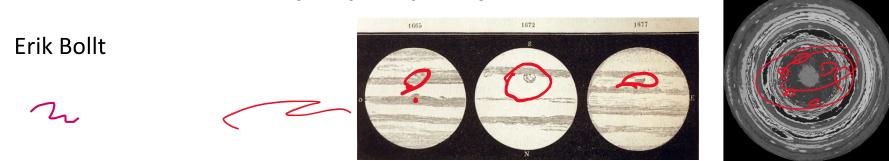
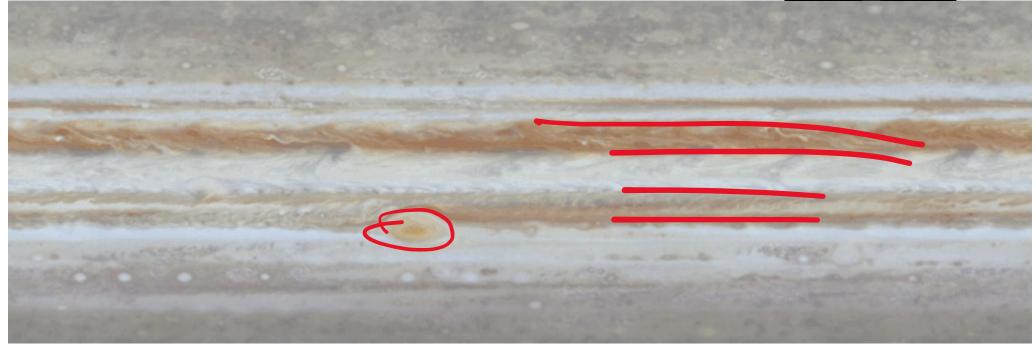
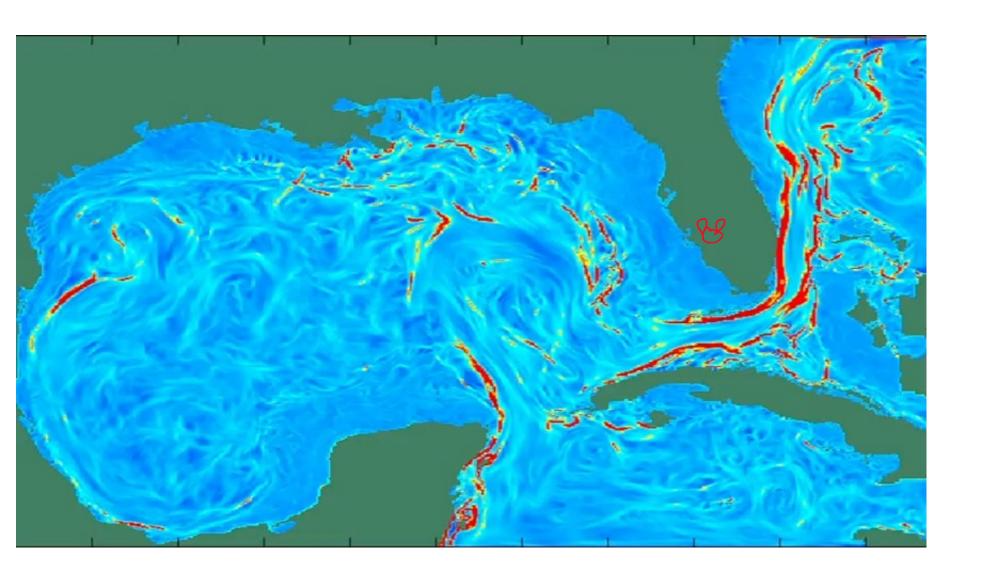
EE520 Data Driven Analysis of Complex Systems







Data as an array

On Matorix Multiplication $(\mathbf{z}): \mathbb{R}^n \to \mathbb{R}^m$ $\mathbf{z} \mapsto \mathbf{z}' = A\mathbf{z},$

$$A = \left(egin{array}{cccc} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ dots & \ddots & dots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{array}
ight), \ ext{and each } a_{i,j} \in \mathbb{C}$$

in terms of the usual matrix times vector multiplication,

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

$$n=2$$
 $\frac{7}{2}$

$$m = 3$$

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

2 - AZ new direction, new length. Eig for square-Cheracterize netrices by knowing Just o Cig. special directions · Aυ = A(α,ν, +α, νz) = α, Αν, +α, Ανz Def(A-)] = D = α, ε, ν, +α, ενz (A-)] = 0 ? Matrix x circle ?! S= {X | 11 X | [= 1, X \in E = 1 | R }; A \cdot S = {y \cdot 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 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:

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

where

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

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

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 ues of A are the scalar values, σ_i , and the columns of U and V have columns that are the corresponding i^{th} left and right singular vectors, u_i and v_i :

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

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

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,

$$4V = U\Sigma V^* V = U\Sigma, (2.8)$$

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

$$\begin{array}{lll} A & = & U \Sigma V^T \\ & = & \left(\begin{array}{ccc} \frac{1}{\sqrt{5}} & \frac{2}{\sqrt{5}} \\ \frac{-2}{\sqrt{5}} & \frac{1}{\sqrt{5}} \end{array} \right) \left(\begin{array}{ccc} \sqrt{70} & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{array} \right) \left(\begin{array}{ccc} \frac{1}{\sqrt{14}} & \sqrt{\frac{2}{7}} & \frac{3}{\sqrt{14}} \\ \frac{-3}{\sqrt{10}} & 0 & \frac{1}{\sqrt{10}} \\ \frac{-1}{\sqrt{35}} & \sqrt{\frac{5}{7}} & \frac{-3}{\sqrt{35}} \end{array} \right). \end{array} \eqno(2.28)$$

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:

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

$$\approx \begin{pmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \end{pmatrix}$$

$$[A] [v_1v_2\cdots v_n] = [u_1u_2\cdots u_n] diag(\sigma_1,\sigma_2,\cdots,\sigma_n).$$



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,

$$A_{m \times n} = \hat{U}_{m \times n} \hat{\Sigma}_{n \times n} V_{n \times n}^*, \tag{2.21}$$

eral SVD Eq. (2.2). $A_{m\times n} = \hat{U}_{m\times n} \hat{\Sigma}_{n\times n} V_{n\times n}^*,$ and $U = [\hat{U}_{m\times n} | \hat{U}_{(n-m)\times n}]$, written in terms of an orthogonal "buffer" matrix

6,7627...36,76,=0

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

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

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

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

 $A_{m \times n} = U_{m \times r} \Sigma_{n \times n} V_{n \times r}^{-},$ 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.

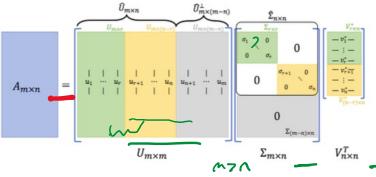


Figure 2.3: m > n tall skinny

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

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

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

$$A_{m \times n} \hat{V}_{n \times n} = \hat{U}_{m \times n} \hat{\Sigma}_{n \times n} \tag{2.25}$$

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:

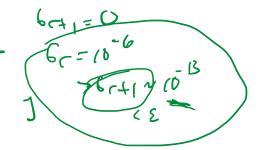
$$Av_1 = \sigma_1 u_1$$

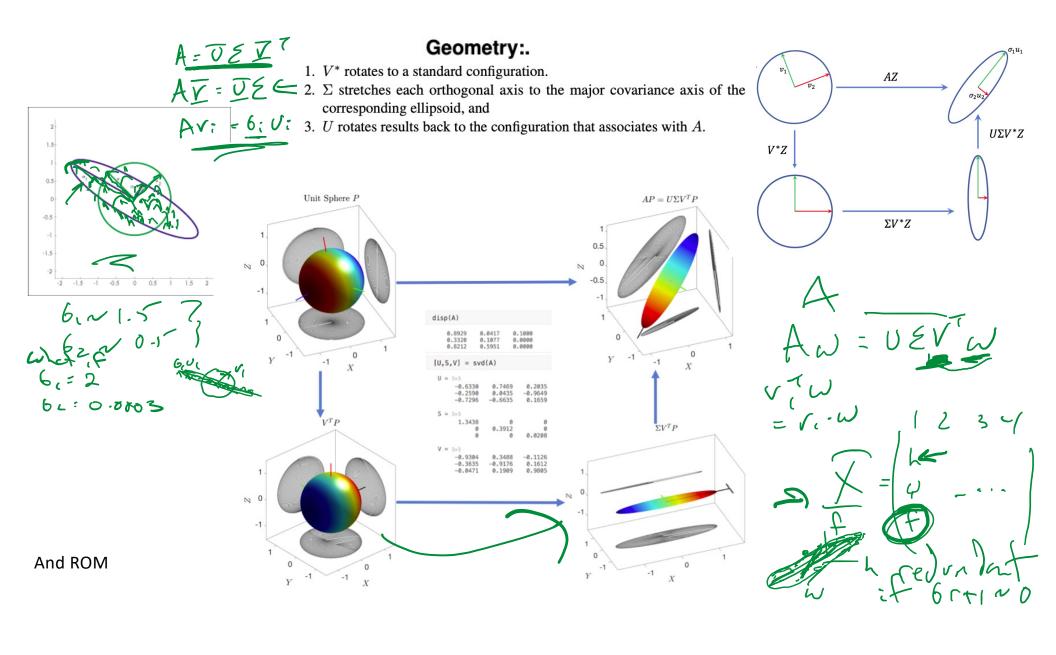
$$Av_2 = \sigma_2 u_2$$

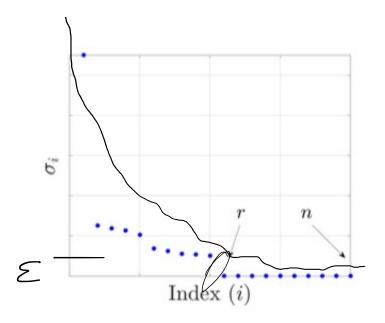
$$\vdots$$

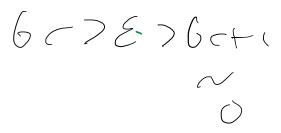
$$Av_n = \sigma_n u_n$$
(2.27)

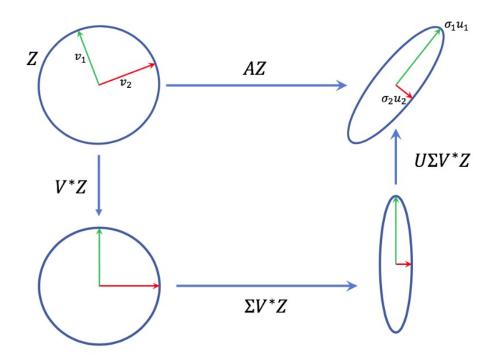
Full, Economy, Truncated SVD











Bunny Compression

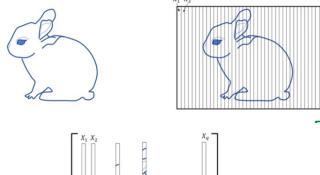
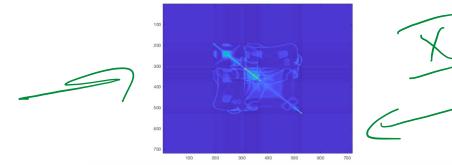


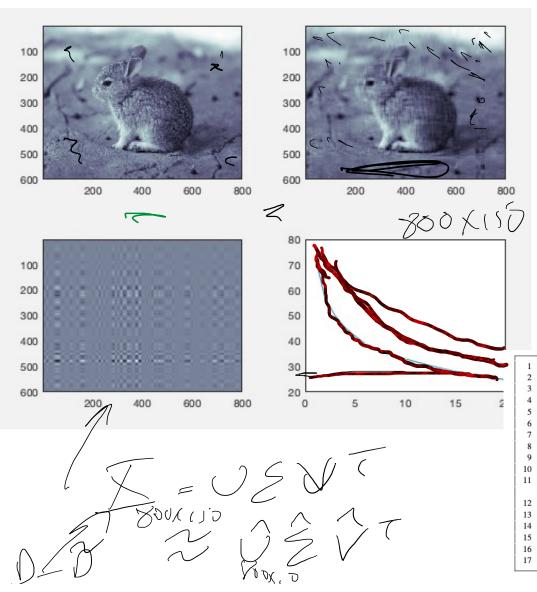
Figure 2.6: Caption



Covariance – notice the demean step

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

XXXXX



Ulxitle 2 a: lt/4:Lx1
[ailt] [i-700]
(ailt] (i-700)
cos fast as
possible.

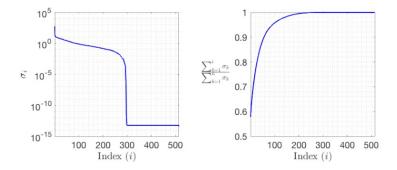
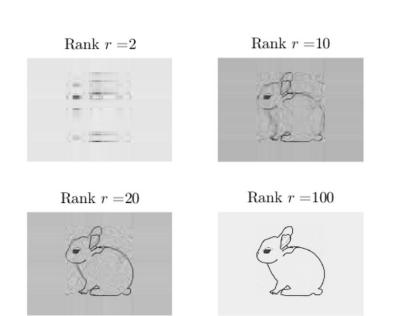
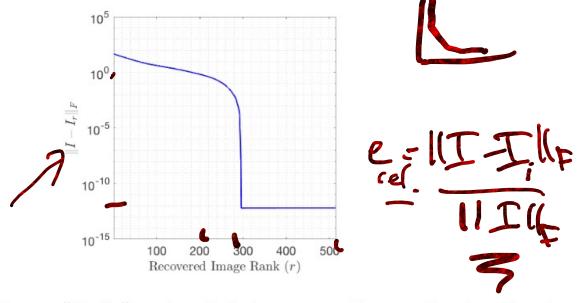


Figure 2.8: (Left) Singular Values. (Right) Energy





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

Code 2.1: Read, convert, and display images.

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

B Date Analysis

15 solving an ill-possed

- Optimize a cost function.

AAT = UEVETUT) - UEETUT (AAT) J= U(SST) = (SET) D

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

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

Some Special (Simple) Matrix Norms

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

• For p = 1:

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

• For $p = \infty$:

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

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

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

where σ_{max} is the maximum singular value of the matrix A.

• The Frobenius norm is given by:

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

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

: 11 (x, x) 11, = 14, (+1421

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

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

Some Special (Simple) Matrix Norms

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

• For p = 1:

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

• For $p = \infty$:

$$||A||_{\infty} = \max_{1 \le i \le m} \sum_{j=1}^{n} |a_{ij}|$$
 (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 σ_{max} is the maximum singular value of the matrix A.

· The Frobenius norm is given by:

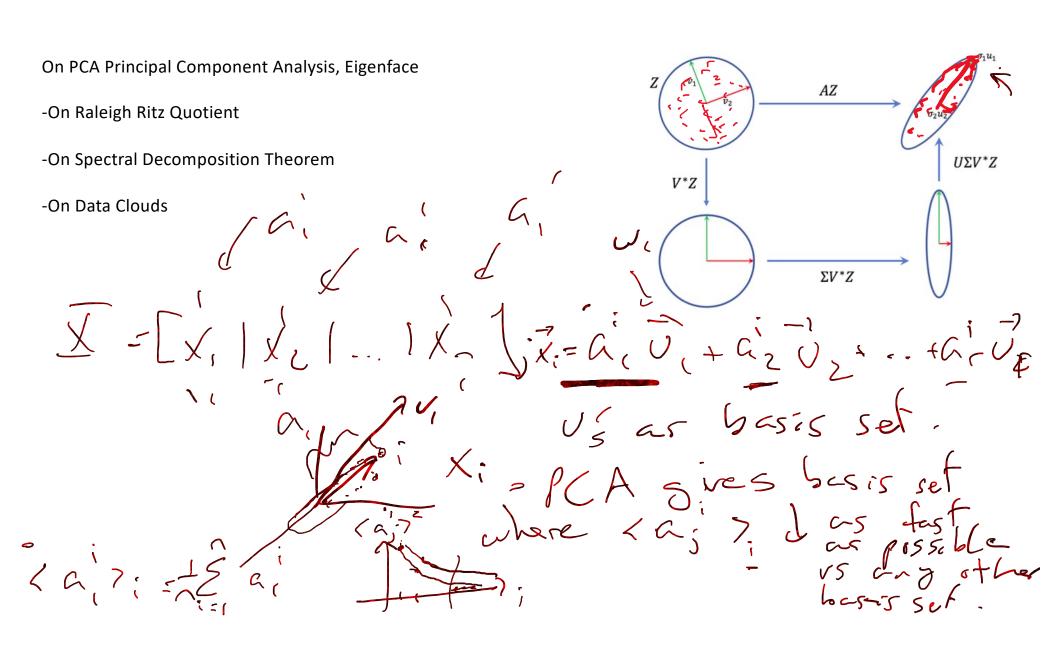
$$||A||_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2}$$
 (2.19)

