

Numerical Analysis for Machine Learning

- 1) $A\vec{x} = \vec{b}$ solvability (4 fundamental spaces)
- 2) $A\vec{x} = \lambda\vec{x}$ generalization for any matrix
- 3) $A\vec{J} = \sigma\vec{w}$ SVD (singular value decomposition), $A = \sigma \vec{w} \vec{v}$
used for PCA, Latent Semantic Analysis
- 4) Minimization quadratic functionals
- 5) Factorization QR factorization

Matrix Vector Multiplication

$$\vec{c} = \underbrace{\begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix}}_A \underbrace{\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}}_x = \begin{bmatrix} x_1 + 4x_2 \\ 2x_1 + 5x_2 \\ 3x_1 + 6x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}x_1 + \begin{bmatrix} 4 \\ 5 \\ 6 \end{bmatrix}x_2$$

$\vec{c} \in C(A)$ column space of A (plane)
 \downarrow
 $\begin{bmatrix} 1 \\ 3 \end{bmatrix}, \begin{bmatrix} 4 \\ 6 \end{bmatrix}$ both linearly independent

$$A_{2 \times 3} \begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\vec{a}_1, \vec{a}_2, \vec{a}_3$$

$$\vec{a}_3 = -\vec{a}_1 + 2\vec{a}_2 \Rightarrow C(A_1) = C(A_2)$$

$$C(A_3) = \mathbb{R}^3$$

$$A_4 = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \\ 3 & 6 & 9 \end{bmatrix} \rightarrow \text{only 1 column vector is linearly independent}$$

$$C(A_4) = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \text{ (line)}$$

dim \Rightarrow rank of a matrix

dim (A_1) = 2
 dim (A_2) = 2
 dim (A_3) = 3
 dim (A_4) = 1

$$A\vec{x} = \vec{b}$$
 (mxn matrix)

$$\begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 10 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}x_1 + \begin{bmatrix} 4 \\ 5 \\ 6 \end{bmatrix}x_2 + \begin{bmatrix} 7 \\ 8 \\ 10 \end{bmatrix}x_3 = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

linear combination of the columns of A $\Rightarrow \vec{b}$ must be in the column space of A
 $\vec{b} \in C(A_3)$

Find $C(A)$ $\vec{a}_1, \vec{a}_2, \dots, \vec{a}_n$ columns of A

- put \vec{a}_2 in $C(A)$
- if $\vec{a}_2 = \alpha \vec{a}_1 \Rightarrow \vec{a}_2 \notin C(A) \Rightarrow$ reiterate
- else put \vec{a}_2 in $C(A)$

$$A_2 = CR$$

$$\begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix} \begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & 2 \end{bmatrix}$$

row-reduced echelon form (rref)
 useful when $\text{rank}(A) \leq n$

combination of \vec{a}_1 and \vec{a}_2 is zero in respective column of A_2

$$A_2 = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix} \quad A_2^T = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ \vec{a}_1 & \vec{a}_2 & \vec{a}_3 \end{bmatrix} \quad \dim(C(A)) = \dim(C(A^T)) = r \leq n$$

same rank

$$\vec{a}_3 = 2\vec{a}_1 - \vec{a}_2$$

Matrix - Matrix Multiplication

$$C = AB = \left[\begin{array}{c|c|c} 1 & & \\ \hline \vec{a}_1 & \cdots & \vec{a}_n \end{array} \right] \left[\begin{array}{c} \vec{b}_1 \\ \vdots \\ \vec{b}_n \end{array} \right] = \vec{a}_1 \vec{b}_1 + \cdots + \vec{a}_n \vec{b}_n$$

