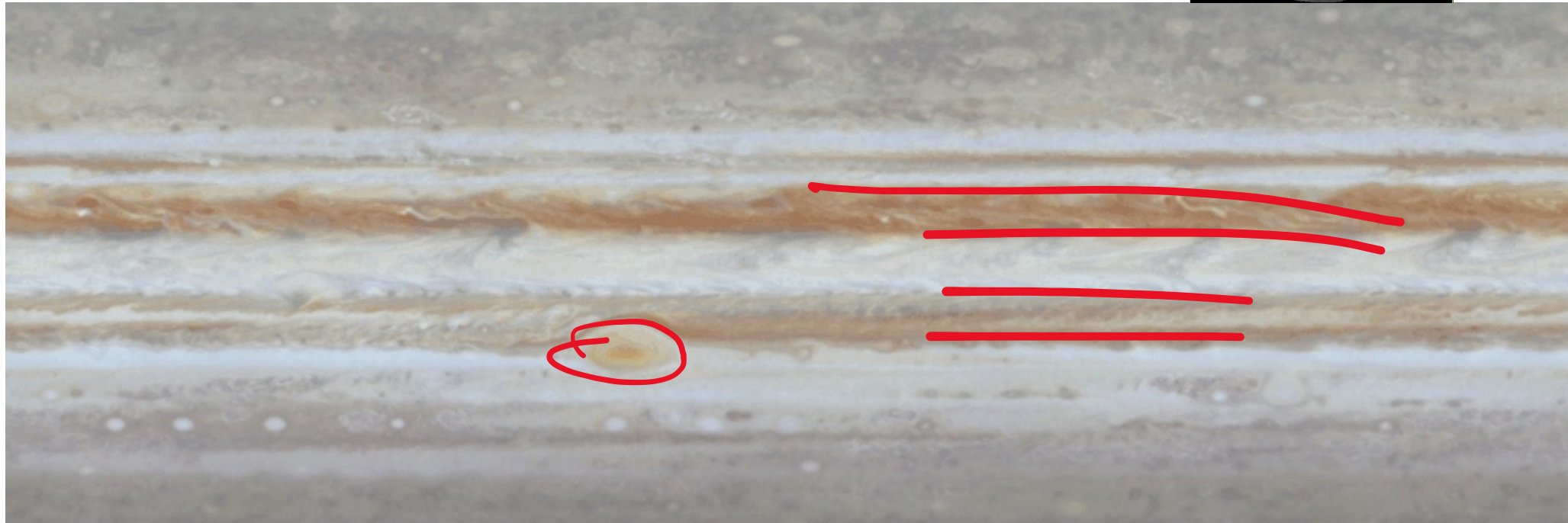# EE520 Data Driven Analysis of Complex Systems
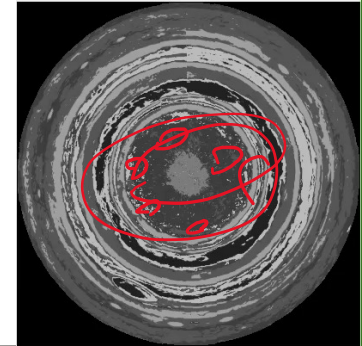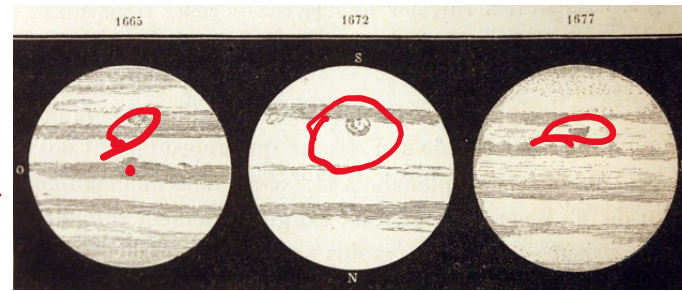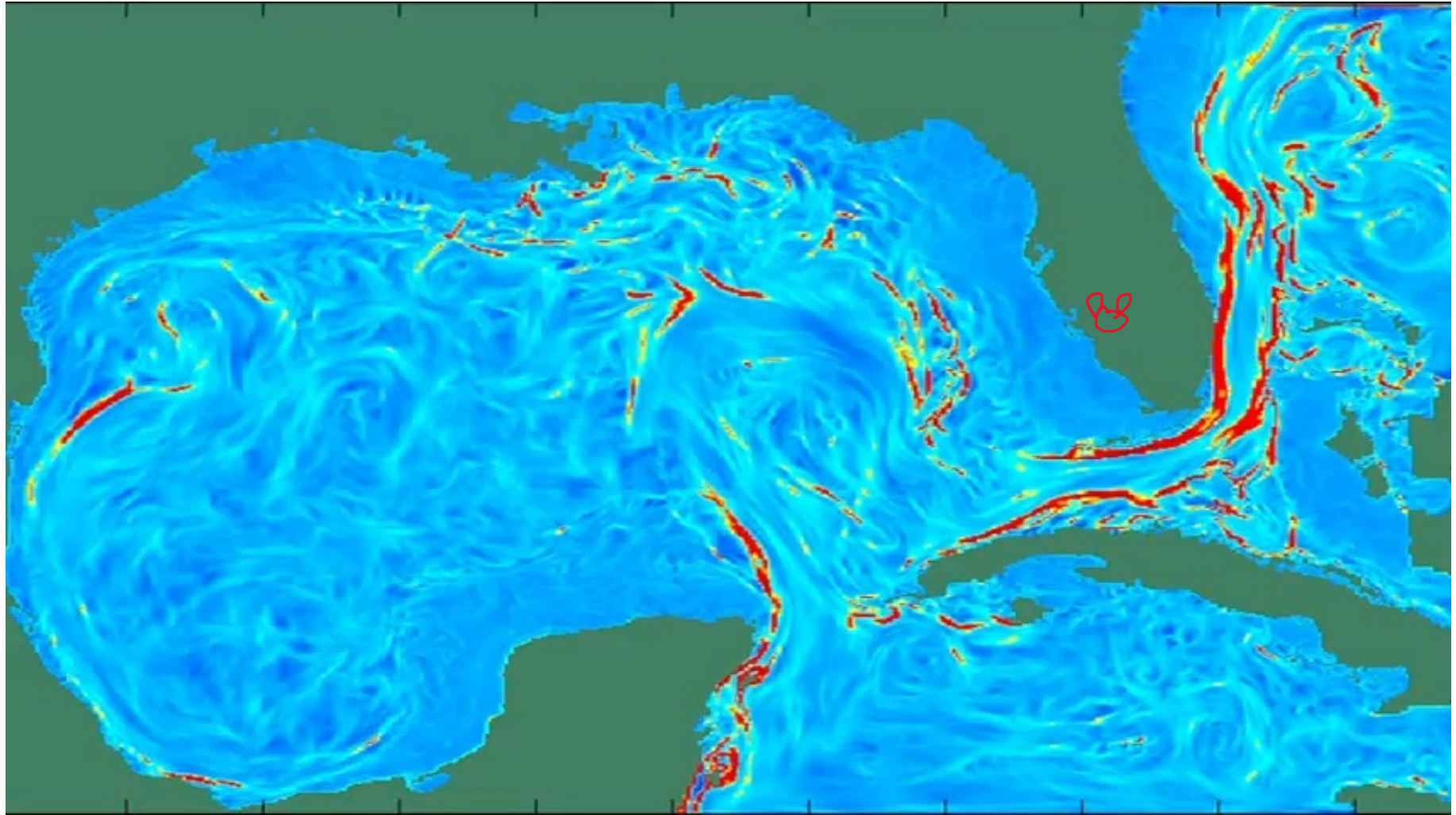
Erik Bollt
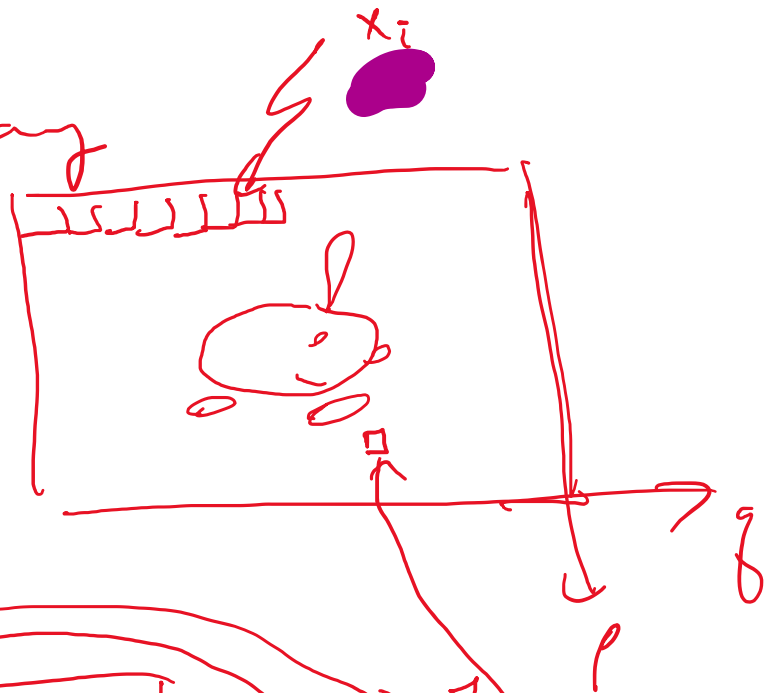
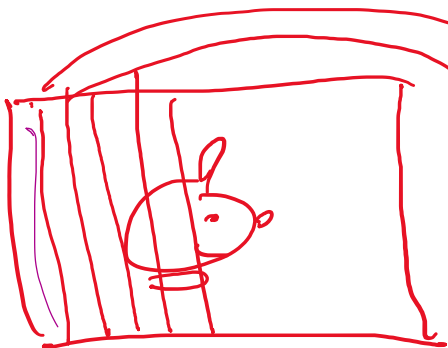Data as
an array