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

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

Fun facts about matrix astronation (late estimation) $\frac{1}{2} A_{1} \qquad 6_{1} Z_{1} - 7_{6} - 7_{6} + 1_{1} = 0$ $\frac{1}{2} A_{1} \qquad 6_{1} Z_{1} - 7_{6} - 7_{6} + 1_{1} = 0$ $\frac{1}{2} A_{1} \qquad 6_{1} Z_{1} - 7_{6} - 7_{6} - 7_{6} + 1_{1} = 0$ $\frac{1}{2} A_{1} \qquad 6_{1} Z_{1} - 7_{6} - 7_{6} - 7_{6} + 1_{1} = 0$ $\frac{1}{2} A_{1} \qquad 6_{1} Z_{1} - 7_{6}$ e lAll z = 6, ; llAll = 56, 2+--+ 62 α A = 26:0: V= = 6, υ, ν + +6, ν + +6

Materix Estimation / Duta Estimation. Amon o let 65 NSC and AN = Z 6; U; V; * (so we very be stripping some if them ...



Date for PCA - "Pretend Pata lords like on ellipsoid"

Ex. X. ~ 4500 X (gene expression talob for each i.

i=(...216 petients

Xi = (xi

Yi = 0 or l

"O" of concer "" of cencer. f: (1R4000)

Z = \$0, 15. Supervised US. unsupervised.

Supervised - Sust input - Sust structure!

Sust Structure!

Sust Cloud ~ Distribution R.V. XN X

* supervised learning is descriptive t: t-> y Chy.

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 λ_j be the eigenvalue corresponding to u_j , that is,

 $Au_{j} \neq \lambda_{j}u_{j}.$ $A = PDP^{-1} = PDP^{T}$

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 $A = PDP^T$ can be rewritten as:

Then

14(. (X1)

The expression

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

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 .

of A=BB is symmetrice T spectral deemp. Hearon i.e. also covariance metrices. of is pos. letinite of 1:70 all I.

1/vill= Ui. Vi = UiTUi scalar = inner protoct

PCA as algorithm o Data = (X, Xz .- Xn) o aht.f. what if $x_i \sim n(x_i x_i)$ coverance metror. $B = X - B - B = \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \end{bmatrix} \times \begin{bmatrix}$ $\overline{X}_{i} = \frac{1}{\sqrt{2}} \sum_{i=1}^{\infty} \overline{X}_{ij}$ B = U & V ; U = [U, U2 . -. U_] Ui is nejor cxis - most vergetic

V2 ic first misor cx.s

i let C= I BTB everywood B = 1X covoriênce motore. B=UEVT U=LUIIUZ-Lond Ul = argmax UB'BU= Raleigh - R- #= g s ot vent - 1 Bu. Bu = ((B)) monex UTBISU 11011 -1 ULU

$$||x_i - proj_w x_i||_2^2 = ||x_i(w \cdot x_i)w||_2^2 = (x_i - (w^T x_i)w)^T (x_i - (w^T x_i)w)$$
$$= (x^T x_i) - (w^T x_i)^2 = ||x_i||_2^2 - (w^T x_i)^2.$$
(2.43)

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,

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

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

CV =

Conclude & That optimes run: XTAX
The x that optimes run: XXX

The x eigenvector and run is its
eigenvalue.

for A = BTB=

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

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

$$||w||_2 = 1, (2.46)$$

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.

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

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

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

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

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

:

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

Theorem 2.2.2 — Spectral Decomposition. If A is a symmetric positive semidefinite 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,

$$u(x,t) = \sum_{k} a_k(t)\varphi_k(x), \qquad (2.57)$$

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

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

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

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

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

$$\mathbf{U} = \begin{pmatrix} | & | & | \\ u(\vec{x}, t_1) & u(\vec{x}, t_2) & \dots & u(\vec{x}, t_T) \\ | & | & | \end{pmatrix}.$$
 (2.61)

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

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

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

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

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

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

$$\langle |a(t)| \rangle = \frac{1}{T} \int_0^T |a(t)| dt$$

$$= \frac{1}{T} \int_0^T |(u, \varphi)| dt$$

$$= \frac{1}{T} \int_0^T \left| \int_{\Omega} u(x, t) \varphi(x) \right| dt$$

$$(2.64)$$

** The goal is to choose a basis with fastest decaying power spectrum, in time average **

Our goal can be summarized by the following loss function.

$$\mathcal{L}(\varphi) = \frac{\langle |(u,\varphi)|^2 \rangle}{\|\varphi\|^2}.$$
 (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,

$$\mathcal{L}(\varphi) = \frac{\frac{1}{T} \int_0^T \left| \int_{\Omega} u(x, t) \varphi(x) dx \right|^2 dt}{\int_0^L \varphi^2(x) dx} \equiv \langle a\varphi^2(t) \rangle$$
 (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:

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

where $a_k = f, \varphi_k$).

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

$$\max_{\left\{ \|\varphi_1\| = 1 \right\}} < |(u,\varphi)| >$$

$$\left\{ \|\varphi_1\| = 1 \right\}$$

$$\varphi \in \mathcal{H}$$

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

$$\max_{\left\{ \|\varphi_2\| = 1 \right\}} < |(u,\varphi)| >$$

$$\left\{ \|\varphi_2\| = 1 \right\}$$

$$\varphi \in \mathcal{H}$$

$$\varphi_2 \perp \varphi_1$$

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 the pictre reshere as vertir M: Pg X=[x,14,1-1/20]nxn

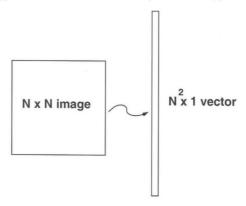
f = 6000/9 = 2008/31/20hemora variated Rogister

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

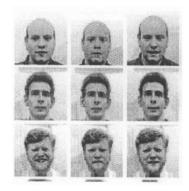
- Suppose Γ is an N^2 x1 vector, corresponding to an NxN face image I.
- The idea is to represent Γ (Φ = Γ mean face) into a low-dimensional space:

$$\hat{\Phi}-mean=w_1u_1+w_2u_2+\cdots w_Ku_K\,(K{<<}N^2)$$

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 Γ_i

Step 3: compute the average face vector Ψ :

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

Step 4: subtract the mean face:

$$\Phi_i = \Gamma_i - \Psi \qquad \checkmark \qquad \checkmark \qquad \checkmark$$

Step 5: compute the covariance matrix C:

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

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

CEXX

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 $A^T A (M \times M \text{ matrix})$

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 us_i and v_i ?

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

$$CAv_i = \mu_i Av_i$$
 or $Cu_i = \mu_i u_i$ where $u_i = Av_i$

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

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^TA (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 : $u_i = Av_i$

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

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



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

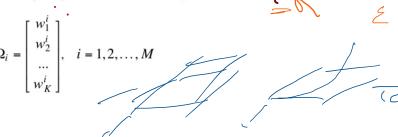
· g = reshape (v(:,i),p,g)

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

 $\int = \left[\begin{array}{c|c} U_1 & V_2 & V_3 & V_4 & V_5 & V_6 & V_$



Each normalized training face Φ_i is represented in this basis by a vector:



 $||a_i||_2^2 ||f||_2^2 - ||f||_2^2 ||f||_2^2$

 $j = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ comp

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

33 () - in 5 min. Signals analysis, Herronice o Historicelle Laveinte boasir set. B= {vi, wi, } (US. energy leverite ocsis set comes from RA) e Taylor Polynomials. (FLM) a catal tall the tax Hilletone (CY) = Cosin + the sin 2 + the sin 3 x + the sin 2 x + the sin 3 x + the sin 2 x + the sin 3 x + the sin 2 x + the sin 3 x + the sin 2 x + the sin 3 x + the sin 2 x + the sin 3 x + the sin 2 x + the sin 3 x + the sin 2 x + the sin 3 x + the sin 2 x + the sin 3 x + the sin 2 x + the sin 3 x + the sin 2 x + the sin 3 x + the sin 2 x + the sin 3 x + the sin 2 x + the sin Sinx = 1x - 43 + 51+ ...