$\vec{a}_k, k \in [3, n]$ is a column vector
 $\vec{b}_j, j \in [1, n]$ is a row vector

$$C = A \sum B = \delta_{ij} \vec{a}_i \vec{b}_j + \cdots + \delta_{jn} \vec{a}_n \vec{b}_n$$

$$\text{rank 2 matrix}$$

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} 2 & 1 \\ 6 & 3 \end{bmatrix} + \begin{bmatrix} 4 & 6 \\ 8 & 12 \end{bmatrix} = \begin{bmatrix} 6 & 7 \\ 14 & 15 \end{bmatrix}$$

$$\vec{a}_1, \vec{a}_2, \vec{b}_1, \vec{b}_2$$

FACTORIZATIONS

$$1) A = LU \quad (PA = LU)$$

$$2) A = QR \quad \begin{cases} Q \text{ is an orthogonal matrix} \\ R \text{ is an upper triangular matrix} \end{cases}$$

$$3) S = ST \quad (\text{non-matrices})$$

$$S = Q \Lambda Q^T \quad \begin{cases} Q \text{ is an orthogonal matrix} \\ \Lambda \text{ is a diagonal matrix} \end{cases}$$

$$4) A = X \Lambda X^{-1} \quad \begin{cases} \Lambda \text{ is a diagonal matrix} \\ A \text{ is a sparse matrix} \end{cases}$$

$$5) A = U \Sigma V^T \quad \begin{cases} U \text{ and } V \text{ are orthogonal matrices} \\ \Sigma \text{ is a pseudo-diagonal matrix} \end{cases}$$

Orthogonal Matrix

$$QQ^T = Q^T Q = I$$

If also implies that $\|Q\vec{x}\|^2 = \|Q^T Q\vec{x}\|^2 = \|\vec{x}\|^2$

$$\begin{aligned} \|Q\vec{x}\|^2 &= (\vec{x}^T) Q^T Q \vec{x} \\ &= \vec{x}^T Q^T Q \vec{x} \\ &= \vec{x}^T \vec{x} = \|\vec{x}\|^2 \end{aligned}$$

Rotation Matrix is orthogonal

$$\begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}$$

Reflection Matrix is orthogonal

$$\begin{aligned} \vec{v}^T \vec{v} &= 1 \\ \vec{v}^T \vec{v} &= \vec{v}^T (\vec{v} - 2(\vec{v}^T \vec{v})\vec{v}) \\ &= (\vec{v} - 2\vec{v}^T \vec{v})^T \vec{v} \\ &= (\vec{I} - 2\vec{v}^T \vec{v})\vec{v} \\ &= R\vec{v} \end{aligned}$$

Reflection Matrix

R

$$S = (Q \Lambda) Q^T = \vec{Q} \vec{Q}^T = \lambda_1 \vec{q}_1 \vec{q}_1^T + \dots + \lambda_n \vec{q}_n \vec{q}_n^T \quad (\text{eigenvalues}) \Rightarrow S \vec{x} = \lambda_1 \vec{q}_1^T \vec{x}$$

Spectral decomposition

$$\boxed{A \in \mathbb{R}^{m \times n}, \quad U \in \mathbb{R}^{m \times m}, \quad V \in \mathbb{R}^{n \times n}, \quad \Sigma \in \mathbb{R}^{m \times n}}$$

Supporting fact $m > n$

$$\sum_i = m \quad \left(\begin{array}{cccccc} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{array} \right)$$

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

$$A \in \mathbb{R}^{m \times n}$$

$$\text{rank}(A) = r$$

$$1) C(A) \subset \mathbb{R}^m, \dim(C(A)) = r$$

$$2) C(A^T) \subset \mathbb{R}^n, \dim(C(A^T)) = r$$

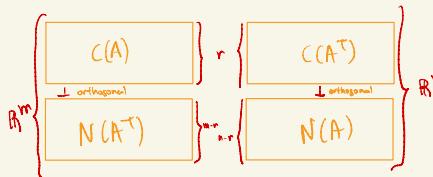
$$3) N(A) \subset \mathbb{R}^n, \dim(N(A)) = ?$$

$$4) N(A^T) \subset \mathbb{R}^m, \dim(N(A^T)) = ?$$

$$A\vec{x} = \vec{0} \quad N(A) = \{ \vec{x} \in \mathbb{R}^n : A\vec{x} = \vec{0} \}$$

$$\text{Ker}(A) = N(A) \quad \text{Null space/Kernel}$$

$$A^T \vec{x} = \vec{0} \quad N(A^T) = \{ \vec{x} \in \mathbb{R}^m : A^T \vec{x} = \vec{0} \}$$



$$1) \vec{x} = \vec{0} \in N(A)$$

$$2) \text{if } \vec{x} \neq \vec{0} \in N(A) \Rightarrow A(\vec{x} + \vec{y}) = \vec{0}$$

