$$p=1 : \|(x_1, x_2)\|_1 = |x_1| + |x_2|$$

$$p=\infty : \|(x_1, x_2)\|_\infty = \max_i |x_i|$$

$$p=2 : \|(x_1, x_2)\|_2 = \sqrt{x_1^2 + x_2^2}$$

$$x = \begin{pmatrix} 1 \\ 3 \end{pmatrix}; \quad \|x\|_1 = 1 + 3 = 4$$

$$\|x\|_\infty = 3$$

$$\|x\|_2 = \sqrt{1^2 + 3^2} = \sqrt{10}$$

**Definition 2.1.2 — Induced Norm.** Suppose a vector norm  $\|\cdot\|$  on  $\mathcal{K}^m$  is given. Any matrix  $A_{m \times n}$  induces a linear operator from  $\mathcal{K}^n$  to  $\mathcal{K}^m$  with respect to the standard basis, and one defines the corresponding induced norm or operator norm on the space  $\mathcal{K}^{m \times n}$  of all  $m \times n$  matrices as follows:

$$\|A\|_p = \sup_{x \neq 0} \frac{\|Ax\|_p}{\|x\|_p} \quad (2.14)$$

or, taking a vector  $x$  such that  $\|x\|_p = 1$ , then we have

$$\|A\|_p = \sup_{\|x\|_p=1} \|Ax\|_p \quad (2.15)$$

### *Some Special (Simple) Matrix Norms*

The first 3 of these are induced norms, but the 4th is not.

- For  $p = 1$ :

$$\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}| \quad (2.16)$$

- For  $p = \infty$ :

$$\|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}| \quad (2.17)$$

- A special case is the spectral norm when  $p = 2$ , in which we have:

$$\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)} = \sigma_{\max} \quad (2.18)$$

where  $\sigma_{\max}$  is the maximum singular value of the matrix  $A$ .

- The Frobenius norm is given by:

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2} \quad (2.19)$$

**Theorem 2.1.2** For a matrix  $A$ , the product of the singular values of  $A$ , equals the absolute value of its determinant:

$$|\det(A)| = \prod_{i=1}^n \sigma_i \quad (2.20)$$

Fun facts about matrix estimation (data estimation)

If  $A$ ,  $b_1 \geq \dots \geq b_r > b_{r+1} = 0$

•  $\text{range}(A) = \text{span}(u_1, u_2, \dots, u_r)$

•  $\text{null}(A) = \text{span}(v_{r+1}, v_{r+2}, \dots, v_n)$

$\} \begin{matrix} A \\ 0 \end{matrix} \rightarrow$   
 $b_2 = 0$   
 $r = 1$

•  $\|A\|_2 = b_1$  ;  $\|A\|_F = \sqrt{b_1^2 + b_2^2 + \dots + b_r^2}$

•  $A = \sum_{i=1}^r b_i u_i v_i^T = \underbrace{b_1 u_1 v_1^T}_{\text{rank-1}} + \underbrace{b_2 u_2 v_2^T}_{\text{outer products}} + \dots + b_r u_r v_r^T$

$w_1^T w_2 = w_1 \cdot w_2$   
 $= \|w_1\| \|w_2\| \cos \theta$

$A = U \Sigma V^T$

$(u_1 \dots u_r \dots u_n) \begin{pmatrix} b_1 & b_2 & \dots & b_r \end{pmatrix} \begin{pmatrix} -v_1^T \\ -v_2^T \\ \vdots \end{pmatrix}^{1 \times n}$

$\sum_{i=1}^n w_i \rightarrow$   
 $\begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix}$

cos θ

# Matrix Estimation / Data Estimation - $A_{m \times n}$

$A$

Let  $0 \leq N \leq r$  and  $A_N = \sum_{i=1}^N \sigma_i u_i v_i^*$

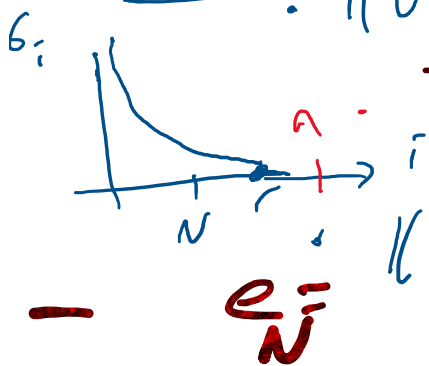
(so we may be skipping some of them ...

$\sum_{i=N+1}^r$

Then

$\|A - A_N\|_2 = \sigma_{N+1}$  (first one skipped)

(what if it zero?)



$\|A - A_N\|_F = \sqrt{\sigma_{N+1}^2 + \sigma_{N+2}^2 + \dots + \sigma_r^2}$

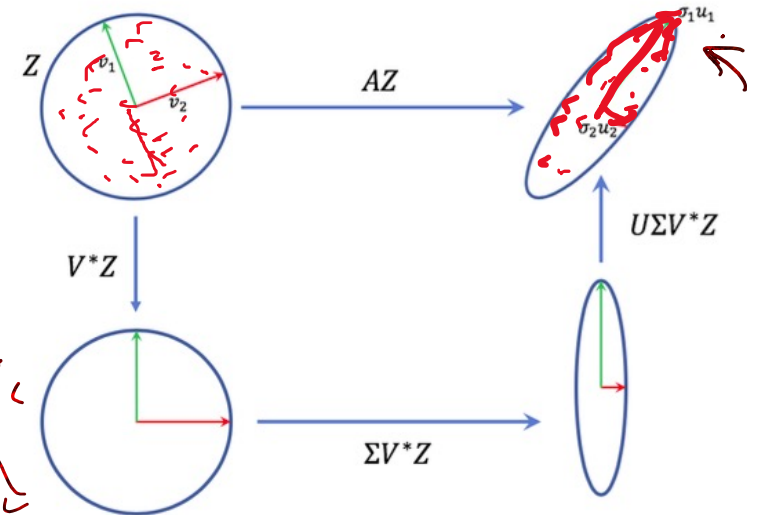


On PCA Principal Component Analysis, Eigenface

-On Raleigh Ritz Quotient

-On Spectral Decomposition Theorem

-On Data Clouds



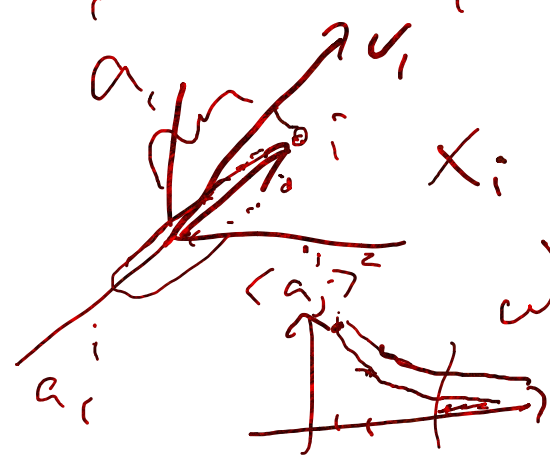
$\bar{X} = [x_1 | x_2 | \dots | x_n]$

$x_i = a_{i1} \vec{v}_1 + a_{i2} \vec{v}_2 + \dots + a_{ir} \vec{v}_r$

$\vec{v}_s$  as basis set.