Changing basis es a soot of coord-rot. 1 (y) = y? 6 2 (y) = y? -1 4 ((g) = s a cy - \$ \$2 (5) = Sin(25) - (3(51=5 in (34)

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

with proporties. 6.ft : you set geometry in E as on angle <(U,V) = <U,V) · vector some. set of objects "like rectors! set et 0 05 cets lite vectors.

That have a + and scalar multiplication - including communitive, associatione, and ident, morese, distribution presenter.

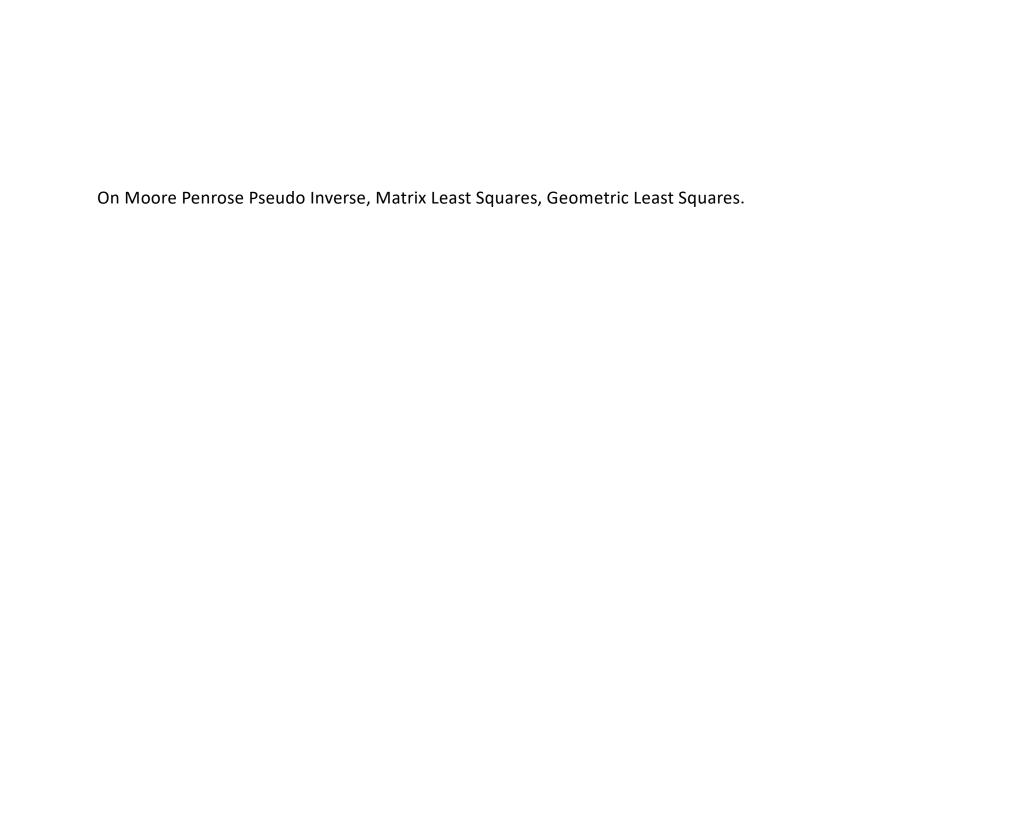
Communitatione, associatione, and ident, morese, distribution presenter. arrays of red numbers that are sx1. (L2(10,13) C C(50,17))

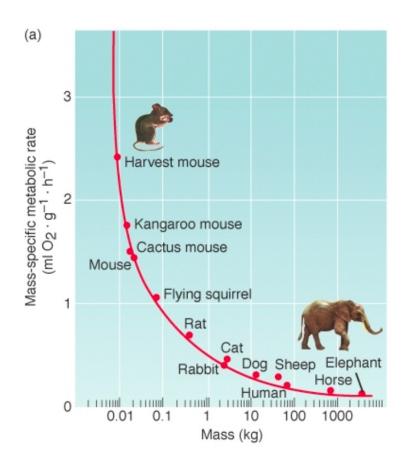
0B= \{-\name{\chi}\nam Former 31, Sin X, Sosx, Sizx, cors 24, --- 3 · B' - 3 - 1 /) , V , . 3 - 5 (x , x , x , x , - 5 Taylor. - Edit | 10 (23 (90° 1 2 1 At 1 At 1. flx1= \(\alpha_{\infty} \frac{\lambda_{\infty}}{\lambda_{\infty}} \), \(\alpha_{\infty} = \lambda_{\infty} \frac{\lambda_{\infty}}{\lambda_{\infty}} \), f(x) & L2 (1R)

On Compressed Sensing and on to Sparsity

7f = Sin x + 3Sin Sx + 9Sin 7x 7f = Sin x + 3Sin Sx + 9Sin 7x 7f = Sin 2x 7f = Sin 2x

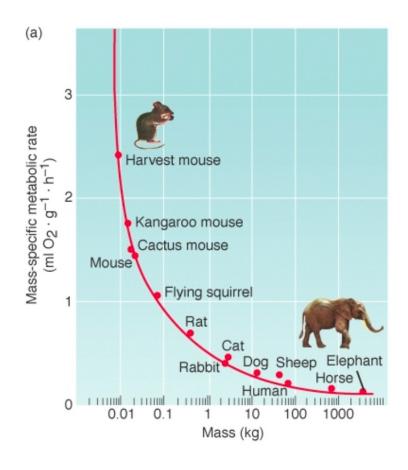
o a vector VEE is K-sperse if [V] has exactly K-nonzero values, K & didnill)





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

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



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

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

$$y_{0} = \beta_{0} + \beta_{1}x_{0} + \epsilon_{0}$$

$$y_{1} = \beta_{0} + \beta_{1}x_{1} + \epsilon_{1}$$

$$\vdots$$

$$y_{N-1} = \beta_{0} + \beta_{1}x_{N-1} + \epsilon_{N} - 1$$

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

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

$$E = \sum_{i=1}^{N} e_i$$

$$= \sum_{i=1}^{N} (f(x_i) - y_i)^2$$

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

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

$$E = \sum_{i=1}^{N} e_{i}$$

$$= \sum_{i=1}^{N} (f(x_{i}) - y_{i})^{2}$$

$$= \sum_{i=1}^{N} (\beta_{0} + \beta_{1}x_{i} - y_{i})^{2}.$$

$$rac{\partial E}{\partial eta_0} \;\; = \;\; 0 \quad ext{ and } \quad rac{\partial E}{\partial eta_1} = 0,$$

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

$$E = \sum_{i=1}^{N} e_{i}$$

$$= \sum_{i=1}^{N} (f(x_{i}) - y_{i})^{2}$$

$$= \sum_{i=1}^{N} (\beta_{0} + \beta_{1}x_{i} - y_{i})^{2}.$$

$$rac{\partial E}{\partial eta_0} \;\; = \;\; 0 \quad ext{ and } \quad rac{\partial E}{\partial eta_1} = 0,$$

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

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

and

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

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

From the above two equations we have:

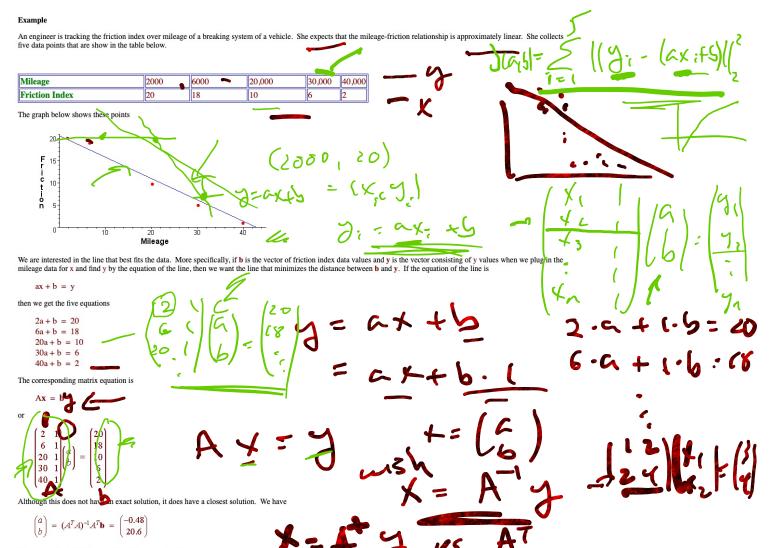
$$\sum_{i=1}^{N} (x_i y_i) = \sum_{i=1}^{N} (\beta_0 x_i) + \sum_{i=1}^{N} (\beta_1 x_i^2)$$
$$\sum_{i=1}^{N} (y_i) = \sum_{i=1}^{N} (\beta_0) + \sum_{i=1}^{N} (\beta_1 x_i)$$

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

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



We can conclude that the equation of the regression line is

$$y = -0.48x + 20.6$$

Least Squares

Definition and Derivations

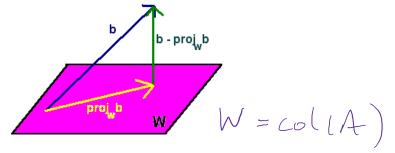
We have already spent much time finding solutions to

$$Ax = b$$

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

$$A\mathbf{x} = \text{proj}_{\mathbf{W}}\mathbf{b}$$

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



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

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

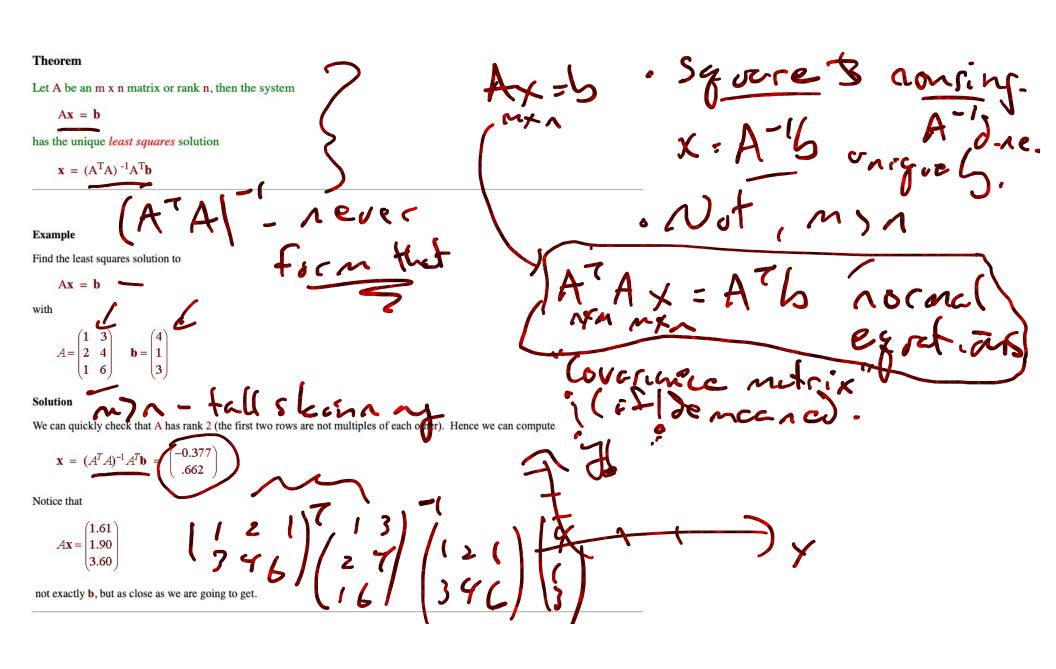
Now we need to show that $A^{T}A$ nonsingular so that we can solve for \mathbf{x} .

Lemma

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

$$\begin{vmatrix} a_{11} & a_{12} & ... & a_{1n} \\ a_{21} & a_{22} & ... & a_{1n} \\ a_{31} & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... \\ a_{mn} & ... & ... & ... & ... \\ a_{mn} & ... & ... \\ a_$$

 $= \sum_{x \in A} (x) + x = \sum_$ => (Ax-b) 1 every rector in CollA) solve "normal egns"



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.

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

$$a + b + c + d = 2.1$$
 $8a + 4b + 2c + d = 3.5$
 $27a + 9b + 3c + d = 4.2$
 $64a + 16b + 4c + d = 3.1$
 $125a + 25b + 5c + d = 4.4$
 $216a + 36b + 6c + d = 6.8$

The corresponding matrix equation is

	1	1	1	1)			(2.1)
	8	4	2	1	(a)		3.5
	27	9	3	1	b		4.2
	64	16	4	1	C	=	3.1
	125 216	25	5	1	(d)		4.4
		36	6	1)	` '		6.8

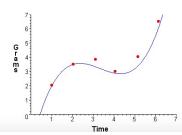
We can use the least squares equation to find the best solution

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

So that the best fitting cubic is

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

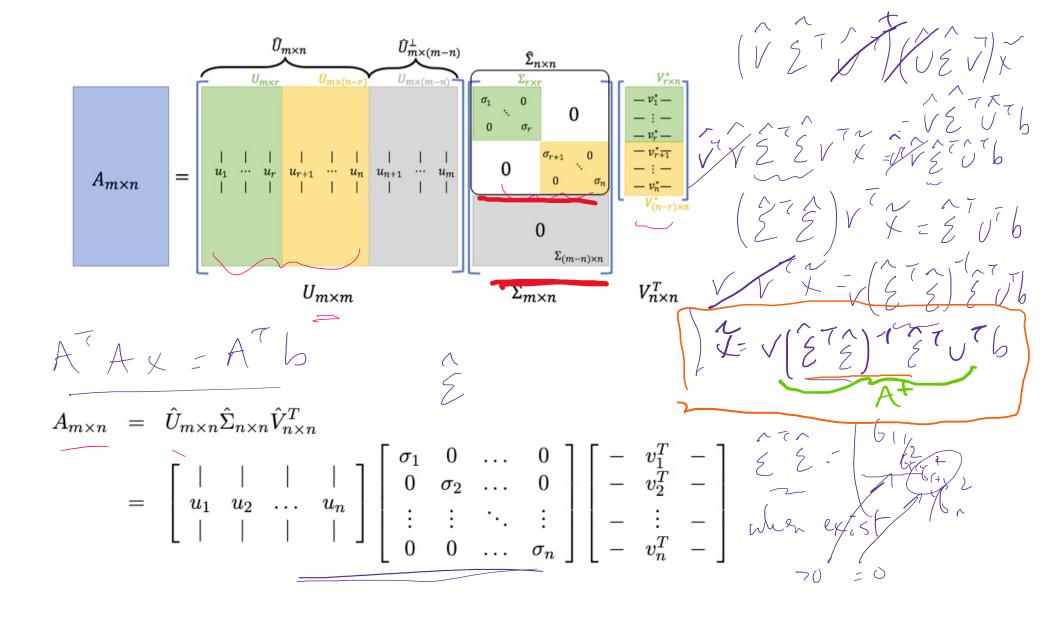
The graph is shown below



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

LS Slick



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

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

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

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

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

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

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

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

$$U_{n \times n} = \hat{U}_{n \times n} \hat{V}_{n \times n}^{T}$$

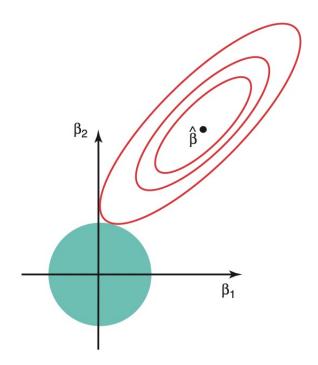
$$U_{n \times n} = \hat{U}_{n \times n} \hat{U}_{n \times n}^{T} \hat{U}_{n \times n} \hat{U}_{n \times n}^{T} \hat{U}$$

LS soln = solve normal eguations ATA X = AT6 When inverse exists Moore-X = (ATA) - (AT) Pseudo-Inverse = A+ b Ax=C o Interns of SVD? o and what it inverse docsit exist.