$$3) \text{if } \vec{x} \in N(A) \Rightarrow \alpha \vec{x}, \alpha \in \mathbb{R} \Rightarrow A(\alpha \vec{x}) = \vec{0}$$

$$\text{Ex: } A = \begin{bmatrix} 1 & 4 & 3 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \Rightarrow \begin{aligned} x_1 + 4x_2 + 3x_3 &= 0 & x_1 - 4x_2 - 7x_3 &= 0 \\ 2x_1 + 5x_2 + 8x_3 &= 0 & 2x_1 + 5x_2 + 8x_3 &= 0 \\ 3x_1 + 6x_2 + 9x_3 &= 0 & 3x_1 + 6x_2 + 9x_3 &= 0 \end{aligned}$$

$$\text{rank } A = 2 \Rightarrow 1 \text{ degree of freedom (dof)} \quad \begin{aligned} -3x_2 - 6x_3 &= 0 \\ -6x_2 - 12x_3 &= 0 \end{aligned}$$

$$m-r = 3-2 = 1$$

$$\text{but } \Rightarrow A = \begin{bmatrix} 1 & 4 & 3 \\ 2 & 5 & 8 \\ 3 & 6 & 10 \end{bmatrix} \Rightarrow \begin{aligned} \text{If } A \text{ is full-rank,} \\ \text{rank } A = 3, 0 \text{ dof} \quad N(A) = \{ \vec{0} \} \end{aligned}$$

$$\begin{aligned} A \in \mathbb{R}^{m \times n} \quad \text{rank}(A) = r < n & \quad A = CR \\ A = [A_1 \ A_2], \quad A_1 \in \mathbb{R}^{m \times r}, \quad A_2 \in \mathbb{R}^{m \times (n-r)} & \\ A = [A_1 \ A_2 \ B], \quad B \in \mathbb{R}^{m \times (n-r)} & \quad A = \begin{bmatrix} 1 & 4 & 3 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{bmatrix} \sim \begin{bmatrix} 1 & 4 & 3 \\ 0 & 1 & 2 \\ 0 & 0 & 1 \end{bmatrix} \sim \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \sim A_1 \end{aligned}$$

$$K = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}, \quad K \in \mathbb{R}^{n \times (n-r)}$$

$$AK = \begin{bmatrix} A_1 \ A_2 \ B \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} = \vec{0}$$

$$-A_2 B + A_2 B = \vec{0} \quad (\text{mv } (n-r))$$

$$AK = \vec{0} \Rightarrow A\vec{x} = \vec{0} \rightsquigarrow \vec{x} \in N(A)$$

$$K\vec{w} = \vec{0} \Rightarrow \vec{w} = \vec{0}$$

$$\Downarrow$$

$$\begin{bmatrix} -8 \\ 1 \\ 1 \end{bmatrix} \vec{w} = \vec{0} \Rightarrow \begin{bmatrix} -8\vec{w}_1 \\ \vec{w}_2 \\ \vec{w}_3 \end{bmatrix} = \begin{bmatrix} \vec{0} \\ \vec{0} \\ \vec{0} \end{bmatrix} \Downarrow \vec{w} = \vec{0}$$

$$A\vec{x} = \vec{0} \quad \text{each } \vec{x} \text{ that satisfy (6) must be a linear combination of the columns of } K$$

$$(6) \quad \vec{A}_1 \vec{x} = \vec{0} \in \mathbb{R}^m \Rightarrow \vec{x} = \vec{0} \in \mathbb{R}^n$$

$$\vec{A}_2 \vec{x} = \vec{0} \in \mathbb{R}^m$$

$$= [A_1 \ A_2 \ B] \begin{bmatrix} \vec{w} \\ \vec{v} \end{bmatrix} = [A_1 \vec{w}_1 + A_2 \vec{w}_2 \vec{v}] = A_1 \begin{bmatrix} \vec{w}_1 + B\vec{w}_2 \vec{v} \end{bmatrix} = \vec{0}$$

$$\vec{w}_1 + B\vec{w}_2 \vec{v} = \vec{0} \Rightarrow \vec{w}_1 = -B\vec{w}_2 \vec{v}$$

$$\vec{w} = \begin{bmatrix} -B\vec{w}_2 \vec{v} \\ \vec{w}_2 \\ \vec{v} \end{bmatrix} = \begin{bmatrix} -B\vec{w}_2 \vec{v} \\ \vec{w}_2 \\ \vec{v} \end{bmatrix} = K\vec{v}$$