$x_i$

$X =$

$\underline{X}_{p \times q} \in \mathbb{R}^{p \times q}$ matrix

$\left[\underline{X}\right]_{\ell, m} =$ one pixel

$q$

$p$

slice & stack

$X_2$

$X_1$

$X_3$

$+$

Data as an array

# On Matrix Multiplication

$$L(\mathbf{z}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$$
$$\mathbf{z} \mapsto \mathbf{z}' = A\mathbf{z},$$

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

in terms of the usual matrix times vector multiplication,

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

Geometric

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

$\mathbb{R}^2$

$\mathbf{z}$

$n = 2$   $\mathbf{z} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$

$\mathbb{R}^3$   $\mathbf{z}'$

$L$
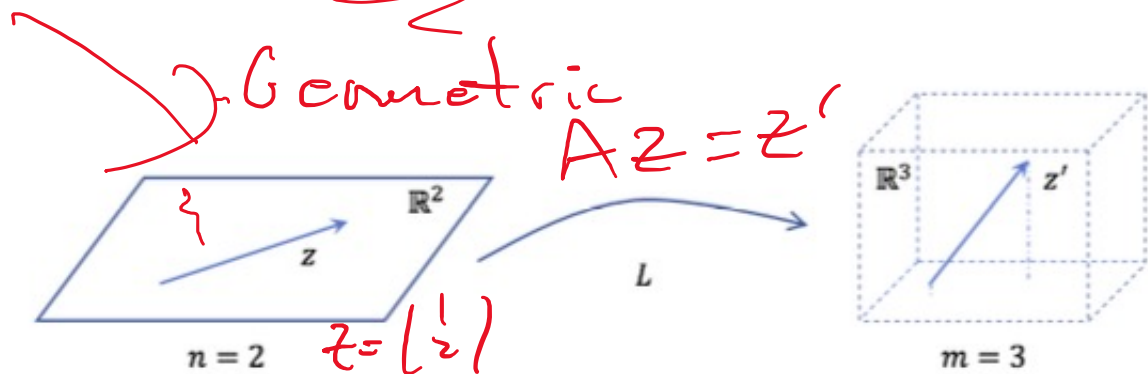
$m = 3$

- a vector has length & direction.

$$A_{2\times 2} = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$$

But as linear algebra; matrice × vectors

$$A\begin{pmatrix} 3 \\ 4 \end{pmatrix} = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}\begin{pmatrix} 3 \\ 4 \end{pmatrix} = \begin{pmatrix} 1\cdot3 + 2\cdot4 \\ 3\cdot3 + 4\cdot4 \end{pmatrix} = \begin{pmatrix} 11 \\ 25 \end{pmatrix}$$

$$z' = Az$$

$$\overline{z' Az}$$

new direction, new length.

Eig for square.

① $A\underset{2\times2}{V} = \lambda V$

Characterize matrices by knowing just these
○ eig. special directions

$$A U = v_2$$
$$v_1$$
$$a_2 \quad A \quad \lambda_2 V_2$$
$$a_1 \quad \lambda_1 V_1$$

○ $Au = A(a_1 v_1 + a_2 v_2) = a_1 A v_1 + a_2 A v_2$   $\det(A - \lambda I) = 0$

$$= a_1 \lambda_1 v_1 + a_2 \lambda_2 v_2 \qquad (A - \lambda I) v = 0$$

# Matrix time circle =

all vectors of length 1.

? Matrix × circle ?! But matrix times vector. ↗ =



$A$

$Ax$

$Ax$

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

**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:

$$A_{m \times n} = U_{m \times m} \Sigma_{m \times n} V^*_{n \times n}, \tag{2.5}$$

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 nonegative values: $\sigma_i \geq 0, i = 1, \cdots, n,$

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

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

**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 nonegative values: $\sigma_i \geq 0, i = 1, \cdots, n,$

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

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

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

$$AV = U\Sigma V^*V = U\Sigma, \tag{2.8}$$

$$\Sigma := diag(\sigma_1, \sigma_2, \cdots, \sigma_p), p = min(m, n),$$
$$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_p \geq 0.$$

Fri 08/21/20

$$A x = b$$

$$\begin{pmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$$

if

$A$ is square

$$A V = U \Sigma$$

$$[A][v_1 v_2 \cdots v_n] = [u_1 u_2 \cdots u_n] \, diag(\sigma_1, \sigma_2, \cdots, \sigma_n).$$

■ **Example 2.1** Let $A = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \end{pmatrix}_{2\times3}$. By SVD of the matrix $A$ we have:

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

$$= \begin{pmatrix} \frac{1}{\sqrt{5}} & \frac{2}{\sqrt{5}} \\ \frac{-2}{\sqrt{5}} & \frac{1}{\sqrt{5}} \end{pmatrix} \begin{pmatrix} \sqrt{70} & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{14}} & \sqrt{\frac{2}{7}} & \frac{3}{\sqrt{14}} \\ \frac{-3}{\sqrt{10}} & 0 & \frac{1}{\sqrt{10}} \\ \frac{-1}{\sqrt{35}} & \sqrt{\frac{5}{7}} & \frac{-3}{\sqrt{35}} \end{pmatrix}. \quad (2.28)$$

We see that the second singular value, $\sigma_2 = 2$, meaning that number of non-zero singular values $r < \min\{m, n\}$. Such matrix is called rank deficient matrix. If we take the economy version (with $r = 1$) of the SVD we will have:

$$u_1 \sigma_1 v_1^T = \begin{pmatrix} \frac{1}{\sqrt{5}} \\ \frac{-2}{\sqrt{5}} \end{pmatrix} (\sqrt{70}) \begin{pmatrix} \frac{1}{\sqrt{14}} & \sqrt{\frac{2}{7}} & \frac{3}{\sqrt{14}} \end{pmatrix}$$

$$\approx \begin{pmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \end{pmatrix} \quad (2.29)$$

$$[A][v_1 v_2 \cdots v_n] = [u_1 u_2 \cdots u_n] \, diag(\sigma_1, \sigma_2, \cdots, \sigma_n).$$

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

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

$$V = [v_1 \; v_2 \cdots v_n]$$

*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,
>
> $$A_{m \times n} = \hat{U}_{m \times n} \hat{\Sigma}_{n \times n} V_{n \times n}^*, \qquad (2.21)$$
>
> and $U = [\hat{U}_{m \times n} | \hat{U}_{(n-m) \times n}]$, written in terms of an orthogonal "buffer" matrix