(ZT)

Tikhonov Regularization for inverse problems and preventing over fitting

- -Ridge regression
- -Lasso regression

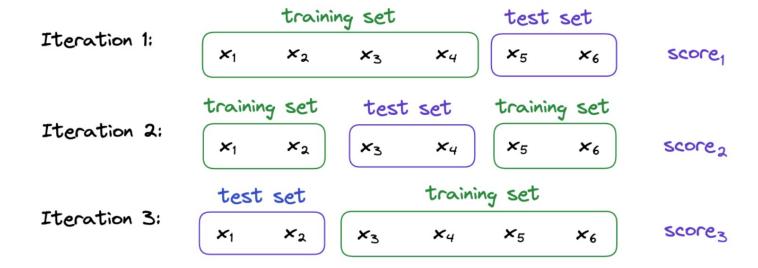


K-fold leave 1 out cross validation – A version cross validation



Then average scores of testing out of training sample

K-fold leave I out cross validation





N=7

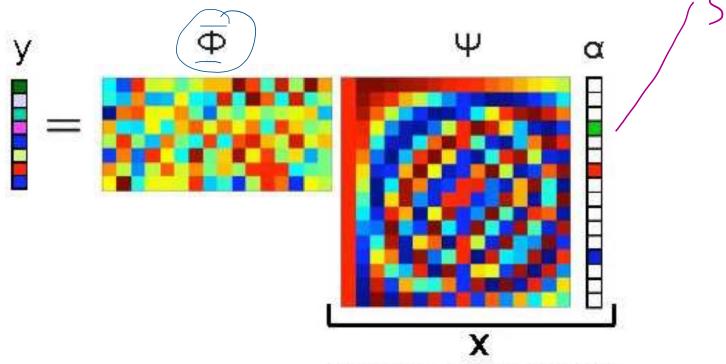
e Tourcet ion impose sparsity. (just tourcate after a large number.) 5.5x + 35,51x + 9 Siblax+9654X just 1,3,4,9 and 7 tell gov where to 8 put them. Truncation at brojection 3 Projection 3 Projection Sy motorix Mult. Fe H $X_i = [X]_i$ $2X \int_{3} = X_{3} = 0$

KE IR 4 50 75

Il Sllo = # of not values. Balls " ((sllp=1) USU, = 15,1+15,1+--+ (5,1) $||S||_2 = (5^2 + 5^2 + - - + 5^2)^2 = (54)^2 = (5)^2 + 5^2 + - - + 5^2 = (5)$ 11 SU = -> max (5:1 1(311) = (51 + 52 + ... + 57) 5 = (R)Wort to solve & w' v- no cm - But have to test all vols NP - complete.

the circ on of time 200d com Basis polsont -- Convex 5 = 05 m cin 1(\$1(, s) = 0.5 m cin 1(\$1)

 $\frac{1}{2} = \frac{1}{2} = \frac{1}$ $7 \in Col(T)$ = ligar, combo = 5, 4, + 5, 4, + --- + 5, 4, 4(laient const et) $X = \sum_{i=1}^{n} x_i + \sum_{i=1$ Random measurements can be used for signals sparse in any basis.



Signal x in sparse basis Ψ

Dogord enough it coherent to high protocold Day randa (PIA)

X= 45 o choose k=m spece for x \$ 4; 5 mx (×m× (Shonnon-Nygvist o Chorse & oll = Chorse & Chorse confect of X o if S is full (non Zeros.)
Then I need at least to trancite "Ide"

lh-sparse than le non zero volus - Cell spense if keem 3 < < 40,000 = 5, 4, +5, 4, + -.. 5 m 4 m sperse basis but like sull in another lasis.

that if x is not sparse with B= {4;}
it may be almost sparse with B 7 = 5,4,+5242+8545+ --+ +5m4a o and all or most $S_i \neq \delta$ but 3. (6

o X = 4 St s coefficients assume s is k esperse or at lsisael least almost k-sparse. · 115-511 < 8 for 570 smell \$ 5 is sparse projection arctif. Observe

for 5 = Solve. X= I 5 (w) 5 pacen

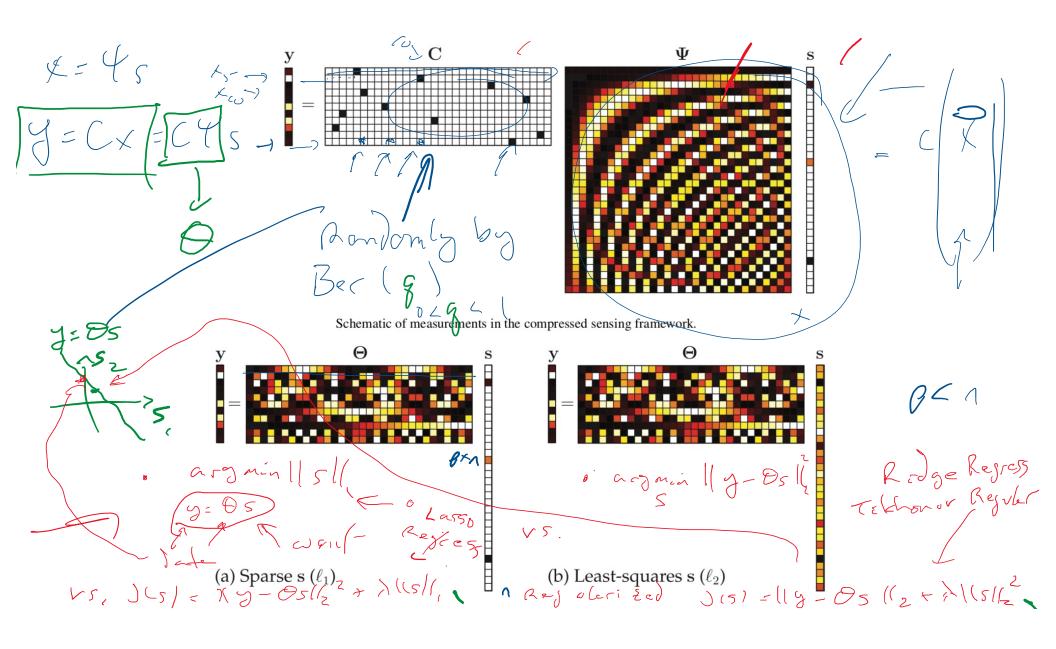
singular problem = no unique silution s inverse problems fre() say this ill-pored. o Hadamand define a well posed problem as on where solution exists else collillposed end cts ω , t. to the determinant of l=(1,2)Extreme what if l=(1,2)Existing t=(1,1) t=(