EIGENVALUES & VECTORS

- A general square matrix $n \times n$

+ symmetry
+ positive definite (spd)

$$A \vec{x}_i = \lambda_i \vec{x}_i \quad i = 1, \dots, n$$

$$X^{-1} A X = \Lambda \quad \text{diagonal matrix of eigenvalues}$$

$$A^2 \vec{v} = \lambda^2 \vec{v} \Rightarrow A(\lambda \vec{v}) = A(\lambda_i \vec{v}_i) = \lambda_i (\lambda_i \vec{v}_i) = \lambda_i^2 \vec{v}_i$$

$$e^{A^2} = e^{\lambda_i^2} \quad \frac{d}{dt} = A^2$$

$$\text{Power method: } \vec{w}^{(n)} = A \vec{w}^{(n-1)} = A^n \vec{w}^{(0)}$$

If A is fullrank then any vector $\vec{v} \in \mathbb{R}^n$ can be written as linear combination of the eigenvectors (\vec{x}_i)

Similar Matrices

$$A \sim B \Leftrightarrow B = M^{-1} A M \quad M \text{ is invertible}$$

$$M^{-1} A M \vec{v} = \lambda \vec{v} \Rightarrow A M \vec{v} = \lambda M \vec{v} \Rightarrow A \vec{v} = \lambda \vec{v} \quad \text{same eigenvectors}$$

Power Method (PM) on A (non-eigenvalues)

Inverse PM on A^{-1} (non-eigenvalues)

PM with shift $(A - \alpha I)^{-1} \quad \alpha \in \mathbb{R}$ (start eigenvalues to ∞)

Deflation method: $\vec{v}_k \vec{v}_k^T$

$$\begin{bmatrix} 0 & 0 \\ 0 & \alpha_1 \end{bmatrix}$$

QR Factorization

$A \in \mathbb{R}^{m \times n}$ $m \geq n$ independent columns
 $A = QR$ $\{Q\}$ has orthogonal columns
 $(m \times m)$ R is upper triangular

GR Iteration

$$\begin{aligned} A &= A^{(0)} Q^{(0)} R^{(0)} \\ A^{(0)} &= Q^{(0)} R^{(0)} \\ A^{(1)} &= Q^{(1)} R^{(1)} \\ A^{(2)} &= \dots \quad A^{(n)} = \text{upper triangular} \quad \text{diagonal values are the eigenvalues} \end{aligned}$$

How can we compute Q ? Gram-Schmidt orthogonalization procedure:

$$\begin{aligned} A &\rightarrow \dots \rightarrow Q \quad \vec{v}_i^T \vec{v}_j = 1 \text{ if } i=j \\ A &\sim \begin{bmatrix} 1 & & & \\ 0 & 1 & & \\ 0 & 0 & 1 & \\ 0 & 0 & 0 & 0 \end{bmatrix} \\ \vec{v}_1 &= \vec{a}_1 \\ \vec{v}_2 &= \vec{a}_2 - \text{proj}_{\vec{v}_1} \vec{a}_2 \\ \vec{v}_3 &= \vec{a}_3 - \text{proj}_{\vec{v}_1} \vec{a}_3 - \text{proj}_{\vec{v}_2} \vec{a}_3 \\ \vdots & \vdots \\ \vec{v}_n &= \vec{a}_n - \sum_{i=1}^{n-1} \text{proj}_{\vec{v}_i} \vec{a}_n \end{aligned}$$

$$\text{Ex: } A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 0 \\ 3 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Symmetric matrix $S \in \mathbb{R}^{n \times n}$ $S = S^T$