$x_i = \text{PCA gives basis set where } \langle a_j \rangle_i \downarrow \text{ as fast as possible vs any other basis set.}$

$\langle a_i \rangle_i = \frac{1}{\sqrt{2}} \sum_{i=1}^n a_i$





Data for PCA - "Pretend data looks like an ellipsoid"

Ex.  $x_i \sim 4000 \times 1$  gene expression table for each  $i$ .  
 $i = 1 \dots 216$  patients

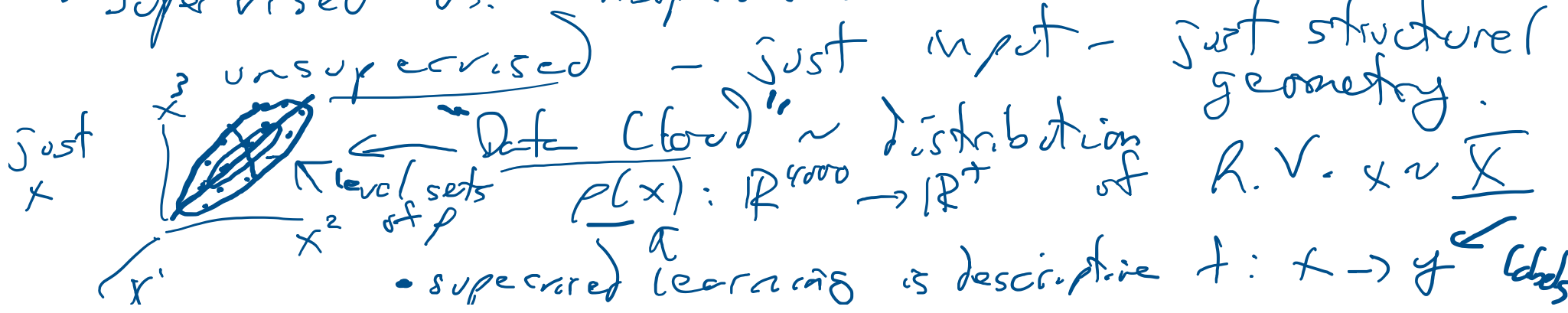
$$x_i = \begin{pmatrix} x_i^1 \\ x_i^2 \\ \vdots \\ x_i^{4000} \end{pmatrix}$$

$y_i = 0$  or  $1$  "0" if not cancer "1" if cancer.



$$\mathbb{Z}_2 = \{0, 1\}$$

Supervised vs. unsupervised.



## THE SPECTRAL DECOMPOSITION

Let  $A$  be a  $n \times n$  symmetric matrix. From the spectral theorem, we know that there is an orthonormal basis  $u_1, \dots, u_n$  of  $\mathbb{R}^n$  such that each  $u_j$  is an eigenvector of  $A$ . Let  $\lambda_j$  be the eigenvalue corresponding to  $u_j$ , that is,

$$A u_j = \lambda_j u_j \quad \leftarrow \text{real}$$

Then

$$A = P D P^{-1} = P D P^T$$

where  $P$  is the orthogonal matrix  $P = [u_1 \ \dots \ u_n]$  and  $D$  is the diagonal matrix with diagonal entries  $\lambda_1, \dots, \lambda_n$ . The equation  $A = P D P^T$  can be rewritten as:

$$\begin{aligned} A &= [u_1 \ \dots \ u_n] \begin{bmatrix} \lambda_1 & & \\ & \dots & \\ & & \lambda_n \end{bmatrix} \begin{bmatrix} u_1^T \\ \vdots \\ u_n^T \end{bmatrix} \\ &= [\lambda_1 u_1 \ \dots \ \lambda_n u_n] \begin{bmatrix} u_1^T \\ \vdots \\ u_n^T \end{bmatrix} \\ &= \lambda_1 u_1 u_1^T + \dots + \lambda_n u_n u_n^T \end{aligned}$$

*order*

$n \times n \cdot (n \times n)$   
 $n \times n$

$$\|v_i\|^2 = v_i \cdot v_i = v_i^T v_i \quad \text{scalar = inner product}$$

The expression

$$A = \lambda_1 u_1 u_1^T + \dots + \lambda_n u_n u_n^T$$

is called the spectral decomposition of  $A$ . Note that each matrix  $u_j u_j^T$  has rank 1 and is the matrix of projection onto the one dimensional subspace spanned by  $u_j$ . In other words, the linear map  $P$  defined by  $P(x) = u_j u_j^T x$  is the orthogonal projection onto the subspace spanned by  $u_j$ .

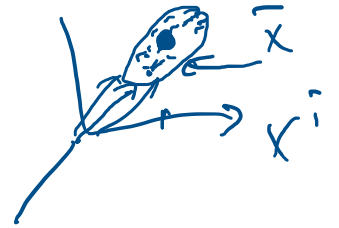
$A = B^T B$  is symmetric  
 $\Rightarrow$  spectral decomp. theorem  
 i.e. also covariance matrices.

$A$  is pos. definite if  
 $\lambda_i > 0$  all  $i$ .

# PCA as algorithm

• Data  $\underline{X} = \begin{pmatrix} x_1 & x_2 & \dots & x_n \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \end{pmatrix}$  <sup>ith</sup>  
 <sub>$m \times n$</sub>

(say  $m = 4000$   
 $n = 216$ )



• what if  $x_i \sim \mathcal{N}(\bar{x}, \Sigma)$  —

covariance matrix.

$$B = \underline{X} - \bar{B}; \quad \bar{B} = \begin{pmatrix} \vdots \\ \vdots \\ \vdots \end{pmatrix}_{n \times 1} \bar{x}^T = \mathbf{1} \bar{x}^T - \bar{X} = \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_n \end{pmatrix}$$

$$\bar{x}_i = \frac{1}{n} \sum_{j=1}^n \underline{X}_{ij}$$



$$B = U \Sigma V^T; \quad U = [u_1, u_2, \dots, u_n]$$

and  $u_1$  is major axis — most energetic — feature that most explains data.  
 $u_2$  is first minor axis

Side note:

$\frac{1}{f^2}$  is slowest converging to zero  $b_i^2$

i.e.  $\frac{1}{f}$  is slowest converging  $b_i$

$$b_i \leq \frac{1}{i}$$

$$\sum_{i=1}^{\infty} \frac{1}{i^p} < \infty \text{ if } p > 1$$

$p < 1$ ,  $p = 1$  harmonic



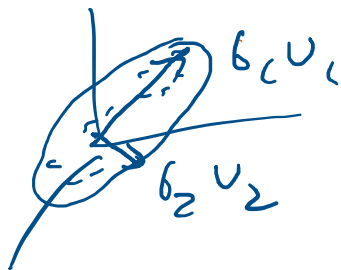
$b_i^2$  the power spectrum

$$= 1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \dots \quad p = 1$$

$B = \frac{1}{n-1} \bar{X}^T$  ; let  $C = \frac{1}{n-1} B^T B$  covariance matrix

$B = U \Sigma V^T$

$U = [u_1, u_2, \dots]$



~~$b_3 < 0$~~   
 $b_2 \geq 0$

$u_1 = \underset{\|u\|=1}{\operatorname{argmax}} u^T B^T B u = \underset{u}{\operatorname{argmax}} \frac{u^T B^T B u}{\underbrace{u^T u}_1}$

Rayleigh-Ritz gradient -  $\|B u\|_2^2$   $\rightarrow$   $B u \cdot B u = \|B u\|_2^2$

$u_2 = \underset{\|u\|=1, u \perp u_1}{\operatorname{argmax}} u^T B^T B u$

The Eigs of  $C=B'B$  give optimal projection – thus PCA and... KL

$C V = V D$  eigenvectors of  $C$  all stacked.

optimize  $x^T A x$ , maximize  $x^T A x$ , maximize  $x^T A x$ , maximize  $x^T A x$  (  $A = B^T B$  )  
 $x = \begin{pmatrix} x^1 \\ \vdots \\ x^j \\ \vdots \\ x^n \end{pmatrix}$   $r(x) = \frac{x^T A x}{x^T x}$

$$\frac{\partial r}{\partial x^j} = \frac{\partial}{\partial x^j} \left( \frac{x^T A x}{x^T x} \right) = \frac{\frac{\partial}{\partial x^j} (x^T A x) - x^T A x \frac{\partial}{\partial x^j} (x^T x)}{(x^T x)^2}$$

$$= \frac{2(Ax)_j}{x^T x} - \frac{(x^T A x) 2x_j}{(x^T x)^2} = \frac{2}{x^T x} (Ax - r(x)x)_j$$

$$\nabla r(x) = \frac{2}{x^T x} (Ax - r(x)x) = \frac{2}{x^T x} (A - r(x)I)x = 0$$

$$Ax = \underbrace{r(x)}_{\lambda} x \quad \Rightarrow \quad Ax = \lambda x$$

Conclude

The  $x$  that optimizes  $r(x) = \frac{x^T Ax}{x^T x}$  is an eigenvector and  $r(x)$  is its eigenvalue.