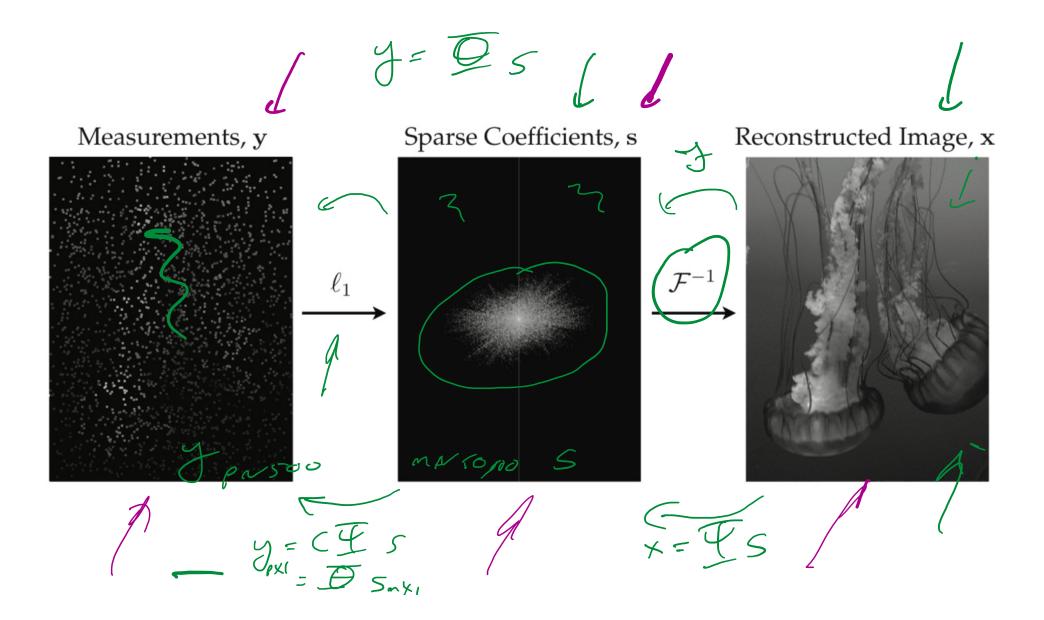
true value

-e

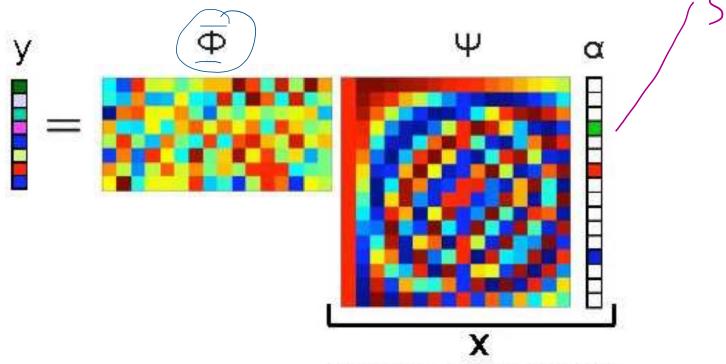
5=

0 $\left(\begin{array}{c} 1 \\ 5 \\ 5 \end{array}\right)$ BJ A tell gov t sis genre tell me . Danc 3 = organin +~. 3/7/10/3 subject to
y=05 =





Random measurements can be used for signals sparse in any basis.



Signal x in sparse basis Ψ

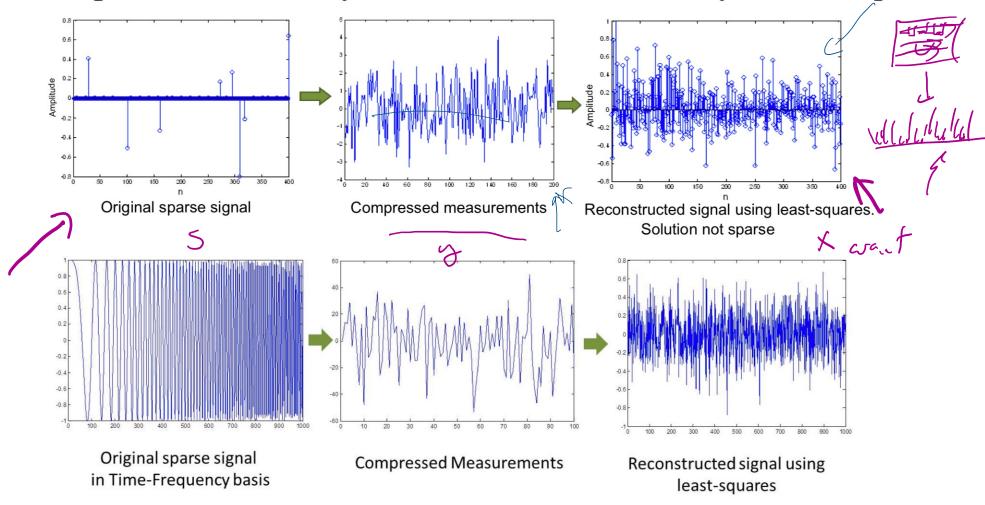
Dogord enough it coherent to high protocold Day randa (PIA)

$$y_k = <\phi_k, x>; k = 1, ..., M; \text{ with } M \ll N$$

- Need to solve an under determined system of equations $y = \Phi x$.
- Infinitely solutions for the system since $M \ll N$.
- A sparse solution x is recovered from y by solving the following inverse problem

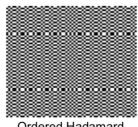
$$(P0): \min_{x} ||x||_{\ell_0} \ s.t. \ y = \Phi x$$

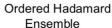
Example of the recovery of an under determined system of equations:



- Sparsity is what makes it possible to recover a signal from undersampled data.
- The number of measurements we need for successful reconstruction depends on the nature of the waveforms ϕ_k , and S

1. Incoherent Orthobasis

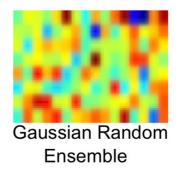






Scrambled Block
Hadamard Ensemble

2. Random waveforms ϕ_k



Incoherent Orthobasis Example

Example of incoherent basis: the "spike" basis (identity) and the Fourier basis.

Consider the case where the dictionary is the union of two orthobasis:

- *I*: the "spike" basis (identity).
- F: the Fourier basis (sinusoids).

$$\Phi = [I; F]$$

where I is a $N \times N$ matrix and F is a $N \times N$ matrix with

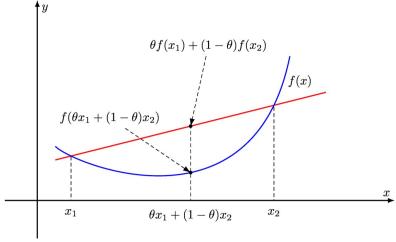
$$f_{m,\ell} = \frac{1}{\sqrt{(N)}} e^{j2\pi(m-1)(\ell-1)/N}$$

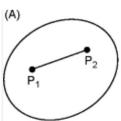
Convex Optimization – optimize a convex objective function over a convex domain

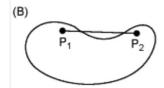
Convex function

$$f\left(tx_{1}+(1-t)x_{2}
ight)\leq tf\left(x_{1}
ight)+(1-t)f\left(x_{2}
ight)$$

Convex Domain







Compressive Sensing

[Candes, Romberg, Tao; Donoho]

- ullet Signal x is K-sparse in basis/dictionary Ψ
 - WLOG assume sparse in space domain

 $\Psi = I$

• Replace samples with few linear projections $y = \Phi x$

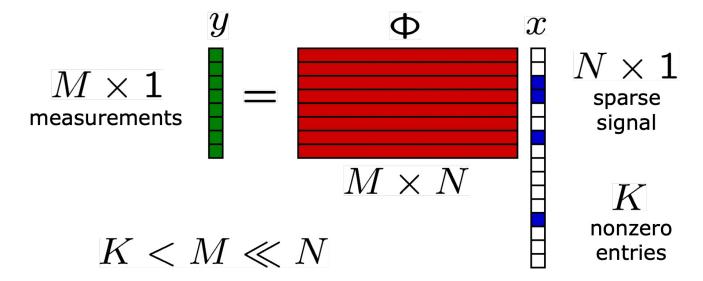
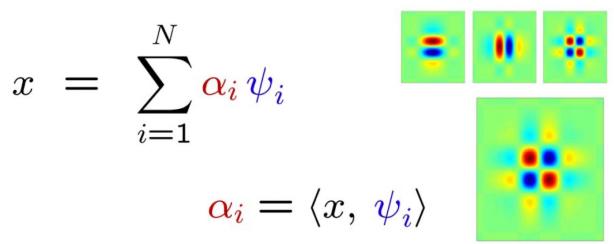
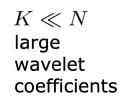