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r \geq \sigma_{c+1} = 0$$

$$= 0 \dots$$

**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

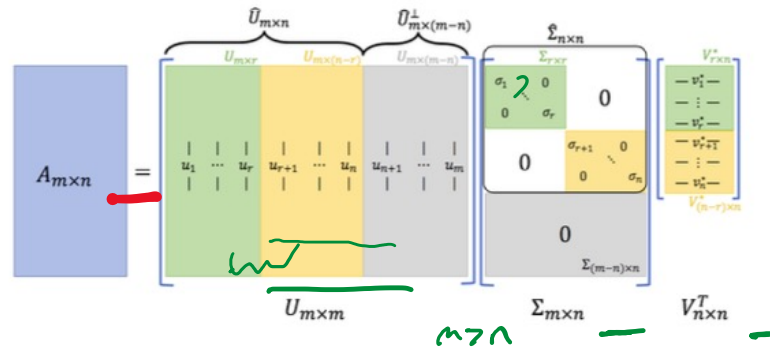$$\sigma_r > \sigma_{r+1} = 0, \text{ for some } r < n. \qquad (2.22)$$

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

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

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

Rank ill conditioned

Figure 2.3: $m > n$ tall skinny

Full, Economy, Truncated SVD

Recall that,

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

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

(2.24)

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

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

(2.25)

so,

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

(2.26)

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

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

(2.27)

## Geometry:

$$A = U\Sigma V^T$$
$$A\overline{V} = U\Sigma$$
$$Av_i = \sigma_i U_i$$

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 $A$.

$AZ$

$V^*Z$

$U\Sigma V^*Z$

$\Sigma V^*Z$

$\sigma_1 u_1$

$\sigma_2 u_2$

$v_1$

$v_2$

Unit Sphere $P$

$AP = U\Sigma V^T P$

$V^T P$

$\Sigma V^T P$

```
disp(A)

    0.8929    0.0417    0.1000
    0.3320    0.1077    0.0000
    0.8212    0.5951    0.0000

[U,S,V] = svd(A)

U = 3x3

   -0.6330    0.7469    0.2035
   -0.2590    0.0435   -0.9649
   -0.7296   -0.6635    0.1659

S = 3x3

    1.3438         0         0
         0    0.3912         0
         0         0    0.0208

V = 3x3

   -0.9304    0.3488   -0.1126
   -0.3635   -0.9176    0.1612
   -0.0471    0.1989    0.9805
```

$\sigma_1 \sim 1.5 ?$
$\text{what if } \sigma_2 \sim 0.5 ?$
$\sigma_1 = 2$
$\sigma_2 = 0.0003$
$\sigma u_i \to v_i$

And ROM

$$A$$
$$Aw = U\Sigma V^T w$$
$$v_i^T w$$
$$= r_i \cdot w$$

$1 \ 2 \ 3 \ 4$

$$X = \begin{pmatrix} h \\ \psi \\ \vdots \\ f \end{pmatrix}$$

$h$ redundant
if $\sigma r + 1 \sim 0$
$w$

$$\sigma_i$$

Index $(i)$

$r$     $n$

$\varepsilon$

$$\sigma_r > \varepsilon > \sigma_{r+1} \sim 0$$



$Z$

$v_1$

$v_2$

$AZ$

$\sigma_1 u_1$

$\sigma_2 u_2$

$V^*Z$

$U\Sigma V^*Z$

$\Sigma V^*Z$

Bunny Compression



Figure 2.6: Caption

Covariance – notice the demean step

$$C_I = \frac{1}{n-1}\left(X - \tilde{X}\right)^T\left(X - \tilde{X}\right)$$

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

$$|a_i(t)| \xrightarrow{\ i\to\infty\ } 0$$

as fast as possible.

$$I = U\Sigma V^T$$
$$\approx U\hat{\Sigma}V^T$$

```matlab
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



Rank $r = 2$

Rank $r = 10$

Rank $r = 20$

Rank $r = 100$

$$e_i = \frac{\|I - I_i\|_F}{\|I\|_F}$$

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

Code 2.1: Read, convert, and display images.