• S.D.T. for  $A = B^T B = \sum_{i=1}^n \lambda_i \underbrace{U_i U_i^T}_{\substack{\text{are} \\ \text{the} \\ v_i^T}}$

let  $B = U \Sigma V^T$  (s.v.d.)

$$\underline{B^T B} = V \Sigma^T \cancel{U^T U} \Sigma V^T = V \underbrace{\Sigma^T \Sigma}_{\text{}} V^T$$

$$D = \begin{pmatrix} b_1 & \dots & b_r & & \\ & & & \dots & 0 \\ & & & & \dots & 0 \\ & & & & & \dots & 0 \end{pmatrix} \quad \left( \begin{array}{c|c} b_1 & \dots & b_r & & \\ \hline & & & \dots & 0 \\ & & & & \dots & 0 \\ & & & & & \dots & 0 \end{array} \right) = \begin{pmatrix} D & & & \\ & b_1^2 & & \\ & & b_2^2 & \\ & & & \dots & b_r^2 \\ & & & & & \dots & 0 \\ & & & & & & \dots & 0 \end{pmatrix}$$

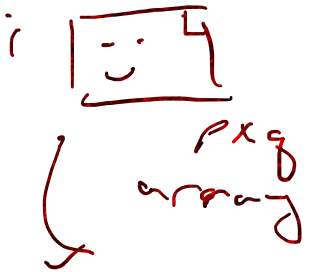
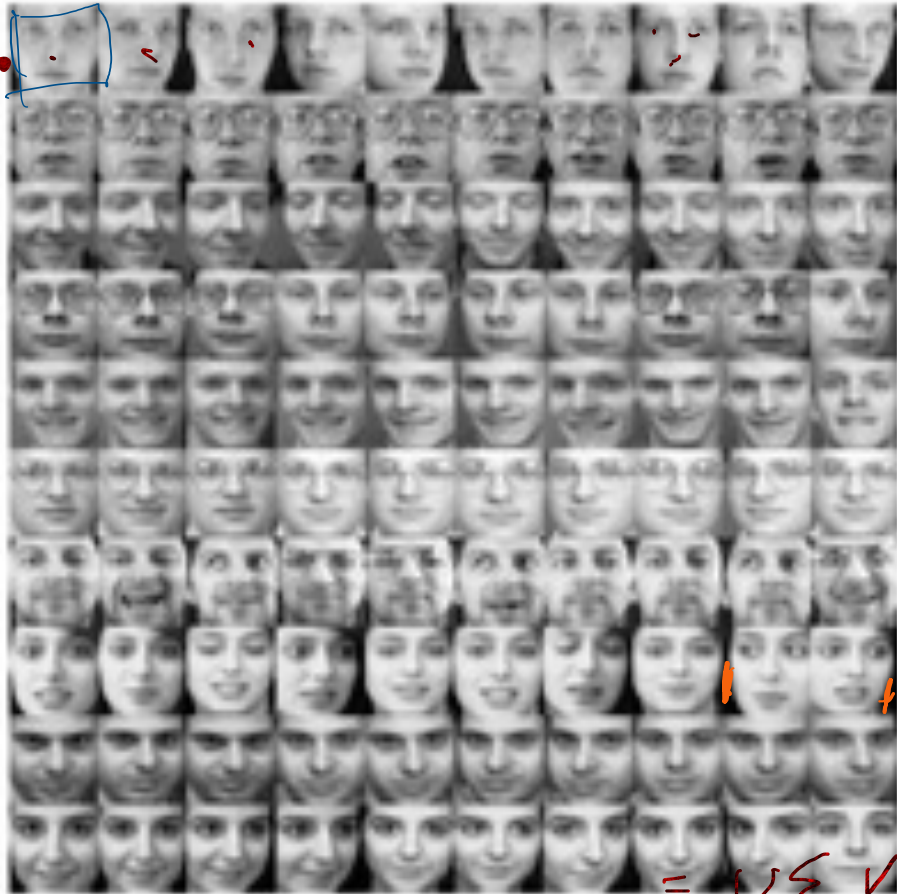
$$= \underline{V D V^T} = \sum_{i=1}^r \underline{b_i^2} \underline{v_i} \underline{v_i^T}$$



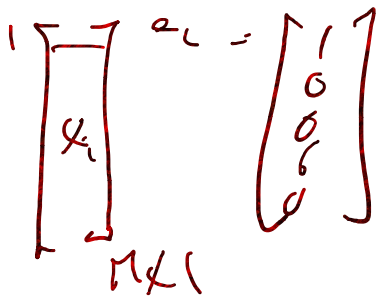
Eigenface



72 picture



reshape as vector



$$M = pq$$

$$\underline{X} = [x_1, x_2, \dots, x_n]_{m \times n}$$

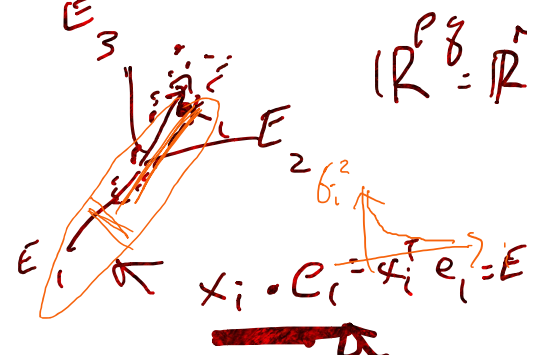
$$= U \Sigma V^T$$

$$p = 6000, q = 2000 \quad 08/31/20$$

remove unwanted variance



Register

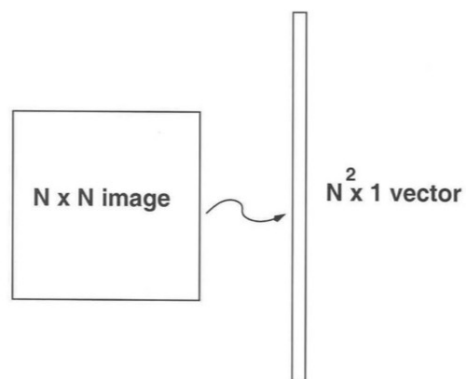


## Eigenfaces for Face Detection/Recognition

(M. Turk and A. Pentland, "Eigenfaces for Recognition", *Journal of Cognitive Neuroscience*, vol. 3, no. 1, pp. 71-86, 1991, hard copy)

### • Face Recognition

- The simplest approach is to think of it as a template matching problem:



- Problems arise when performing recognition in a high-dimensional space.

- Significant improvements can be achieved by first mapping the data into a *lower-dimensionality* space.

- How to find this lower-dimensional space?

### • Main idea behind eigenfaces

- Suppose  $\Gamma$  is an  $N^2 \times 1$  vector, corresponding to an  $N \times N$  face image  $I$ .