$S = Q \Lambda Q^T$ Q orthogonal matrix

$S \vec{v} = \lambda \vec{v}$, $S \vec{v} = \alpha \vec{v}$ one real

$$\begin{aligned} \vec{v} \in C(S) &\Rightarrow \vec{v} \in N(S) \quad \text{implies} \quad \vec{v} \perp \vec{v} \\ &= C(S^T) \quad \vec{v} \in N(S) \quad \text{implies} \quad \vec{v} \perp S \vec{v} \\ &= C(S) \quad \vec{v} \in N(S) \quad \text{implies} \quad \vec{v} \perp S^T \vec{v} \\ &\text{defining } S^T = S \quad \vec{v} \in N(S) \quad \text{implies} \quad \vec{v} \perp S \vec{v} \\ &\vec{v} \in C(S) \quad \vec{v} \in N(S) \quad \text{implies} \quad \vec{v} \perp S \vec{v} \end{aligned}$$

$$\vec{v}^T S \vec{v} = \lambda \vec{v}^T \vec{v} \quad \text{Rayleigh's quotient}$$

$$(a+i b)(a-i b) = (a^2 + b^2)$$

(Symmetric) Positive-definite (spd)

$$(i) \lambda_i > 0 \quad i = 1, \dots, n$$

$$(ii) \vec{v}^T S \vec{v} > 0, \forall \vec{v} \in \mathbb{R}^n, \vec{v}^T S \vec{v} = \vec{v}^T \vec{v} = \vec{v}^2 \quad \vec{v} \neq 0$$

(iii) Leading determinants are positive

$$(iv) S = B^T B \quad (\text{Cholesky factorization})$$

upper triangular

v) All pivot elements are positive

$$\lambda > 0$$

$$S \vec{v} = \lambda \vec{v}$$

$$\vec{v}^T S \vec{v} = \lambda \vec{v}^T \vec{v} = \lambda \| \vec{v} \|^2 > 0$$

$$(2) \vec{v} = (c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_n \vec{v}_n)$$

$$(c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_n \vec{v}_n)^T S (c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_n \vec{v}_n)$$

$$(i) c_1 \vec{v}_1^T S \vec{v}_1 + c_2 \vec{v}_2^T S \vec{v}_2 + \dots + c_n \vec{v}_n^T S \vec{v}_n = c_1^2 \lambda_1 \| \vec{v}_1 \|^2 > 0$$

$$(ii) c_1 c_2 \vec{v}_1^T S \vec{v}_2 = (c_1 \lambda_1) c_2 \vec{v}_1^T \vec{v}_2 = 0$$

$$(3) S = B^T B$$

$$\begin{aligned} J(B^T) \vec{v} &= (J^T B^T)(B \vec{v}) \\ &= (B^T)^T (B \vec{v}) = \| B \vec{v} \|^2 > 0 \end{aligned}$$

SVD (Single Value Decomposition)

Applications

- **Least-squares approximation**
finding the regression line representing the trend of data points
need to find slope intercept of y-axis
 - Via SVD, we can introduce a pseudo-inverse matrix
 - **Low-rank approximation**
(Eigen-Value)
based on the usage of the SVD

✓ works for every matrix (always possible)

$A \in \mathbb{R}^{m \times n}$
 $\text{rank}(A) = r < n$ $\Rightarrow A = U \Sigma V^T$, where
 1) U has orthogonal columns
 2) V has orthogonal columns
 3) Σ is at most diagonal matrix

(rank(A)) $\leq r$

$\text{rank}(A) = r$, then
 $A = U \Sigma V^T$
 $A \vec{x}_j = \sigma_j \vec{u}_j$

$\text{rank}(A) = r$, then
 $\begin{cases} \vec{u}_1 \perp \vec{u}_2 \perp \dots \perp \vec{u}_r \\ \vec{u}_1, \vec{u}_2, \dots, \vec{u}_r \neq 0 \end{cases}$

$A = U \Sigma V^T$
 $A \vec{v}_i = \sigma_i \vec{v}_i$
 $A \vec{v}_i = \vec{0}$
 $A \vec{v}_i = 0^\top$

$\vec{v}_1, \vec{v}_2, \dots, \vec{v}_n \in \text{C}(A)$, $\forall i: i \in r$
 $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_n \in N(A)$, $\forall i: i \in n-r$