```matlab
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')
```

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.

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

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

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

$$\overline{U} = (\dot{U}_1 \ \dot{U}_2 \ \cdots \ \dot{U}_m)$$

**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:

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

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

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

### Some Special (Simple) Matrix Norms

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

- For $p = 1$:

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

- For $p = \infty$:

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

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

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

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

- The Frobenius norm is given by:

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

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

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



$A x = x'$

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

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

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

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

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

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

**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:

$$\|A\|_p = \sup_{x \neq 0} \frac{\|Ax\|_p}{\|x\|_p} \tag{2.14}$$

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

$$\|A\|_p = \sup_{\|x\|_p = 1} \|Ax\|_p \tag{2.15}$$

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

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

### Some Special (Simple) Matrix Norms
The first 3 of these are induced norms, but the 4th is not.
- For $p = 1$:

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

- For $p = \infty$:

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

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

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

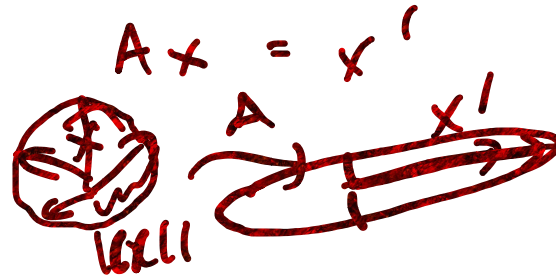where $\sigma_{max}$ is the maximum singular value of the matrix $A$.
- The Frobenius norm is given by:

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

Fun facts about matrix estimation (late estimation)

If $A$, $\quad \sigma_1 \geq \cdots \geq \sigma_r > \sigma_{r+1} = 0$

- range$(A) = $ span$(U_1, U_2, \cdots, U_r)$
- null $(A) = $ span$(V_{r+1}, V_{r+2}, \cdots, V_n)$

$$\left. \right\} \; 0 \rightarrow \overset{A}{\diagdown}$$

$\sigma_2 = 0$

$r = 1$

- $\|A\|_2 = \sigma_1 \quad ; \quad \|A\|_F = \sqrt{\sigma_1^2 + \sigma_2^2 + \cdots + \sigma_r^2}$

- $A = \sum\limits_{i=1}^{r} \sigma_i U_i V_i^{\top} = \sigma_1 U_1 V_1^{\top} + \sigma_2 U_2 V_2^{\top} + \cdots + \sigma_r V_r^{\top}$

  rank$-1$ outer products

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

$\begin{pmatrix} U_1 & U_2 & \cdots & U_n \end{pmatrix}_{\!\!(m \times q)} \begin{pmatrix} \sigma_1 \sigma_2 \cdots \sigma_r \end{pmatrix} \begin{pmatrix} - V_1^{\top} - \\ - V_2^{\top} - \\ \vdots \end{pmatrix}^{1 \times n}$

$\omega_2 \rightarrow \omega_2$

$\begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix}$

$\omega_1^{\top} \omega_2 = \omega_1 \cdot \omega_2$

$= \|\omega_1\| \|\omega_2\|$

$\cos \theta$

# 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?)

$\sigma_i$

$e_N$

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

On PCA Principal Component Analysis, Eigenface

-On Raleigh Ritz Quotient

-On Spectral Decomposition Theorem

-On Data Clouds



$$\overline{X} = [X_1 | X_2 | \dots | X_n]$$

$$\vec{x_i} = a_1^i \vec{v_1} + a_2^i \vec{v_2} + \dots + a_r^i \vec{v_r}$$

$\vec{u_i}$'s are basis set.

$x_i$ = PCA gives basis set

where $\langle a_j \rangle_i \downarrow$ as fast as possible vs any other basis set.

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

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

Ex. $\underline{X_i} \sim 4900 \times 1$ gene expression table for each $i$.

$i = 1 \ldots 216$ patients

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

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

$$f: \mathbb{R}^{4000} \longrightarrow \mathbb{Z}_2 \qquad \mathbb{Z}_2 = \{0, 1\}.$$

$x \in \mathbb{R}^{4000}$

- Supervised vs. unsupervised.

unsupervised — just input — just structural geometry.



just $x$

← "Data Cloud" ~ distribution of R.V. $x \sim \underline{X}$

level sets of $p$

$p(x): \mathbb{R}^{4000} \longrightarrow \mathbb{R}^+$

- supervised learning is descriptive $f: x \rightarrow y$ ← labels

# THE SPECTRAL DECOMPOSITION

Let A be a $n \times n$ symmetric matrix. From the spectral theorem, we know that there is an orthonormal basis $u_1, \cdots, u_n$ 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,

$$A u_j = \lambda_j u_j. \quad \longleftarrow \text{ real}$$

Then

$$A = PDP^{-1} = PDP^\mathsf{T}$$

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

$$A = [u_1 \cdots u_n] \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} \begin{bmatrix} u_1^\mathsf{T} \\ \vdots \\ u_n^\mathsf{T} \end{bmatrix}$$

$$= [\lambda_1 u_1 \cdots \lambda_n u_n] \begin{bmatrix} u_1^\mathsf{T} \\ \vdots \\ u_n^\mathsf{T} \end{bmatrix}$$

outer

$n \times 1 \cdot 1 \times n$
$n \times n$

$$= \lambda_1 u_1 u_1^\mathsf{T} + \cdots + \lambda_n u_n u_n^\mathsf{T}.$$

The expression

$$A = \lambda_1 u_1 u_1^\mathsf{T} + \cdots + \lambda_n u_n u_n^\mathsf{T}.$$

is called the spectral decomposition of A. Note that each matrix $u_j u_j^\mathsf{T}$ 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 $P(x) = u_j u_j^\mathsf{T} x$ is the orthogonal projection onto the subspace spanned by $u_j$.

○ $A = B^\mathsf{T} B$ is symmetric

⟹ spectral decomp. theorem

i.e. also covariance matrices.

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

$\|v_i\|^2 = v_i \cdot v_i = v_i^\mathsf{T} v_i$ scalar - inner product

# PCA as algorithm

- Data $\underline{X} = \begin{pmatrix} | & | & & | \\ x_1 & x_2 & \cdots & x_n \\ | & | & & | \end{pmatrix}^{ith}_{m \times n}$

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

covariance matrix.

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

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

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

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

$u_2$ is first minor axis

(Sigma = 4000
n = 216)



$\frac{1}{\sqrt{2}}$

Side Note:

$\frac{1}{i^2}$ is slowest converging to zero $\sigma_i^2$

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

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

$$\sum_{i=1}^{\infty} \frac{1}{i^p} < \infty \quad \text{if prob bad.}$$

$p < 1, \quad p = 1 \text{ harmonic}$

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

$\sigma_i^2$ the power spectrum

$\sigma_i \leq \frac{1}{i} \in$

$$B = \bar{I} \bar{X}^T \qquad ; \qquad let \quad C = \frac{1}{n-1} B^T B \quad \overline{\text{everywhere}}$$

is covariance matrix

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

$$U = [\, U_1 \mid U_2 \dots U_n \,]$$



$\sigma_c U_c$

$\sigma_2 U_2$

$b_3 \sim 0$

$b_2 \gtrsim 0$

$\circ \quad U_1 = \operatorname{argmax}_{\|U\|=1} U^T B^T B \, U = \operatorname{argmax}_{U} (U^T B^T B U)$

$\underline{\|BU\|_2^2}$

$(U^T U)$

Raleigh - Ritz quotient -

$Bu \cdot Bu = \|Bu\|_2^2$

$\circ \quad U_2 = \operatorname{argmax}_{\substack{\|U\|=1 \\ U \perp U_1}} U^T B^T B U$

$$\begin{aligned} ||x_i - proj_w x_i||_2^2 \quad &= \quad ||x_i(w \cdot x_i)w||_2^2 = (x_i - (w^T x_i)w)^T (x_i - (w^T x_i)w) \\ &= \quad (x^T x_i) - (w^T x_i)^2 = ||x_i||_2^2 - (w^T x_i)^2. \qquad (2.43) \end{aligned}$$

To minimize this residual with respect to the unknown vector $w$, averaged across the data set, it is sufficient to maximize the second term since the first term does not depend on $w$. Thus we wish to maximize,

$$\mathcal{L}_1(\mathcal{D}; \Theta) = \frac{1}{n} \sum_{i=1}^{N} (w^T x_i)^2, \qquad (2.44)$$

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

$$o \quad CV = VD \qquad \text{eigenvectors of } C \text{ all} \quad -$$
$$\text{stacked.}$$

$$Ax \qquad \text{optimize} \qquad x^T A x \quad , \quad \text{maximize.}$$

$$x = \begin{pmatrix} x^1 \\ \vdots \\ x^n \end{pmatrix}_j \quad r(x) = \frac{x^T A x}{x^T x} \qquad \not{A} \qquad \begin{pmatrix} A = B^T B \\ x = 0 \end{pmatrix}$$

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

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

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

$$Ax = r(x) \, x \qquad \Rightarrow \qquad Ax = \lambda x$$

$$\underset{\lambda}{\underbrace{}}$$

Conclude

The $x$ that optimizes $r(x) = \dfrac{x^T A x}{x^T x}$

is an eigenvector and $r(x)$ is its

eigenvalue.

• S.D.T. for $A = B^T B = \sum_{i=1}^{n} \lambda_i U_i U_i^T$ ← ( $U_i$ ) are the $v_s$!

let $B = \widehat{U} \Sigma V^T$ svd.

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

$$D = \begin{pmatrix} \sigma_1 & & \\ & \ddots & \sigma_r \\ & & \ddots & 0 \\ & & & & 0 \end{pmatrix} \begin{pmatrix} \sigma_1 & & & \\ & \ddots & \sigma_r & \\ & & \ddots & 0 \\ & & & 0 \end{pmatrix} = \begin{pmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & \ddots & \\ & & \sigma_r^2 & \\ & & & 0 \end{pmatrix}$$

$= V \underline{D V^T} = \sum_{i=1}^{l} \sigma_i^2 \, v_i v_i^T$

$$\mathcal{L}_1(\mathcal{D}; \Theta) = \frac{1}{n}(X^T w)^T (X^T w) = \frac{1}{n} w^T (XX^T) w, \tag{2.45}$$

and the matrix $\frac{1}{n}(XX^T)$ is familiar in statistics as a covariance matrix. To optimize $\mathcal{L}_1$, subject to a constraint,[12]

$$\|w\|_2 = 1, \tag{2.46}$$

we can use the Lagrange multiplier method by defining an expanded loss function(cost function) with the equality constraint built in with a Lagrange multiplier. Let,

$$\mathcal{L}(\mathcal{D}; \Theta, \lambda) = \mathcal{L}_1(\mathcal{D}; \Theta) - \lambda(w^T w - 1) = \frac{1}{n} w^T (XX^T) w - \lambda(w^T w - 1). \tag{2.47}$$

To minimize this, we take derivatives and set them equal to zero.

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{2}{n} XX^T w - 2\lambda w \implies \frac{1}{n}(XX^T)w = \lambda w \qquad (2.48)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = w^T w - 1 \implies \|w\|_2 = 1. \qquad (2.49)$$

**Theorem 2.2.1 — PCA foundations.** Let A be a symmetric $d \times d$ matrix. Then its (real) eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$, associate with orthogonal eigenvectors $w_1, w_2, \dots, w_d$. Furthermore,

$$\lambda_1 = \max_{\|w\|=1} w^T A w, \text{ with } w_1 = \arg \max_{\|w\|=1} w^T A w.r \qquad (2.50)$$

$$\lambda_2 = \max_{\|w\|=1, w \perp w_1} w^T A w \text{ with } w_2 = \arg \max_{\|w\|=1, w \perp w_1} w^T A w.$$

$$\vdots$$

$$\lambda_d = \max_{\|w\|=1, w \perp w_1, w_2, .., w_{d-1}} w^T A w \text{ with } w_d = \arg \max_{\|w\|=1, w \perp w_1, w_2, .., w_{d-1}} w^T A w.$$

**Theorem 2.2.2 — Spectral Decomposition.** If $A$ is a symmetric positive semi-definite matrix, then there is an orthogonal set of eigenvectors, $u_i$, each with non-negative eigenvalues, $\lambda_i \geq 0$. Furthermore, the decomposition of $A$ has the following representation by rank one matrices that describe the action of $A$ as a weighted sum of simple projections onto the subspaces spanned by each $u_i$,

$$A = \sum_{i=1}^{N} \lambda_i u_i u_i^T \tag{2.51}$$

**Which functions are most efficient?**

That is, we write a linear combination,

$$u(x,t) = \sum_k a_k(t)\varphi_k(x), \tag{2.57}$$

of functions $\varphi_k(x)$, where the time varying (component projection) values,

$$a_k(t) = \frac{(u(x,t), \varphi_k(x)}{\|\varphi_k(x)\|_2}, \tag{2.58}$$

or better yet, we can effectively skip the denominator by choosing the basis set of functions such that,

$$\|\varphi_k(x)\|_2^2 = (\varphi_k(x), \varphi_k(x)) = 1. \tag{2.59}$$

**Given a spatiotemporal data** sample as an array, (for example, typically from a solution derived from a computational solver for a PDE):

$$\mathbf{U} = \left( \begin{array}{cccc} | & | & & | \\ u(\vec{x}, t_1) & u(\vec{x}, t_2) & \ldots & u(\vec{x}, t_T) \\ | & | & & | \end{array} \right). \tag{2.61}$$

then develop the demeaned array $U - \overline{U}$. **Find**:

$$\mathbf{\Phi} = \left( \begin{array}{cccc} | & | & & | \\ \varphi_1(x) & \varphi_2(x) & \ldots & \varphi_k(x) \\ | & | & & | \end{array} \right). \tag{2.62}$$

by the singular value decomposition. This basis yields the fastest decaying power spectrum, *in time average*, versus all other possible basis.

$$
\begin{aligned}
a(t) &= \frac{(u, \varphi)}{\|\varphi\|^2} = \frac{\|u\|\|\varphi\|\cos\theta}{\|\varphi\|^2} = \|u\|\cos\theta \\
&= \int_\Omega u(x, t)\varphi(x)dx
\end{aligned}
\tag{2.63}
$$

when $\|\varphi\| = 1$.

So, we will also write time average using brackets $< \cdot, \cdot >$, so define:

$$
\begin{aligned}
< |a(t)| > &= \frac{1}{T}\int_0^T |a(t)|dt \\
&= \frac{1}{T}\int_0^T |(u, \varphi)|dt \\
&= \frac{1}{T}\int_0^T \left|\int_\Omega u(x, t)\varphi(x)\right| dt
\end{aligned}
\tag{2.64}
$$

> $\star\star$ The goal is to choose a basis with fastest decaying power spectrum, in time average $\star\star$