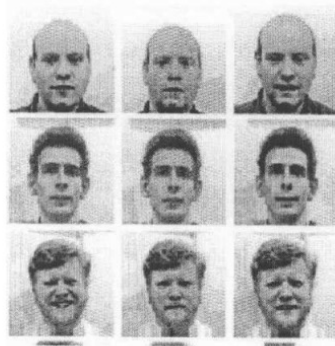
- The idea is to represent  $\Gamma$  ( $\Phi = \Gamma$  - mean face) into a low-dimensional space:

$$\hat{\Phi} - \text{mean} = w_1 u_1 + w_2 u_2 + \dots + w_K u_K \quad (K \ll N^2)$$

## Computation of the eigenfaces

Step 1: obtain face images  $I_1, I_2, \dots, I_M$  (training faces)

(**very important:** the face images must be *centered* and of the same *size*)



Step 2: represent every image  $I_i$  as a vector  $\Gamma_i$

Step 3: compute the average face vector  $\Psi$ :

$$\Psi = \frac{1}{M} \sum_{i=1}^M \Gamma_i$$

Step 4: subtract the mean face:

$$\Phi_i = \Gamma_i - \Psi$$

Step 5: compute the covariance matrix  $C$ :

$$C = \frac{1}{M} \sum_{n=1}^M \Phi_n \Phi_n^T = AA^T \quad (N^2 \times N^2 \text{ matrix})$$

$$\text{where } A = [\Phi_1 \ \Phi_2 \ \dots \ \Phi_M] \quad (N^2 \times M \text{ matrix})$$

$$C = \frac{1}{M} X^T X$$

Step 6: compute the eigenvectors  $u_i$  of  $AA^T$

The matrix  $AA^T$  is very large --> not practical !!

Step 6.1: consider the matrix  $A^T A$  ( $M \times M$  matrix)

Step 6.2: compute the eigenvectors  $v_i$  of  $A^T A$

$$A^T A v_i = \mu_i v_i$$

What is the relationship between  $u_i$  and  $v_i$ ?

$$A^T A v_i = \mu_i v_i \Rightarrow AA^T A v_i = \mu_i A v_i \Rightarrow$$

$$C A v_i = \mu_i A v_i \text{ or } C u_i = \mu_i u_i \text{ where } u_i = A v_i$$

Thus,  $AA^T$  and  $A^T A$  have the same eigenvalues and their eigenvectors are related as follows:  $u_i = A v_i$  !!

Note 1:  $AA^T$  can have up to  $N^2$  eigenvalues and eigenvectors.

Note 2:  $A^T A$  can have up to  $M$  eigenvalues and eigenvectors.

Note 3: The  $M$  eigenvalues of  $A^T A$  (along with their corresponding eigenvectors) correspond to the  $M$  largest eigenvalues of  $AA^T$  (along with their corresponding eigenvectors).

Step 6.3: compute the  $M$  best eigenvectors of  $AA^T$ :  $u_i = A v_i$

**(important:** normalize  $u_i$  such that  $\|u_i\| = 1$ )

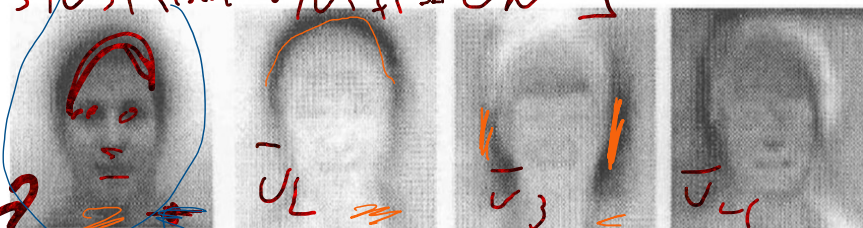
Step 7: keep only  $K$  eigenvectors (corresponding to the  $K$  largest eigenvalues)

### Representing faces onto this basis

- Each face (minus the mean)  $\Phi_i$  in the training set can be represented as a linear combination of the best  $K$  eigenvectors:

$$\hat{\Phi}_i - \text{mean} = \sum_{j=1}^K w_j u_j, \quad (w_j = u_j^T \Phi_i)$$

(we call the  $u_j$ 's eigenfaces)



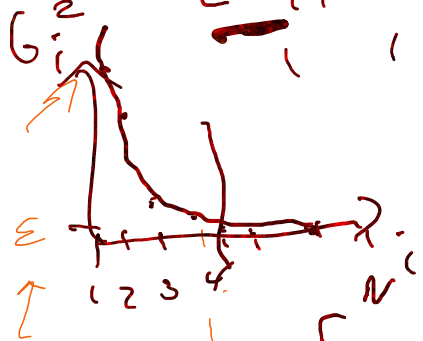
Each normalized training face  $\Phi_i$  is represented in this basis by a vector:

$$\Omega_i = \begin{bmatrix} w_1^i \\ w_2^i \\ \dots \\ w_K^i \end{bmatrix}, \quad i = 1, 2, \dots, M$$

$\cdot$  images  $\text{sc}(\mathcal{I})$   
 $\cdot \mathcal{I} = \text{reshape}(U(:, i), p, q)$

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

$$U = [u_1 | u_2 | \dots | u_s | u_{s+1} | \dots | u_r | u_{r+1} | \dots | u_n]$$



$$\|\vec{a}\|_2^2 = \sum_i a_i^2$$

Preserved 'neg.'  
 $\| \vec{a} \|_2^2 \approx \| \vec{a} \|_2^2 = \| \vec{a} \|_2^2$

$\vec{u} = a_1 \vec{u}_1 + a_2 \vec{u}_2 + a_3 \vec{u}_3 + \dots$   
 Fourier





On basis, functions, and Hilbert space.  
Fourier, Taylor, Wavelet, POD-KL

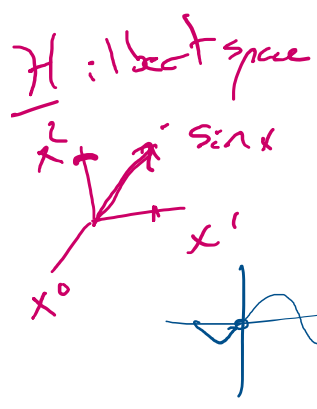
lead in 3 \$ (2) - in 5 min. Signals analysis, Harmonic.

Historically favorite basis set.  $B = \{\omega_1, \omega_2, \dots\}$   
 (vs. energy favorite basis set comes from PCA)

Taylor polynomials.  $f(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_n x^n$

Fourier modes.  $B = \{x^0, x^1, x^2, \dots\}$   
 $a_i = \frac{f^{(i)}(0)(x-0)^i}{i!}$

$f(x) = a_1 \sin x + a_2 \sin 2x + a_3 \sin 3x + \dots$   
 $B = \{\sin x, \sin 2x, \dots\}$



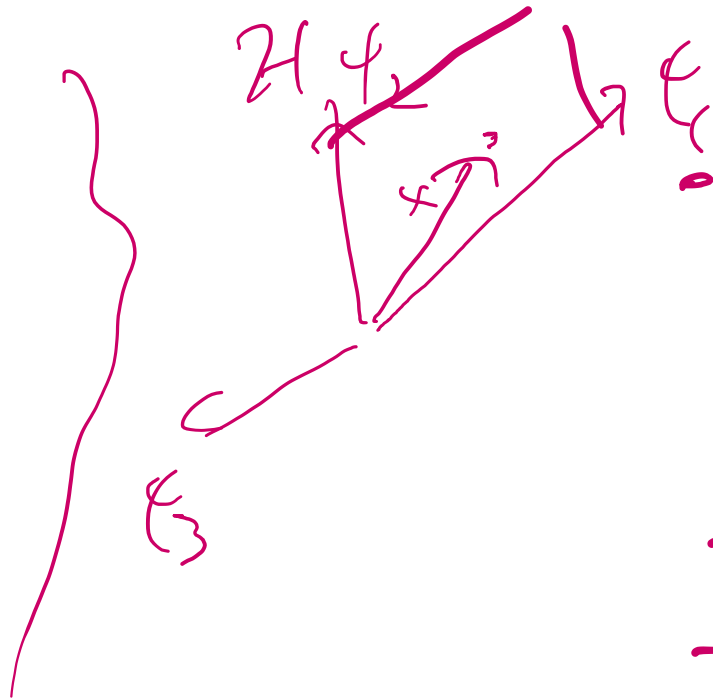
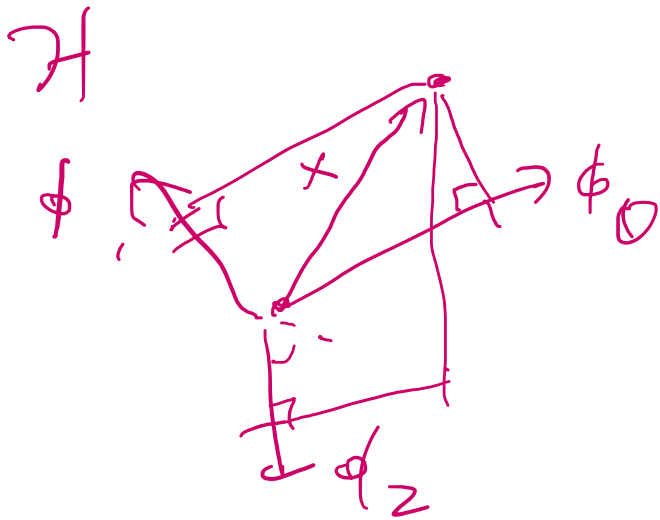
$a_1 = (f(x), \sin x) = \frac{1}{2\pi} \int_0^{2\pi} f(x) \sin x dx$