ECONOMY SVD

$m = \begin{matrix} \text{compact/reduce} \\ \text{representation} \end{matrix}$

$(m \times r) (r \times n)$

$(m \times r) (r \times r) (r \times n)$

Proof of existence of SVD

$$\begin{aligned}
 & \text{symmetric} \quad \text{rank}(A) = r \\
 & \text{positive definite} \quad i) (A^T A) = A^T A \\
 & \quad ii) x^T (A^T A) x = (x^T A)^T A x = \\
 & \quad \quad \quad = (A x)^T A x = \|Ax\|^2 \geq 0 \\
 & A^T A = V \Lambda V^T = \\
 & \quad \quad \quad = A \tilde{x} = \tilde{\delta} \Rightarrow \tilde{x} \in N(A^T A) \\
 & \quad \quad \quad A^T (A \tilde{x}) = A^T \tilde{\delta} \Rightarrow \tilde{x} \in N(A^T A) \\
 & \quad \quad \quad \text{consequence} \\
 & - A^T A \tilde{x} = \tilde{\delta} \Rightarrow \tilde{x} \in N(A^T A) \\
 & \quad \quad \quad \tilde{x}^T A^T A \tilde{x} = \|\tilde{x}\|^2 = 0, \quad x \in N(A^T A) \\
 & \quad \quad \quad \text{rank}(A^T A) = \text{rank}(A)
 \end{aligned}$$

$$\text{(iii) } \vec{w}_i \text{ are eigenvectors of } A^T A \text{ with eigenvalues } \sigma_i^2$$

$$A^T \vec{w}_i = A^T \left(\frac{\sigma_i \vec{v}_i}{\sigma_i} \right) = \frac{A^T A \vec{v}_i}{\sigma_i} =$$

$$= \frac{A \sigma_i^2 \vec{v}_i}{\sigma_i} = \sigma_i^2 \left(\frac{A \vec{v}_i}{\sigma_i} \right) = \sigma_i^2 \vec{w}_i \quad \xrightarrow{\text{by defn}}$$

$$\text{Ex. 1}$$

$$A = \begin{bmatrix} 4 & 4 \\ -3 & 3 \end{bmatrix}$$

(rank(A) = 2)

$$X = A^{-T} A = \begin{bmatrix} 25 & 7 \\ 7 & 25 \end{bmatrix}$$

$$\text{eigs}(X) = \begin{bmatrix} \lambda_1 = 43 & \vec{v}_1 = \begin{bmatrix} \frac{\sqrt{15}}{2} \\ \frac{1}{2} \end{bmatrix} \\ \lambda_2 = 32 & \vec{v}_2 = \begin{bmatrix} \frac{1}{2} \\ \frac{\sqrt{15}}{2} \end{bmatrix} \end{bmatrix}$$

$$Y = AA^T = \begin{bmatrix} 32 & 0 \\ 0 & 49 \end{bmatrix}$$

$$\vec{u}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \vec{u}_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$U = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad V = \begin{bmatrix} -\frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \sqrt{18} & 0 \\ 0 & \sqrt{32} \end{bmatrix}$$

$$\text{Ex-2}$$

$$A = \begin{bmatrix} 4 & 3 \\ 8 & 6 \end{bmatrix} \quad X = A^{-1} = \begin{bmatrix} 60 & -60 \\ 60 & 45 \end{bmatrix}$$

$$\text{rank}(A) = 1$$

$$\text{rank}(X) = ?$$

$$\begin{cases} \lambda_1 = 125 \\ \lambda_2 = 0 \end{cases} \Rightarrow \begin{cases} v_1 = \begin{bmatrix} -4/3 \\ 2/3 \end{bmatrix} \\ v_2 = ? \end{cases} \Rightarrow v_2 = \begin{bmatrix} 3/5 \\ -3/5 \end{bmatrix}$$

$$A^{-1} = \frac{1}{\det(A)} \begin{bmatrix} 60 & -60 \\ 60 & 45 \end{bmatrix}$$

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

$$A \in \mathbb{R}^{n \times n}$$

$A = A^T$

$$A = QS$$

↑ polar decomposition

? symmetric positive

$$A = U\Sigma V^T = \underbrace{(UV^T)}_Q \underbrace{\Sigma}_{S} \underbrace{(V\Sigma V^T)}_S$$