Our goal can be summarized by the following loss function.

$$
\mathcal{L}(\varphi) = \frac{< |(u, \varphi)|^2 >}{\|\varphi\|^2}.
\tag{2.65}
$$

This very compact notation, encodes two integrations. We remind that the round brackets describe the inner product, $(f, g)$, meaning integration in the "space" variable. Now we have introduced the pointy brackets to describe time average. So,

$$
\mathcal{L}(\varphi) = \frac{\frac{1}{T}\int_0^T |\int_\Omega u(x, t)\varphi(x)dx|^2\, dt}{\int_0^L \varphi^2(x)dx} \equiv < a\varphi^2(t) >
\tag{2.66}
$$

or

$$
\max_{\|\varphi\|=1} \frac{1}{T}\int_0^T \left|\int_0^L u(x, t)\varphi(x)dx\right|^2 dt
\tag{2.67}
$$

**Theorem 2.3.1 — Parseval's Like Idenfity.** If $f \in L^2([0, L])$, then if $\{\varphi_k(x)\}$ is an orthonormal basis set, then:

$$\|f\|_2^2 \leq \sum_{k=0}^{\infty} |a_k|^2 \tag{2.68}$$

where $a_k = f, \varphi_k)$.

$$\frac{\partial}{\partial \delta}\mathcal{L}(\varphi + \delta\psi) = 0:$$

$$\mathcal{L}(\varphi) = <|(u,\varphi)|^2> -\lambda\left(\|\varphi\|^2 - 1\right)$$

$$\begin{cases} \max & <|(u,\varphi|> \\ \|\varphi_1\| = 1 \\ \varphi \in \mathcal{H} \end{cases}$$

$$C\vec{\varphi}(x) = \lambda\vec{\varphi}(x), C = \mathbf{U}\mathbf{U}^T,$$

and $\varphi_2(x)$ solves:

$$\begin{cases} \max & <|(u,\varphi|> \\ \|\varphi_2\| = 1 \\ \varphi \in \mathcal{H} \\ \varphi_2 \perp \varphi_1 \end{cases}$$

---

**Theorem 2.3.2 — Spectral Decomposition.** If $A$ is a symmetric positive semi-definite matrix (i.e. $\forall u \in \mathbb{R}^n, u^T A u \geq 0$), then there is an orthogonal set of (column) eigenvectors, $v_i$, each with non-negative eigenvalues, $\lambda_i \geq 0$, and furthermore the decomposition of $A$ as the following rank one matrices describes the action of $A$ as a weighted sum of simple projections onto the subspaces spanned by each $v_i$,

$$A = \sum_{i=1}^{N}\lambda_i v_i v_i^T \tag{2.71}$$

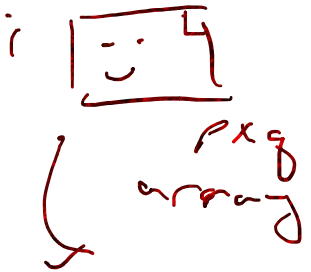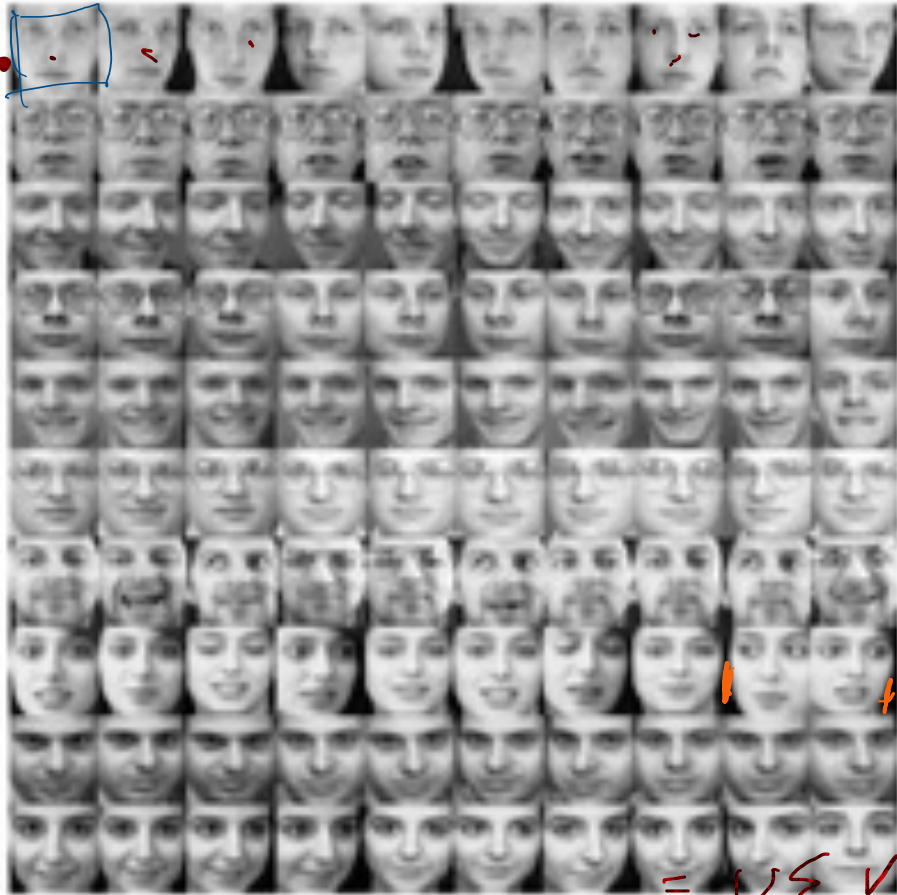Solve $\mathbf{U}\mathbf{U}^t$ for eigs; solve fastest decaying time averaged power spectrum

$$\mathbf{U}\mathbf{U}^t V^t = \Lambda V^t, \tag{2.74}$$

we use $\mathbf{U} = U\Sigma V^t$, with $\Lambda = \Sigma^2$

Eigenface



$;12$ picture

$p = 1000, q = 2000$

08/31/20

○ Remove correlated variance

Register

$E_3$

$\mathbb{R}^{p \cdot q} = \mathbb{R}^r$

$E_2$

$E_1$

$x_i \cdot e_i = x_i^T e_i = E$

$i$ $p \times q$ array

reshape as vector

$e_i = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$

$\begin{bmatrix} x_i \end{bmatrix}$

$M \times 1$

$M = pq$

$X = [x_1, x_2, \ldots, x_N]_{M \times N} = U \Sigma V^T$

# Eigenfaces for Face Detection/Recognition

## • 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 $N^2 \mathrm{x} 1$ vector, corresponding to an $N \mathrm{x} N$ face image $I$.

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

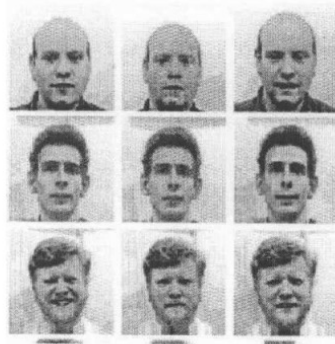$$\hat{\Phi} - mean = w_1 u_1 + w_2 u_2 + \cdots w_K u_K \ (K << N^2)$$

# Computation of the eigenfaces

Step 1: obtain face images $I_1, I_2, ..., I_M$ (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 $C$:

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

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

$C = X^T X$

<u>Step 6</u>: compute the eigenvectors $u_i$ of $AA^T$

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

<u>Step 6.1</u>: consider the matrix $A^T A$ ($M$x$M$ matrix)

<u>Step 6.2</u>: compute the eigenvectors $v_i$ of $A^T A$

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

<u>What is the relationship between $us_i$ and $v_i$?</u>

$$A^T Av_i = \mu_i v_i => AA^T Av_i = \mu_i Av_i =>$$

$$CAv_i = \mu_i Av_i \text{ or } Cu_i = \mu_i u_i \text{ where } u_i = Av_i$$

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

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

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

<u>Note 3</u>: 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).

<u>Step 6.3</u>: compute the $M$ best eigenvectors of $AA^T$: $u_i = Av_i$

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

<u>Step 7</u>: 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:

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

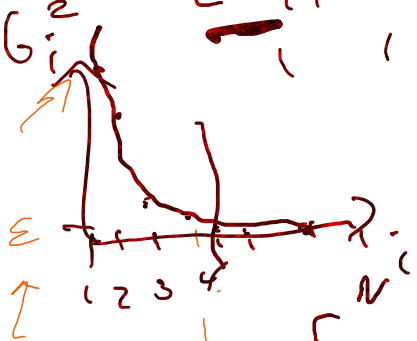(we call the $u_j$'s *eigenfaces*)



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

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

Bunny

$$j = \boxed{\phantom{bunny matrix}} \; r \times n \longrightarrow X = [\; X_1 \mid X_2 \mid ... \mid X_N \;]$$

$\underbrace{\phantom{xxx}}$
$x_1 \; x_2 \quad x_N$

← one column

$$= U \Sigma V^T$$

PCA of bunny



$$X_i = a_1^i U_1 + a_2^i U_2 + ... + a_N^i U_N$$

bases
full of as
possible for
the data set

$$a_1^i = \boxed{X_i^T U_1} = comp_{U_1} X_i = \frac{\boxed{U_1 \cdot X_i}}{|U_1|} = \frac{|U_1||X_i|\cos\theta}{1}$$

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

leadin
3 B ② — in 5 min. Signals analysis, Harmonic.

○ Historically favorite basis set. $B = \{\omega_1, \omega_2, ...\}$

(us. energy favorite basis set comes from PCA)

○ Taylor Polynomials. — $f(x) = a_0 + a_1 x + a_2 x^2 + ... + a_r x^r$

$B = \{x^0, x^1, x^2, ...\}$

$a_i = \dfrac{f^{(i)}(0)(x-0)^i}{i!}$

○ Fourier modes.

Hilbert space $f(x) = a_1 \sin x + a_2 \sin 2x + a_3 \sin 3x + ...$

$B = \{\sin x, \sin 2x, ...\}$

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

○ Wavelet basis.

○ chebychev polys, legendre ...

$\sin x = 1 + x - \dfrac{x^3}{3!} + \dfrac{x^5}{5!} + ...$

$\text{comp}_x \sin x = 1$

# Changing bases is a sort of coord. rot.

$\mathcal{H}$



$\phi_1$    $x$    $\rightarrow \phi_0$

$\phi_2$

$\mathcal{H}$ $\psi_1$    $\xi_1$

$x$   $\vec{x}'$

$\xi_2$

$\xi_3$

$\phi_0(y) = y^0$

$\phi_1(y) = y^1$

$\phi_2(y) = y^2$

$\vdots$

$\rightarrow \psi_1(y) = \sin(y)$

$\psi_2(y) = \sin(2y)$

$\psi_3(y) = \sin(3y)$

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

vs. 2.1

Hilbert space — a complete inner product space.

- An inner product space is a "vector space" $E$ together with ~that limits
  a function called "inner product" $\langle \cdot, \cdot \rangle : E \times E \to \mathbb{C}$ / $\mathbb{R}$.
  with properties.

Gift: you set geometry in $E$
as an angle $\langle (u,v) := \frac{\langle u,v \rangle}{\|u\| \|v\|}$ cos
○ and projection.
- vector space.

  - Conjugate $\langle u, v \rangle = \overline{\langle v, u \rangle}$ $\forall u, v \in E$ (symm.)
  - linear $\langle au_1 + bu_2, v \rangle = a\langle u_1, v \rangle + b\langle u_2, v \rangle$ $\forall a, b \in \mathbb{C}$ field.
  - pos. def: $\langle u, u \rangle > 0$ $\forall u \in E, u \neq 0$. $u_1, u_2, v \in E$