Wavelet basis. ~~Legendre~~  
 Chebyshev polys, Legendre ...

$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots$   
 $\text{complex } \sin x = 1$



Changing bases is a sort of coord-rot.



$$\begin{aligned} \phi_0(y) &= y^0 \\ \phi_1(y) &= y^1 \\ \phi_2(y) &= y^2 \\ &\vdots \end{aligned}$$

$$\begin{aligned} \rightarrow \phi_c(y) &= \sin(y) \\ \rightarrow \phi_2(y) &= \sin(2y) \\ \rightarrow \phi_3(y) &= \sin(3y) \end{aligned}$$

On basis, functions, and Hilbert space. Fourier, Taylor, Wavelet, POD-KL

ks. 2.1

Hilbert space - a complete inner product space.

• An inner product space is a "vector space"  $E$  together with a function called "inner product"  $\langle \cdot, \cdot \rangle : E \times E \rightarrow \mathbb{C}$  with properties.

GFT: you set geometry in  $E$  as on a circle  $\langle u, v \rangle = \frac{\langle u, v \rangle}{\|u\| \|v\|}$

• and projection vector space.

set of objects "like vectors" that have a + and scalar multiplication - including commutative, associative, add. ident, inverse, distributive w/ scalar, identity for scalar.

- conjugate (sym.)  $\langle u, v \rangle = \overline{\langle v, u \rangle} \quad \forall u, v \in E$
- linear  $\langle au + bv, v \rangle = a \langle u, v \rangle + b \langle v, v \rangle \quad \forall a, b \in \mathbb{C}, u, v \in E$
- pos. defn:  $\langle u, u \rangle \geq 0 \quad \forall u \in E, u \neq 0$

Ex: arrays of real numbers that are  $2 \times 1$ .  $\begin{bmatrix} 1 \\ 2 \end{bmatrix} + c \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 3c+1 \\ 4c+2 \end{bmatrix}$

(norm:  $\|x\|^2 = \langle x, x \rangle$ )

Ex functions in  $C([0,1])$  e.g.  $3x^2 + 4x^3 + 7 \sin x = f(x) \in C([0,1])$   
 $\phi_1(x) \quad \phi_2(x) \quad \phi_3(x)$

Ex. Norm on inner product space  $L^2([0,1]) = \{f \mid \int_0^1 |f(x)|^2 dx < \infty\}$   
 $(L^2([0,1]) \subset C([0,1]))$

let  $\langle f, g \rangle_{L^2([0,1])} = \int_0^1 f(x) \overline{g(x)} dx$   $\exists$  btw  $\| \cdot \|_{L^2([0,1])} = \left( \int_0^1 |f(x)|^2 dx \right)^{1/2}$

Basis of unit vectors.  $B = \{ \phi_i \mid \phi_i \in B \subset E \text{ and } \langle \phi_i, \phi_j \rangle = \delta_{ij} \text{ and } \|\phi_i\| = 1 \}$   
 • orthogonal

$u = \sum a_i \phi_i$   $\leftarrow$  projection!

Ex:  $L([0,1])$  ;  $B = \{ \cos \frac{2\pi k x}{L}, \sin \frac{2\pi k x}{L}, 1 \mid k \in \mathbb{N} \}$   
 $\cdot f \in L^1$  ;  $f(x) = \frac{a_0}{2} + \sum_{k=1}^{\infty} a_k \cos kx + b_k \sin kx$  ;  $a_k = \frac{2}{L} \int_0^L f(x) \cos kx dx$  ;  $b_k = \frac{2}{L} \int_0^L f(x) \sin kx dx$



$\phi_k(x) = \frac{2 \cos kx}{L}$  or  $\frac{2 \sin kx}{L}$

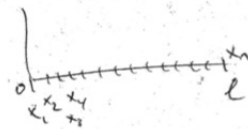
• similar if  $\phi_k(x) = c e^{ikx}$

finite if "trig. poly."

or maybe  $B = \sum x^k \cdot \xi$  Taylor poly.

or maybe  $B = \sum \Delta$  wavelets.  
 Hat "2.13" function plot

Infinite dimensional inner product space of functions.



$x_1 < x_2 < \dots < x_n$  grid on  $[0, L]$   
 $G = \sum_{i=1}^n \xi_i \delta_{x_i}$

and "function values."  $f: G \rightarrow \mathbb{R}$ .  
 let  $f_i = f(x_i) \Rightarrow \vec{f}$   
 • connect dots if you like.

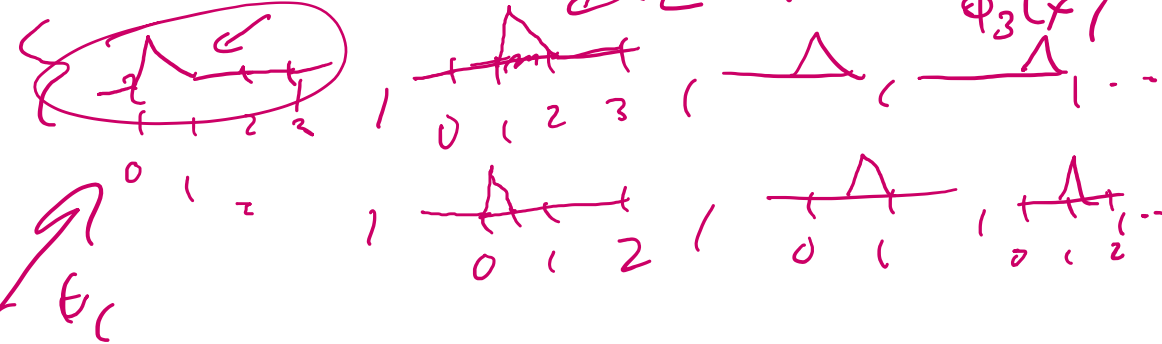
Ex let  $\langle \vec{f}, \vec{g} \rangle = \sum_{k=1}^n f_k \overline{g_k}$

• this inner product space is "isomorphic" to vectors in  $\mathbb{C}^n$ .

• separable. there exists a countable basis  $B$ .  
 Ex:  $L^2([0,1])$  vs.  $L^2(\mathbb{R})$

$\circ B = \{ \text{---}, \text{~}, \text{~}, \text{~}, \text{~}, \text{~} \}$  [old] Fourier  
 $\Rightarrow$  Fourier  $\{ 1, \sin x, \cos x, \sin 2x, \cos 2x, \dots \}$

$\circ B' = \{ \text{---}, /, \text{~}, \text{~}, \text{~}, \dots \} = \{ 1, x, x^2, x^3, \dots \}$   
 Taylor ... Legendre

$\circ B'' =$  

$f(x) \in L^2(\mathbb{R})$

$f(x) = \sum a_i \phi_i(x);$

$a_i = \langle f, \phi_i(x) \rangle$   
 $\Rightarrow \|\phi_i(x)\|$

$\circ$   Hence.....