Image Representation



$$N$$
 pixels

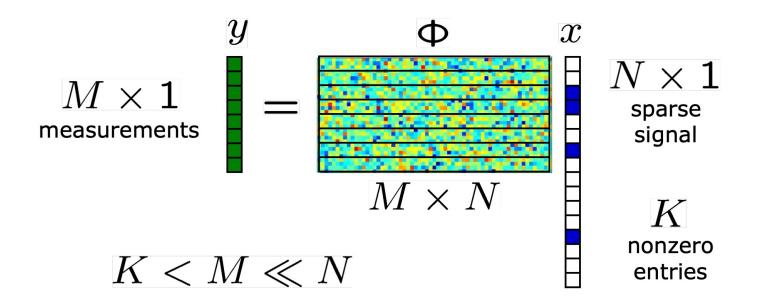






CS Signal Recovery

• Reconstruction/decoding: given $y = \Phi x$ (ill-posed inverse problem) find x



Reconstruction/decoding: given
 (ill-posed inverse problem) find

$$y = \Phi x$$

L₂ fast, wrong

$$\widehat{x} = \arg\min_{y = \Phi x} \|x\|_2$$

L₀ correct, slow

only M=K+1measurements required to perfectly reconstruct K-sparse signal

[Bresler; Wakin et al]

$$\widehat{x} = \arg\min_{y = \Phi x} \|x\|_0$$

$$\uparrow$$

$$number\ of$$

$$nonzero$$

$$entries$$

$$\begin{array}{ll} \text{given} & y = \Phi x \\ \text{find} & x \end{array}$$

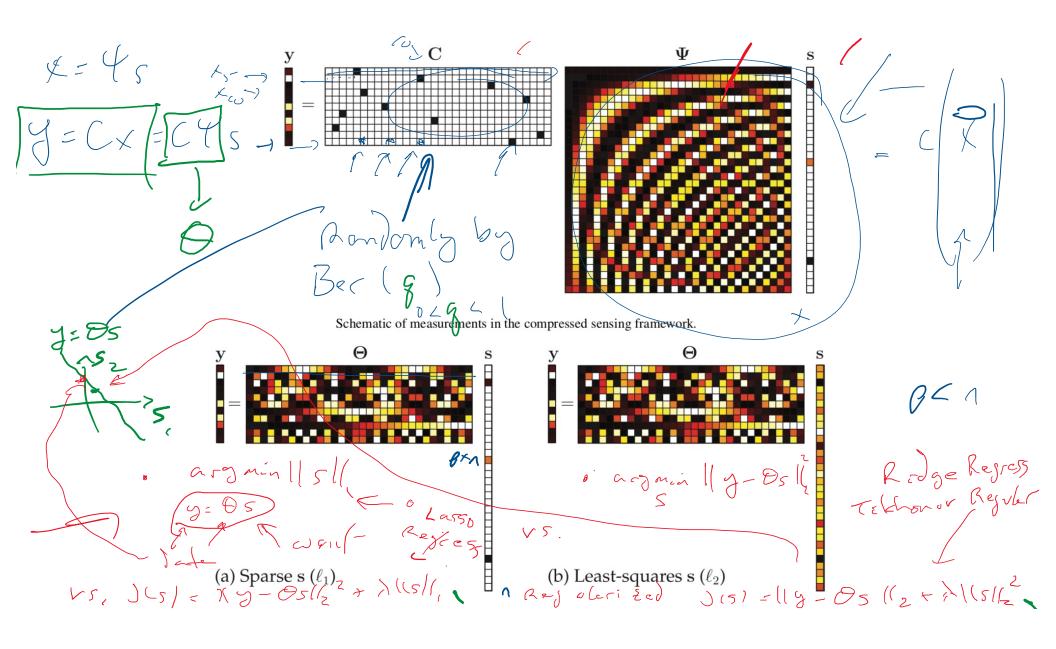
$$\widehat{x} = \arg\min_{y = \Phi x} \|x\|_2$$

$$\widehat{x} = \arg\min_{y = \Phi x} \|x\|_0$$

$$\widehat{x} = \arg\min_{y = \Phi x} \|x\|_1$$

[Candes et al, Donoho]

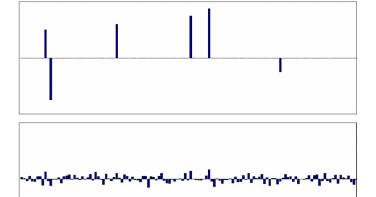
$$M \approx K \log N \ll N$$



CS Signal Recovery

- Reconstruction/decoding: given $y = \Phi x$ (ill-posed inverse problem) find x
- L₂ fast, wrong

$$\widehat{x} = \arg\min_{y = \Phi x} \|x\|_2$$

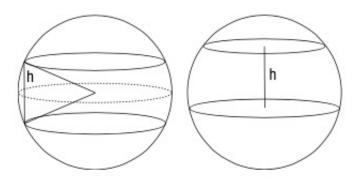


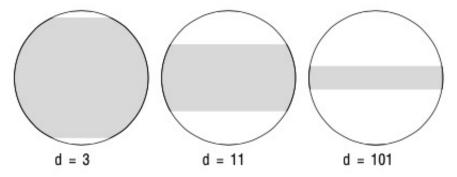
 \boldsymbol{x}

$$\hat{x} = (\Phi^T \Phi)^{-1} \Phi^T y$$

Another view on where all the random goes ... geometry of data in high dimensions highly important topic in Machine Learning, Data Science, and ROM

- -story one random gets soaked up by time.
- -story two its about random projection (Johnson-Lindenstrass) which is an application of concentration of measure





a strip of width h

(width of the) strip around the equator that contains 90% of the area

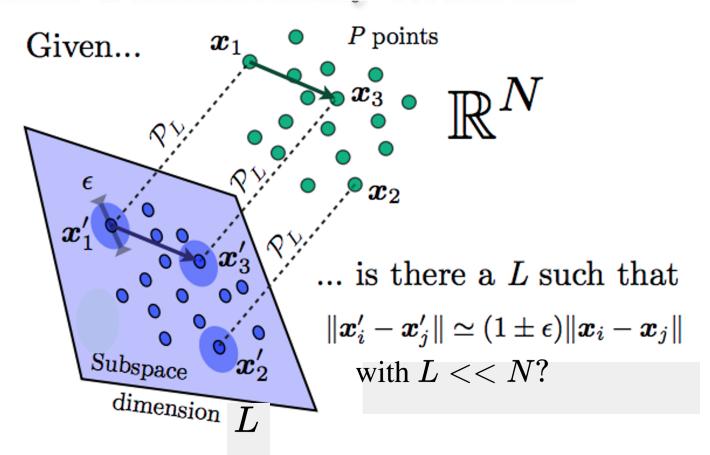
Several different important concepts follow from this simple idea

- -Markov inequalities, Chernoff inequality
- -for us in high dimensions, samples of data "look" almost like they live in a "flatter" space Johnson-Lindenstrass theorem

See Ortega for CoM/JL random projection interpretation of RC

Matous ek's book [Mat02]

Random Projection



Source: https://perso.uclouvain.be/laurent.jacques/uploads/Main/buffon-present-120314-MLG.pdf

Random Projection – on projection and approximate isometry

• The random projection method is based on the Johnson-Lindenstrauss lemma.

Theorem[Johnson-Lindenstrauss Lemma]:

For any $0<\epsilon<1$ and any integer M>1, let L be a positive integer such that $L\geq L_0$ with $L_0=\frac{C\ln M}{\epsilon^2}$,

where C is a suitable constant ($C \approx 8$ in practice,C = 2 is good enough). Then for any set X of M data points in \mathbb{R}^N ,

there exists a map $f:\mathbb{R}^N o\mathbb{R}^L$ such that

$$\text{ for all } x_1,x_2 \in X, (1-\epsilon)||x_1-x_2||^2 \leq ||f(x_1)-f(x_2)||^2 \leq (1+\epsilon)||x_1-x_2||^2.$$

[Johnson et. al, 1984]

Theorem[Random Projection]

For any $0 < \epsilon, \delta < \frac{1}{2}$ and positive integer N,

there exits a random matrix of B of size L imes N such that for $L \geq L_0$ with $L_0 = rac{C \ln(1/\delta)}{\epsilon^2}$

and for any unit-length vector $x \in R^N, \Pr\{|||Bx||^2 - 1| > \epsilon\} \leq \delta$

or
$$Pr\{||Bx||^2 - 1| > \epsilon\} \le e^{-CL\epsilon^2}$$

[Papadimitriou et. al., 1998]

CS and linprog

ait convex

CVX examples

To solve

 $\min_{x} \|x\|_{1}$ Ax = b Converged (Converged)

Use CVX code

```
cvx_begin
  variable x(n);
  variable t(n);
  minimize(sum(t));
  subject to
    x <= t;
    -x <= t;
    A*x == b;
cvx_end</pre>
```

1(X l(= 1 x (+1 x 1 + - + 1 X n)

$$\sum_{i=1}^{n} t_i$$

$$x \le t$$

$$-x \le t$$

$$Ax = b$$

$$x \le t$$

$$x \le t$$