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

Ex: arrays of real numbers that are $\infty \times 1$
$\quad 3 \begin{bmatrix} 2 \\ 1 \end{bmatrix} + 6 \begin{bmatrix} 4 \\ 1 \end{bmatrix} = \begin{bmatrix} 30 \\ 15 \end{bmatrix}$

Norm: $\|x\|^2 = \langle x, x \rangle^{1/2}$

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

Ex: Vector space of functions
Name an inner product space $E = L^2([0,1]) = \{ f \mid \int_0^1 |f(x)|^2 dx < \infty \}$
$( L^2([0,1]) \subset C([0,1]) )$

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

Basis of unit vectors. $B = \{ \phi_i \mid \phi_i \in B \subset E \text{ and } \langle \phi_i, \phi_j \rangle = \delta_{ij}$
- orthogonal $\quad$ and $\|\phi_i\| = 1 \}$
$\quad u = \sum a_i \phi_i$ (projection)

Ex: $L([0,\ell])$ ; $B = \{ \cos \frac{2\pi k x}{\ell} ; \sin \frac{2\pi k x}{\ell}, \frac{1}{\sqrt{\ell}} \mid k \in \mathbb{N} \} = \{ -\infty, \ldots \}$
if $\ell = 1$ ; $f(x) = \frac{a_0}{2} + \sum_{k=1}^\infty a_k \cos kx + b_k \sin kx$ ; $a_k = \frac{\langle f, \cos kx \rangle}{\|\cos kx\|} = \frac{1}{\pi} \int_{-\pi}^\pi f(x) \cos kx \, dx$
$\quad b_k = \quad \sin kx$

$f(x) = \frac{a_0}{2} + a_1 \quad + b_1 \quad + a_2 \quad + b_2 \quad + \cdots$

$\phi_k(x) = \frac{2\cos kx}{\ell} \cdot \frac{1}{\ell}$ or $\frac{2}{\ell} \cdot \sin \frac{\pi k x}{\ell}$.
○ similar if $\phi_k(x) = c e^{ikx}$ $\quad$ finite if "trig poly"

or maybe $B' = \{ x^k \cdot \}$ Taylor poly.

or maybe $B'' = \{ \quad \}$ wavelets.

"Hat" 2.13
functions 1st

Finite dimensional inner product space of functions.

$x_1 < x_2 < \ldots < x_n$ a grid on $[0, \ell]$
$G = \{ x_i \}_{i=1}^n$

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

Ex: let $\langle \vec{f}, \vec{g} \rangle = \sum_{k=1}^n f_k \overline{g_k}$ $\quad$ ○ 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. $L^2(\mathbb{R})$

$\odot \ B = \{ \ -, \ \sim, \ \dot{\sim}, \ \dots, \ \sim, \ \sim, \ \dots \} \quad \underset{Fourier}{[0,2\pi]}$

$= \{ 1, \ \sin x, \ \cos x, \ \sin 2x, \ \cos 2x, \ \dots \} \quad \Rightarrow$
Fourier

$\odot \ B' = \{ -, \ /, \ \smile, \ \frown, \ \vee, \dots \} = \{ 1, x, x^2, x^3 \dots \}$

Taylor $\uparrow \uparrow \uparrow$

... Legendre,

Taylor.

$\odot \ B'' = \{$  $, \ \dots \}$  $\phi_3(x)$

Wavelet.

$\theta_c$

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

$f(x) = \sum a_i \ \phi_i(x) ;$  $a_i = \langle f, \ \phi_i(x) \rangle$

$\odot \ \underset{Haar \dots}{\frown}$  $\dfrac{}{\| \phi_i(x) \|}$

On Compressed Sensing and on to Sparsity

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

$$= \sum_{n=1}^{\infty} a_n \sin nx \qquad \frac{\sin kx}{d}$$

$$\hat{f}(x) = \begin{pmatrix} \boxed{a_1} \\ \boxed{a_3} \\ a_7 \\ \vdots \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 3 \\ 0 \\ 0 \\ 0 \\ 0 \\ 4 \\ 0 \end{pmatrix}$$

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

$$\hat{f}_\tau = \begin{pmatrix} 58 \\ 0 \\ c_3 \\ 0 \\ c_5 \end{pmatrix}$$

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

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

$$\Rightarrow 38x - \left( \frac{1}{3!} + \frac{3^3}{c_3 \ 2!} + \frac{4 \cdot 7^3}{3!} \right) x^3 + \cdots$$

$\mathcal{H}$

$\sin 5x$

$\phi_3$

$\phi_2$

$\sin 2x$

$\phi_1$

$\sin x$

$f = \underline{\sin x} + \underline{3 \sin 3x} + \underline{4 \sin 7x}$

$\mathcal{H}$

$s$

- 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.

(a)

Mass-specific metabolic rate (ml $O_2 \cdot g^{-1} \cdot h^{-1}$)

Harvest mouse

Kangaroo mouse

Cactus mouse

Mouse

Flying squirrel

Rat

Cat

Rabbit    Dog  Sheep   Elephant

Horse

Human

Mass (kg)

$$y_i = \beta_0 + \beta_1 x_i$$

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

(a)

$$y_i = \beta_0 + \beta_1 x_i$$

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

$$
\begin{aligned}
y_0 &= \beta_0 + \beta_1 x_0 + \epsilon_0 \\
y_1 &= \beta_0 + \beta_1 x_1 + \epsilon_1 \\
&\vdots \\
y_{N-1} &= \beta_0 + \beta_1 x_{N-1} + \epsilon_{N-1}
\end{aligned}
$$

$$Y = X\beta + \epsilon$$

$$e_i = (f(x_i) - y_i)^2$$

$$
\begin{aligned}
E \;\; &= \;\; \sum_{i=1}^{N} e_i \\
&= \;\; \sum_{i=1}^{N} (f(x_i) - y_i)^2 \\
&= \;\; \sum_{i=1}^{N} (\beta_0 + \beta_1 x_i - y_i)^2 \,.
\end{aligned}
$$

$$e_i = (f(x_i) - y_i)^2$$

$$
\begin{aligned}
E &= \sum_{i=1}^{N} e_i \\
&= \sum_{i=1}^{N} (f(x_i) - y_i)^2 \\
&= \sum_{i=1}^{N} (\beta_0 + \beta_1 x_i - y_i)^2 .
\end{aligned}
$$

$$\frac{\partial E}{\partial \beta_0} = 0 \quad \text{and} \quad \frac{\partial E}{\partial \beta_1} = 0,$$

$$e_i = (f(x_i) - y_i)^2$$

$$
\begin{aligned}
E \;&=\; \sum_{i=1}^{N} e_i \\
&=\; \sum_{i=1}^{N} (f(x_i) - y_i)^2 \\
&=\; \sum_{i=1}^{N} (\beta_0 + \beta_1 x_i - y_i)^2 \,.
\end{aligned}
$$

$$\frac{\partial E}{\partial \beta_0} \;=\; 0 \quad \text{and} \quad \frac{\partial E}{\partial \beta_1} = 0,$$

$$\frac{\partial E}{\partial \beta_1} = \sum_{i=1}^{N} 2x_i \left( \beta_0 + \beta_1 x_i - y_i \right)$$

$$= \sum_{i=1}^{N} (2\beta_0 x_i) + \sum_{i=1}^{N} (2\beta_1 x_i^2) - \sum_{i=1}^{N} (2x_i y_i) = 0$$

and

$$\frac{\partial E}{\partial \beta_0} = \sum_{i=1}^{N} 2 \left( \beta_0 + \beta_1 x_i - y_i \right)$$

$$= \sum_{i=1}^{N} (2\beta_0) + \sum_{i=1}^{N} (2\beta_1 x_i) - \sum_{i=1}^{N} (2y_i) = 0.$$

From the above two equations we have:

$$\sum_{i=1}^{N} (x_i y_i) = \sum_{i=1}^{N} (\beta_0 x_i) + \sum_{i=1}^{N} (\beta_1 x_i^2)$$

$$\sum_{i=1}^{N} (y_i) = \sum_{i=1}^{N} (\beta_0) + \sum_{i=1}^{N} (\beta_1 x_i)$$

again, which can be written in matrix form as:

$$\begin{pmatrix} \sum_{i=1}^{N} (x_i y_i) \\ \sum_{i=1}^{N} (y_i) \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^{N} x_i & \sum_{i=1}^{N} x_i^2 \\ \sum_{i=1}^{N} 1 & \sum_{i=1}^{N} x_i \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$$
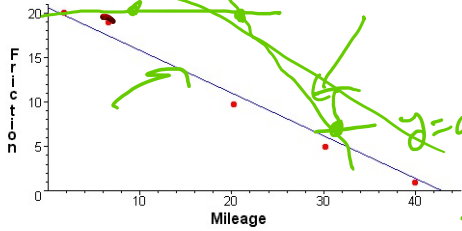
and then

$$\begin{pmatrix} b_0 \\ b_1 \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^{N} x_i & \sum_{i=1}^{N} x_i^2 \\ \sum_{i=1}^{N} 1 & \sum_{i=1}^{N} x_i \end{pmatrix}^{-1} \begin{pmatrix} \sum_{i=1}^{N} (x_i y_i) \\ \sum_{i=1}^{N} (y_i) \end{pmatrix}.$$

**Example**

An engineer is tracking the friction index over mileage of a breaking 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 **b** is the vector of friction index data values and **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 **b** and **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

$$Ax = b$$

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

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

We can conclude that the equation of the regression line is

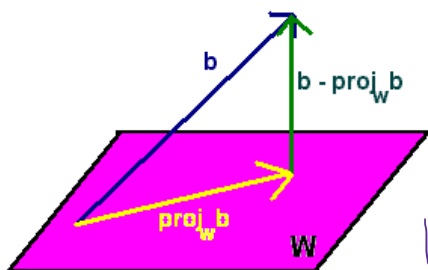$$y = -0.48x + 20.6$$