On Compressed Sensing and on to Sparsity

$$f(x) = \sin x + 3 \sin 3x + 4 \sin 7x \in L^2([0, \pi])$$

$$= \sum_{n=1}^{\infty} a_n \sin nx$$

$$|f(x)| = \begin{pmatrix} a_1 \\ a_3 \\ a_7 \\ \dots \end{pmatrix}$$

$$= \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \dots \end{pmatrix}$$

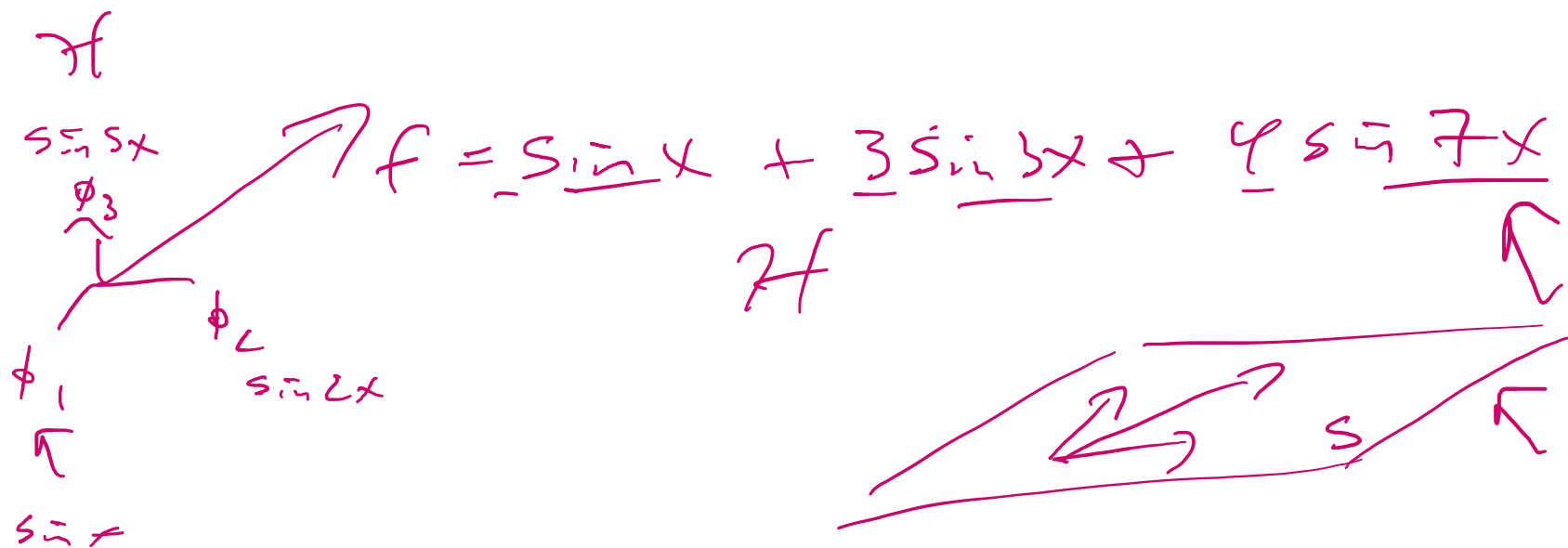
$$|f| = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \dots \end{pmatrix}$$

$$|f(x)| = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots$$

$$+ 3 \left( 3x - \frac{(3x)^3}{3!} + \frac{(3x)^5}{5!} - \dots \right)$$

$$+ 4 \left( 7x - \frac{(7x)^3}{3!} + \frac{(7x)^5}{5!} - \dots \right)$$

$$= 38x - \left( \frac{1}{5!} + \frac{3^3}{3!} + \frac{4 \cdot 7^3}{3!} \right) x^3 + \dots$$



• a vector  $v \in E$  is  $k$ -sparse, if  $[v]$  has exactly  $k$ -nonzero values, and  $k \leq \dim(E)$

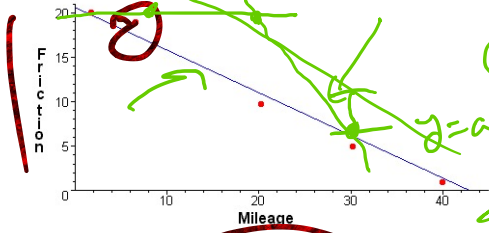
On Moore Penrose Pseudo Inverse, Matrix Least Squares, Geometric Least Squares.

**Example**

An engineer is tracking the friction index over mileage of a braking system of a vehicle. She expects that the mileage-friction relationship is approximately linear. She collects five data points that are show in the table below.

Mileage	2000	6000	20,000	30,000	40,000
Friction Index	20	18	10	6	2

The graph below shows these points



We are interested in the line that best fits the data. More specifically, if  $\mathbf{b}$  is the vector of friction index data values and  $\mathbf{y}$  is the vector consisting of  $y$  values when we plug in the mileage data for  $x$  and find  $y$  by the equation of the line, then we want the line that minimizes the distance between  $\mathbf{b}$  and  $\mathbf{y}$ . If the equation of the line is

$$ax + b = y$$

then we get the five equations

- $2a + b = 20$
- $6a + b = 18$
- $20a + b = 10$
- $30a + b = 6$
- $40a + b = 2$

The corresponding matrix equation is

$$A\mathbf{x} = \mathbf{y}$$

or

$$\begin{pmatrix} 2 & 1 \\ 6 & 1 \\ 20 & 1 \\ 30 & 1 \\ 40 & 1 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 20 \\ 18 \\ 10 \\ 6 \\ 2 \end{pmatrix}$$

Although this does not have an exact solution, it does have a closest solution. We have

$$\begin{pmatrix} a \\ b \end{pmatrix} = (A^T A)^{-1} A^T \mathbf{y} = \begin{pmatrix} -0.48 \\ 20.6 \end{pmatrix}$$

We can conclude that the equation of the regression line is

$$y = -0.48x + 20.6$$

$$\text{Cost} = \sum_{i=1}^n \left( (y_i - (ax + b))^2 \right)$$



$$\begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

$$\begin{aligned} 2 \cdot a + 1 \cdot b &= 20 \\ 6 \cdot a + 1 \cdot b &= 18 \\ &\vdots \end{aligned}$$

$$\begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 3 \\ 4 \end{pmatrix}$$

$$\mathbf{x} = \mathbf{A}^+ \mathbf{y} \quad \text{vs} \quad \mathbf{A}^T$$

## Least Squares

### Definition and Derivations

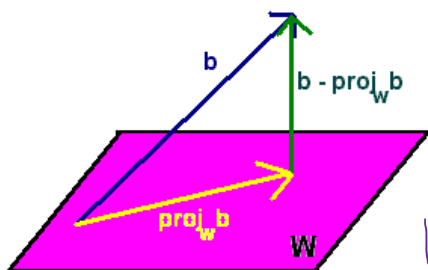
We have already spent much time finding solutions to

$$Ax = b$$

If there isn't a solution, we attempt to seek the  $x$  that gets closest to being a solution. The closest such vector will be the  $x$  such that

$$Ax = \text{proj}_W b$$

where  $W$  is the column space of  $A$ .



$$W = \text{col}(A)$$

Notice that  $b - \text{proj}_W b$  is in the orthogonal complement of  $W$  hence in the null space of  $A^T$ . Hence if  $x$  is a this closest vector, then

$$A^T(b - Ax) = 0 \quad A^T Ax = A^T b$$

Now we need to show that  $A^T A$  nonsingular so that we can solve for  $x$ .

### Lemma

If  $A$  is an  $m \times n$  matrix of rank  $n$ , then  $A^T A$  is nonsingular.