# Least Squares

Definition and Derivations

We have already spent much time finding solutions to

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

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

$$A\mathbf{x} = \text{proj}_W\mathbf{b}$$

where W is the column space of A.



Notice that **b** - 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

$$A^T(\mathbf{b} - A\mathbf{x}) = 0 \qquad A^TA\mathbf{x} = A^T\mathbf{b}$$

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

---

**Lemma**

If A is an m x n matrix of rank n, then $A^TA$ is nonsingular.

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

$m \times n$

m - eqns

n - unknowns

x

$$\implies a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1$$
$$\vdots$$
$$a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n = b_n$$

---

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

How many?

How many?



$x_2 \uparrow$ $\longrightarrow x_1$

col (A)

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

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

$$A x = b \implies x_1 \vec{a_1} + x_2 \vec{a_2} + \cdots + x_n \vec{a_n} = \vec{b}$$

$A x - b = 0$

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

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

$A\tilde{x} = \text{proj}_w b$ ; $\tilde{x} = \underset{x}{\text{argmin}} \|A x - b\|_2^2$

$\implies$ LS solution

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

vs.

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

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

$$\begin{pmatrix} - a_1 - \\ - a_2 - \\ \vdots \\ - a_n - \end{pmatrix} (A x - b) = 0 \leftarrow \begin{vmatrix} 0 \\ 0 \\ 0 \end{vmatrix}$$

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

$$\implies \vec{a_i} \perp (A\tilde{x} - b) \; \forall i$$

$$\implies (A\tilde{x} - b) \perp \text{every vector in Col}(A) \impliedby \text{Solve "normal eqns"}$$

$A^T A \tilde{x}$

## Theorem

Let A be an m x n matrix or rank n, then the system

$$Ax = b$$

has the unique *least squares* solution

$$x = (A^TA)^{-1}A^Tb$$

## 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

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

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.

---

*Handwritten annotations:*

$Ax = b$  
$m \times n$

$(A^TA)^{-1}$ — never form that

$m > n$ — tall skinny

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

· Square & nonsing.  
$x = A^{-1}b$  unique.  
$A^{-1}$ d-ne.

· Not, $m > n$

$A^TA x = A^Tb$  normal equations  
$n \times m$  $m \times n$

"covariance matrix"  
(if demeaned).

$y$

# 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 is it 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

$$\begin{pmatrix} 1 & 1 & 1 & 1 \\ 8 & 4 & 2 & 1 \\ 27 & 9 & 3 & 1 \\ 64 & 16 & 4 & 1 \\ 125 & 25 & 5 & 1 \\ 216 & 36 & 6 & 1 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \\ d \end{pmatrix} = \begin{pmatrix} 2.1 \\ 3.5 \\ 4.2 \\ 3.1 \\ 4.4 \\ 6.8 \end{pmatrix}$$
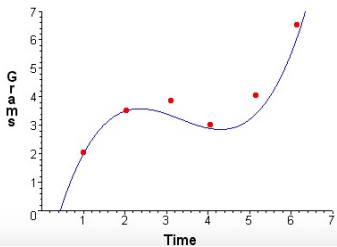
We can use the least squares equation to find the best solution

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

So that the best fitting cubic is

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

The graph is shown below



<!-- Handwritten annotations -->

Matrix Formulation of

LS slick

Grams = f(time) for general
y = f(x)  Models

$$2.1 = a \cdot 1^3 + b \cdot 1^2 + C \cdot 1 + d \qquad i = 1$$
$$3.5 = a \cdot 2^3 + b \cdot 2^2 + C \cdot 2 + d \qquad i = 2$$

$$A\vec{x} = b$$
$$X\beta = \vec{y}$$

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

$$\beta \qquad = y$$

$$\frac{X^T}{X^T A}$$

Top right handwritten:

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

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

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

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

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

$$\underbrace{\qquad}_{A^+}$$

$$\hat{\Sigma}^T\hat{\Sigma} = \begin{pmatrix} \sigma_1^2 & & & \\ & \ddots & & \\ & & \sigma_r^2 & \\ & & & \sigma_r \end{pmatrix}$$

when exist   $> 0$   $= 0$

Lower left handwritten:

$$A^TAx = A^Tb$$

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

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

$\hat{\Sigma}$

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

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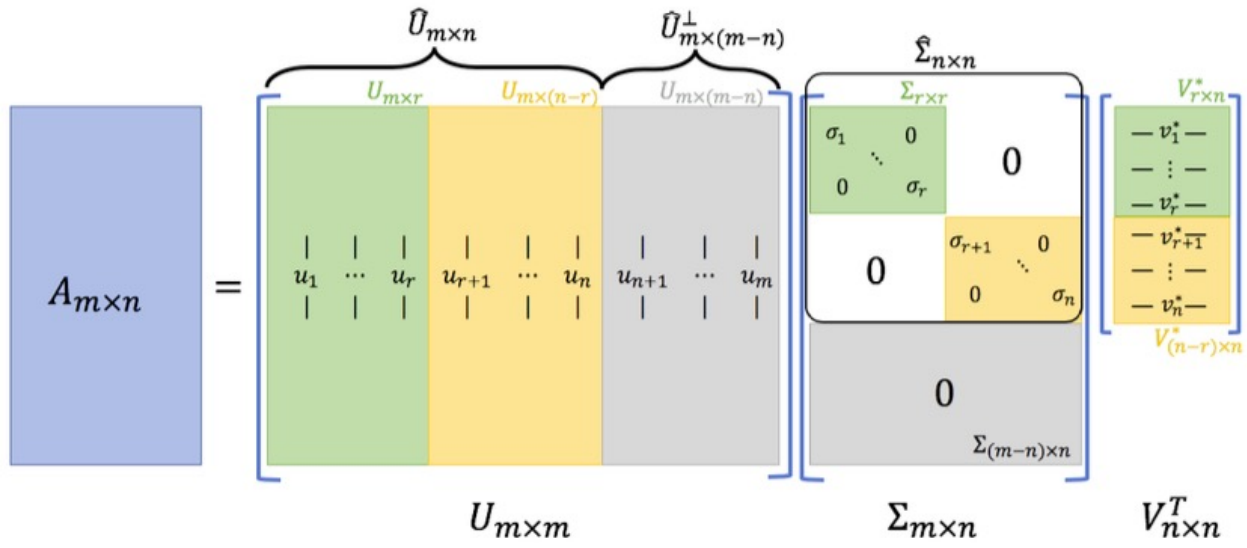
Handwritten notes (orange):

$$\left(\hat{\Sigma}^+\right) = \left(\hat{\Sigma}^T\hat{\Sigma}\right)^{-1}\hat{\Sigma}$$

if $\hat{\Sigma}$ is invertible

○ if $r = n$

if $\sigma_i > 0$ $\forall\, i < n$

But if $r < n$

$\sigma_r > 0,\ \sigma_{r+1} = 0$

$$\left(\hat{\Sigma}^T\hat{\Sigma}\right) = \begin{pmatrix} 1/\sigma_1^2 \dots 1/\sigma_r^2 & 0 \\ 0 & \end{pmatrix}$$

$1/\sigma_1''$  $0/\sigma \dots$

$= 1/\sigma''$

boxed: $= \frac{1}{\sigma}'' := 0$

$1/\sigma''$  $0/\sigma$

LS soln = solve normal equations

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

When inverse exists

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

$$= \boxed{A^+} b$$

Moore -
Penrose
; Pseudo-Inverse

$$A^+$$

$$Ax = b$$

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

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

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

$$Ax = b$$

$$U^T U \Sigma V^T x = U^T b$$

$$(\Sigma^T \Sigma) V^T x = \Sigma^T V^T b$$

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

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

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

$$x = V \Sigma^{-1} U^T b := V\Sigma^+ U^T b \quad \text{if exists}$$

$$\underbrace{\quad}_{A^+}$$

$$A^+$$

$$\Sigma^+ = \begin{pmatrix} 1/\sigma_1 & & \\ & \ddots & 1/\sigma_r \\ & & & 0 \cdots 0 \end{pmatrix}$$

$$\sigma_r \neq 0$$

$$\sigma_{r+1} = 0 :\Rightarrow \frac{1}{0} = 0$$

Exact sparsity vs. approximate

$$f(x) = \sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \frac{x^9}{9!} \Big/ - \frac{x^{11}}{11!} + \frac{x^{13}}{13!} + \dots$$

$\sin x$

$f_a(x)$  $N = 9$

$\phi_a(x) = \frac{x^n}{n}$

$$\sin x = \sum_{n=1}^{\infty} a_n \phi_n(x) \approx \sum_{n=1}^{\infty} a_n \phi_n(x)$$

$f(x)$  $f_N(x)$

$$\| f - f_N(x) \|$$

$$= \left( \left( \sin x - \left( x - \frac{x^3}{3!} \right) \dots \right) \right) \|$$

$$N = 3$$

$$\hat{f} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \frac{1}{6} \\ 0 \\ \frac{1}{120} \\ 0 \\ \frac{1}{5040} \end{pmatrix} \quad \text{vs.} \quad \hat{f}_N = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \frac{1}{6} \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

trunc