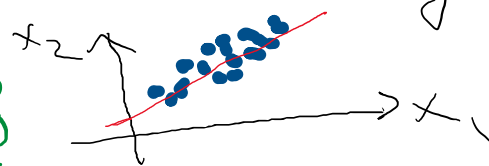
$$\begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & & \\ a_{31} & & & \\ \vdots & & & \\ a_{m1} & \dots & \dots & a_{mn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix}$$

$m > n$   
 $n$ -eqns  
 $n$ -unknowns  
 $x$

$$\Leftrightarrow \begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1 \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n = b_m \end{cases}$$

$m > n$  tall skinny,  $n=2$

How many?



How many?



$$x_1 \begin{pmatrix} a_{11} \\ a_{21} \\ \vdots \\ a_{m1} \end{pmatrix} + x_2 \begin{pmatrix} a_{12} \\ \vdots \\ a_{m2} \end{pmatrix} + \dots + x_n \begin{pmatrix} a_{1n} \\ \vdots \\ a_{mn} \end{pmatrix} = \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix}$$



**Theorem**

Let  $A$  be an  $m \times n$  matrix of rank  $n$ , then the system

$$Ax = b$$

has the unique *least squares* solution

$$x = (A^T A)^{-1} A^T b$$

**Example**

Find the least squares solution to

$$Ax = b$$

with

$$A = \begin{pmatrix} 1 & 3 \\ 2 & 4 \\ 1 & 6 \end{pmatrix} \quad b = \begin{pmatrix} 4 \\ 1 \\ 3 \end{pmatrix}$$

**Solution**

We can quickly check that  $A$  has rank 2 (the first two rows are not multiples of each other). Hence we can compute

$$x = (A^T A)^{-1} A^T b = \begin{pmatrix} -0.377 \\ .662 \end{pmatrix}$$

Notice that

$$Ax = \begin{pmatrix} 1.61 \\ 1.90 \\ 3.60 \end{pmatrix}$$

not exactly  $b$ , but as close as we are going to get.

$(A^T A)^{-1}$  - never form that

$Ax = b$   $m \times n$  • Square & solving.  $A^{-1} b$  d.n.e. uniquely. • Not,  $m > n$

$A^T A x = A^T b$  normal equations  
 Covariance matrix (if demeaned)

$$\begin{pmatrix} 1 & 1 & 2 & 1 \\ 3 & 4 & 6 \end{pmatrix}^T \begin{pmatrix} 1 & 3 \\ 2 & 4 \\ 1 & 6 \end{pmatrix}^{-1} \begin{pmatrix} 1 & 2 & 1 \\ 3 & 4 & 6 \end{pmatrix} \begin{pmatrix} 4 \\ 1 \\ 3 \end{pmatrix} = x$$

### Best Fitting Curves

Often, a line is not the best model for the data. Fortunately the same technique works if we want to use other nonlinear curves to fit the data. Here we will explain how to find the least squares cubic. The process for other polynomials is similar.

#### Example

A bioengineer is studying the growth of a genetically engineered bacteria culture and suspects that it is approximately follows a cubic model. He collects six data points listed below

Time in Days	1	2	3	4	5	6
Grams	2.1	3.5	4.2	3.1	4.4	6.8

He assumes the equation has the form

$$ax^3 + bx^2 + cx + d = y$$

This gives six equations with four unknowns

$$\begin{aligned} a + b + c + d &= 2.1 \\ 8a + 4b + 2c + d &= 3.5 \\ 27a + 9b + 3c + d &= 4.2 \\ 64a + 16b + 4c + d &= 3.1 \\ 125a + 25b + 5c + d &= 4.4 \\ 216a + 36b + 6c + d &= 6.8 \end{aligned}$$

The corresponding matrix equation is

$$\begin{pmatrix} 1 & 1 & 1 & 1 \\ 8 & 4 & 2 & 1 \\ 27 & 9 & 3 & 1 \\ 64 & 16 & 4 & 1 \\ 125 & 25 & 5 & 1 \\ 216 & 36 & 6 & 1 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \\ d \end{pmatrix} = \begin{pmatrix} 2.1 \\ 3.5 \\ 4.2 \\ 3.1 \\ 4.4 \\ 6.8 \end{pmatrix}$$

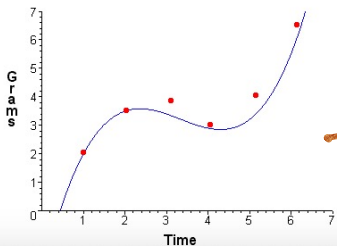
We can use the least squares equation to find the best solution

$$\begin{pmatrix} a \\ b \\ c \\ d \end{pmatrix} = (A^T A)^{-1} A^T \mathbf{b} = \begin{pmatrix} 0.2 \\ -2.0 \\ 6.1 \\ -2.3 \end{pmatrix}$$

So that the best fitting cubic is

$$y = 0.2x^3 - 2.0x^2 + 6.1x - 2.3$$

The graph is shown below



Matrix Formulation of

LS slide

for general Models



$$\begin{aligned} 2.1 &= a \cdot 1^3 + b \cdot 1^2 + c \cdot 1 + d \\ 3.5 &= a \cdot 2^3 + b \cdot 2^2 + c \cdot 2 + d \end{aligned}$$

$$\begin{aligned} Ax &= b \\ X^T P &= I \end{aligned}$$

$$\begin{pmatrix} x_1^3 & x_1^2 & x_1 & 1 \\ x_2^3 & x_2^2 & x_2 & 1 \\ \vdots & \vdots & \vdots & \vdots \\ x_n^3 & x_n^2 & x_n & 1 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \\ d \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

$Ax = b$   $\Rightarrow x_1 \vec{a}_1 + x_2 \vec{a}_2 + \dots + x_n \vec{a}_n = \vec{b}$

$A\tilde{x} = \text{proj}_W b$       $\tilde{x} = \underset{\text{min}}{\text{argmin}} \|Ax - b\|_2^2$

$\Rightarrow$  LS solution

$W = \text{Col}(A)$

$\vec{e} = A\tilde{x} - b$

Recall  $u \perp v$  iff  $u \cdot v = u^T v = 0$

vs.  $u \cdot v = \|u\| \|v\| \cos \theta$

$\Rightarrow$   $A^T(A\tilde{x} - b) = 0$

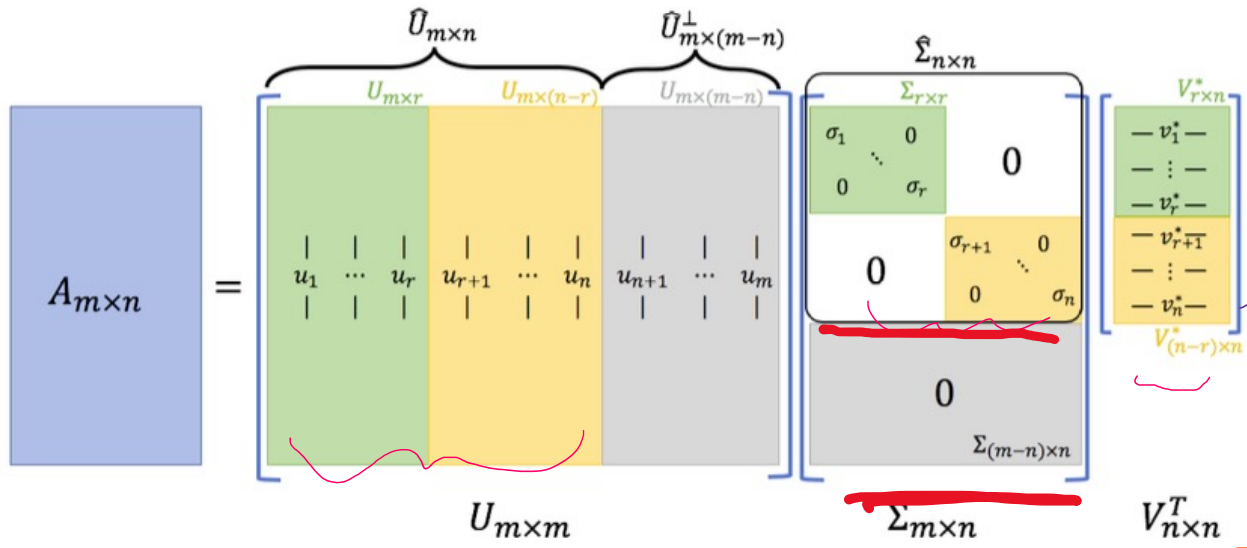
$\begin{pmatrix} -a_1' \\ -a_2' \\ \vdots \\ -a_n' \end{pmatrix} (Ax - b) = 0 \ll \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$

$A^T(Ax - b) = 0$

$\Rightarrow \vec{a}_i \perp (Ax - b) \forall i$

$\Rightarrow$   $(A\tilde{x} - b) \perp$  every vector in  $\text{Col}(A)$   $\Leftrightarrow$  solve "normal eqns"

$A^T A \tilde{x}$



Handwritten notes and equations:

$$(\hat{V} \hat{\Sigma}^T \hat{V}^T) (U \Sigma V^T) x$$

$$\hat{V} \hat{\Sigma}^T \hat{V}^T x = \hat{V} \hat{\Sigma}^T U^T b$$

$$(\hat{\Sigma}^T \hat{\Sigma}) V^T x = \hat{\Sigma}^T U^T b$$

$$V^T x = V (\hat{\Sigma}^T \hat{\Sigma})^{-1} \hat{\Sigma}^T U^T b$$

$A^T A x = A^T b$

$A_{m \times n} = \hat{U}_{m \times n} \hat{\Sigma}_{n \times n} \hat{V}_{n \times n}^T$

Expanded SVD decomposition:

$$= \begin{bmatrix} | & | & & | \\ u_1 & u_2 & \dots & u_n \\ | & | & & | \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n \end{bmatrix} \begin{bmatrix} - & v_1^T & - \\ - & v_2^T & - \\ - & \vdots & - \\ - & v_n^T & - \end{bmatrix}$$

Handwritten boxed equation:

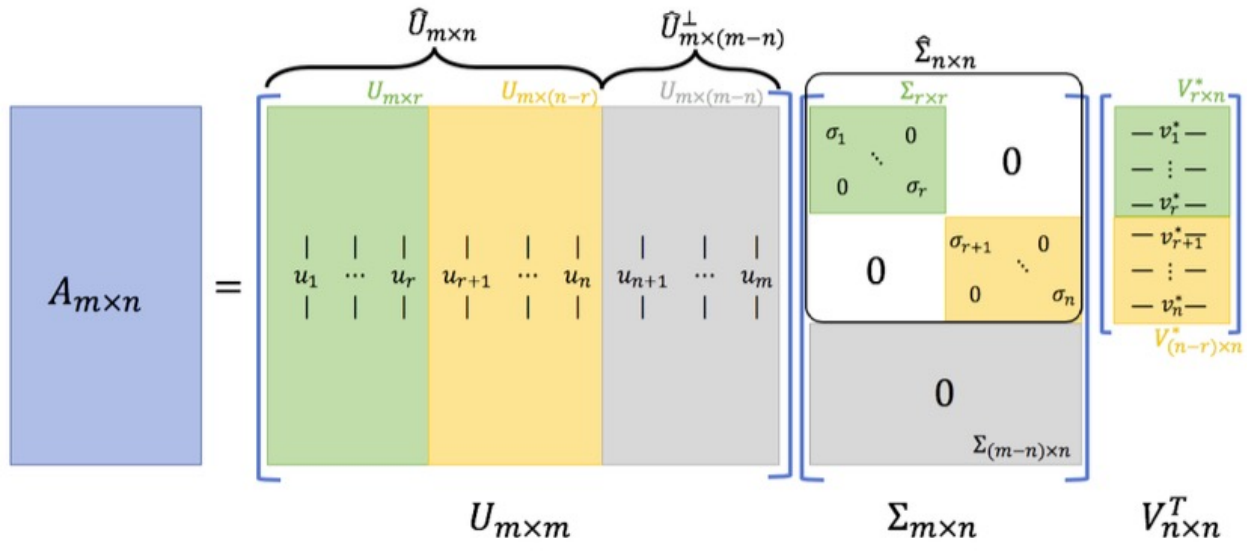
$$x = V \left[ \hat{\Sigma}^T \hat{\Sigma} \right]^{-1} \hat{\Sigma}^T U^T b$$

The term  $\left[ \hat{\Sigma}^T \hat{\Sigma} \right]^{-1} \hat{\Sigma}^T U^T b$  is labeled  $A^+$ .

Handwritten notes:

$\hat{\Sigma} \hat{\Sigma}^T =$  (matrix with  $\sigma_i^2$  on diagonal)

when exist  $\rightarrow 0 = 0$



$\hat{\Sigma}^+ = (\hat{\Sigma}^T \hat{\Sigma})^{-1} \hat{\Sigma}$   
 if  $\hat{\Sigma}$  is invertible  
 if  $r = n$   
 if  $\sigma_i > 0, 1 \leq i < n$

But if  $r < n$   
 $\sigma_r > 0, \sigma_{r+1} = 0$

$$A_{m \times n} = \hat{U}_{m \times n} \hat{\Sigma}_{n \times n} \hat{V}_{n \times n}^T$$

$$= \begin{bmatrix} | & | & & | \\ u_1 & u_2 & \dots & u_n \\ | & | & & | \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n \end{bmatrix} \begin{bmatrix} - & v_1^T & - \\ - & v_2^T & - \\ - & \vdots & - \\ - & v_n^T & - \end{bmatrix}$$

$(\hat{\Sigma}^T \hat{\Sigma}) = \begin{pmatrix} \sigma_1^2 & & & 0 \\ & \sigma_2^2 & & \\ & & \dots & \\ & & & \sigma_r^2 \\ & & & & 0 \\ & & & & & \dots \\ & & & & & & \sigma_{n-r}^2 \\ & & & & & & & 0 \end{pmatrix}$   
 Annotations:  $\frac{1}{\sigma} \dots = 0$  (circled),  $\frac{1}{\sigma} \dots = 0$  (circled),  $\frac{1}{\sigma} \dots = 0$  (circled)

LS soln = solve normal equations

$$A^T A \tilde{x} = A^T b$$

When inverse exists

$$\tilde{x} = (A^T A)^{-1} A^T b$$

$$\equiv A^+ b$$

$A^+$   
=

Moore -  
Penrose  
Pseudo-Inverse

$$Ax = b$$

- In terms of SVD?
- and what if inverse doesn't exist.

$$Ax = b$$

$$U^T \cancel{U} \Sigma V^T x = U^T b$$

$$(\Sigma^T \cancel{U}^T V^T) x = \Sigma^T V^T b$$

$$V^T x = \underbrace{(\Sigma^T \Sigma)^{-1}}_{\Sigma^{-1}} \Sigma^T V^T b$$

$$= \Sigma^{-1} U^T b$$

$$x = \underbrace{V \Sigma^{-1} U^T}_{A^+} b := \underbrace{V \Sigma^+ U^T}_{A^+} b$$

$x$  exists

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

$$\Sigma^+ := (\Sigma^T \Sigma)^{-1} \Sigma$$

$$\equiv \begin{pmatrix} 1/\sigma_1 & \dots & 1/\sigma_r & & \\ & & & & 0 \\ & & & & \vdots \\ & & & & 0 \end{pmatrix}$$

$$\begin{cases} \sigma_i \neq 0 & \Rightarrow 1/\sigma_i \\ \sigma_i = 0 & \Rightarrow 1/0 = 0 \end{cases}$$