- $$\text{If } A \text{ is orthogonal from } \sigma_1 = 1 \quad A^T A = I$$

(i) All eigenvalues of A are unitary and ≤ 1

$A^T A = I \Rightarrow \lambda_i^2 = 1 \quad \forall i$

$\Rightarrow \|A\|_F \leq \sqrt{\|\lambda\|_F^2}$

$\|A\|_F^2 = \sqrt{\lambda_1^2 + \lambda_2^2 + \dots + \lambda_n^2} = \sqrt{n} \leq \sqrt{1^2 + 1^2 + \dots + 1^2} = \sqrt{n}$

A^T { which is better? AA^T (n²n) \underline{ATA} (n²n) method a. samples

$\|I - B\|_2 = \text{rank}(B) = K$
 $N(B) \subset \mathbb{R}^d \quad \dim(N(B)) = d - K$

$V_{K+1} = \begin{bmatrix} v_{K+1}^1 & \dots & v_{K+1}^d \end{bmatrix}$ First $K+1$ columns of V

$C(V_{K+1}) \subset \mathbb{R}^d \quad \dim(C(V_{K+1})) = K + 1$
 $\dim(N(B)) + \dim(C(V_{K+1})) = d - K + K + 1 = d + 1$

since b is
 $N(B)$ and
 $C(V_{K+1})$
 $\Rightarrow b \in C(V_{K+1})$

$\tilde{w} \in N(B) \cap C(V_{K+1}) \quad \|\tilde{w}\|_2 = 1$

$\tilde{w} = \sum_{i=1}^d c_i v_{i+1}$

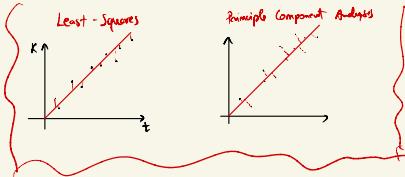
$= V_{K+1} \tilde{c}$

$\|A - B\tilde{w}\|_2 = \max_{1 \leq i \leq d} \|A_i - B_i\tilde{w}\|_2 \geq \|A_i - B_i\tilde{w}\|_2 \geq \|(A_i - B_i)\tilde{w}\|_2$

$B_i \tilde{w} = \vec{0} \Rightarrow \|A_i - B_i\tilde{w}\|_2 \geq \|A_i \tilde{w}\|_2 = \tilde{w}^\top A_i^\top A_i \tilde{w} = \tilde{w}^\top V_i^\top V_i \tilde{w}$
 $= \tilde{c}^\top V_{K+1}^\top V_{K+1} \tilde{c} \stackrel{?}{=} \tilde{c}^\top V_{K+1} V_{K+1} \tilde{c} \geq \|\tilde{c}\|_2^2 = \|\tilde{w}\|_2^2$

PCA

- Project the dataset in a new space where
 - variance is maximized
 - covariance is minimized
 - dimensionality reduction of the dataset



i) Center the matrix

\bar{A} is mean centered with respect to columns

$$H = I_n - \frac{1}{n} \vec{1}_n \vec{1}_n^T$$

(centering matrix)

$$\bar{A} = HA$$

ii) Build the covariance matrix

$$S = \frac{\bar{A}^T \bar{A}}{n-1} \Rightarrow S \text{ is spd, thus } SV = VD \text{ holds}$$

$$\delta V = VD \Rightarrow S = VD\sqrt{V^T}, D = \sqrt{V^T S V} \quad \text{vectors in } V \text{ are the principal components}$$

$$\bar{A} = V\Sigma V^T$$

$$S = \frac{1}{n-1} \bar{A}^T \bar{A} = \frac{1}{n-1} V\Sigma V^T V\Sigma V^T = \frac{1}{n-1} V\Sigma^2 V^T = \frac{1}{n-1} V\Sigma^2 V^T$$

$$D = \frac{1}{n-1} \Sigma^2 \Rightarrow \lambda_k = \frac{\sigma_k^2}{n-1}$$

Image Compression

$$A = A_{true} + \gamma A_{noise}$$

Gaussian noise with zero mean, unit variance

τ threshold $\sigma_e > \tau$

When γ is known:

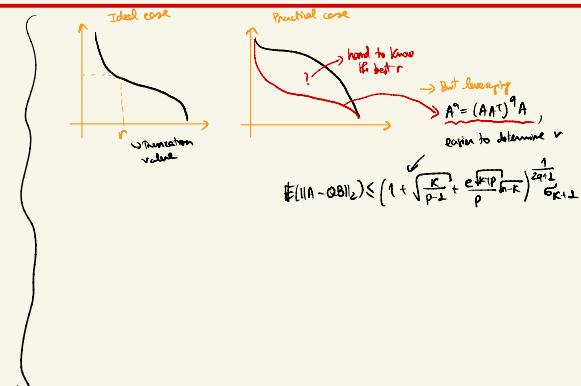
$$\text{i) } A \in \mathbb{R}^{n \times n} \quad \tau = \frac{4}{\sqrt{3}} \sqrt{n} \gamma$$

$$\text{ii) } A \in \mathbb{R}^{m \times n} \quad \begin{cases} \text{(a) } n \ll m, \beta = \frac{m}{n} \\ \text{(b) } m \ll n, \beta = \frac{n}{m} \end{cases} \quad \lambda(\beta) = \left(\frac{2(\beta+1)}{(\beta-1)} + \frac{8\beta}{(\beta-1)(\beta^2-14\beta+12)} \right)^{1/2}$$

$$\tau = \lambda(\beta) \sqrt{n} \gamma$$

when β is not known:

$$\tau = w(B) \sigma_{\text{med}} \quad w(p) = \lambda(\beta) / p_B, \text{ when } p_B = \int_{(-p)^2}^{w_B} \frac{((1+t\beta^2)^2 - t)(t - (1-t\beta^2)^2)}{2t^2} dt$$



Least Squares Approximation (LS)

$$X \in \mathbb{R}^{n \times p} \quad n \Rightarrow \# \text{samples}$$

rank(X) = p $p \Rightarrow \# \text{features}$

(features are independent)

$$y \in \mathbb{R}^n \quad \Rightarrow \text{"labels" of each sample}$$

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} \quad \vec{x}_j \Rightarrow j^{\text{th}} \text{ column of matrix } (feature \text{ column})$$

\vec{x} \rightsquigarrow \vec{y} \rightsquigarrow weight vector
 $y = \vec{x}^T \cdot \vec{w}$ instead $\vec{y} = X\vec{w}$ $\Rightarrow \vec{y} \in C(X)$
? $\vec{y} \in C(X)$

Approaches to minimize $\|\vec{r}(\vec{\omega})\|_2^2$

1) Geometric Approach

Rem. $y_i \in \{-1, +1\} \quad i=1, \dots, m$ $\vec{w} = (X^T X)^{-1} X^T \vec{y}$
 $\vec{y} = X \vec{w}$ $\vec{q} = \text{sign}(\vec{y})$ \rightarrow classification problems tend to
 w.r.t. require post-processing of the
 results

2) Optimization

$$\begin{aligned}
 \vec{w}^* &= \underset{\vec{w}}{\text{argmin}} \| \vec{q}^T (\vec{w}) \|^2_2 = \quad \vec{q}^T (\vec{w}) = \vec{q}^* - X\vec{w} \\
 &= \underset{\vec{w}}{\text{argmin}} (\vec{q}^* - X\vec{w})^T (\vec{q}^* - X\vec{w}) = \boxed{X \text{ is full rank}} \\
 &= \underset{\vec{w}}{\text{argmin}} \vec{q}^* \vec{q}^* - (\vec{q}^* \vec{w})^T \vec{q}^* - \vec{q}^T X \vec{w} + (\vec{q}^* \vec{w})^T X \vec{w} = \boxed{X^T X \text{ is positive definite}} \\
 &= \underset{\vec{w}}{\text{argmin}} \vec{q}^* \vec{q}^* - 2\vec{w}^T X^T \vec{q}^* + \vec{w}^T X^T X \vec{w} = F(\vec{w}) \quad \textcolor{red}{\downarrow} \quad \nabla_{\vec{w}} F(\vec{w}) = \vec{0} \\
 &= \underset{\vec{w}}{\text{argmin}} \vec{q}^* \vec{q}^* - 2\vec{w}^T X^T \vec{q}^* + \vec{w}^T X^T X \vec{w} = F(\vec{w})
 \end{aligned}$$

From the deductions, we have
 $\vec{q} = \vec{x} - (\vec{x}^T \vec{x})^{-1} \vec{x}^T \vec{y}$ $\Rightarrow P_x = \vec{x} \vec{x}^T$, being a projection matrix + $\vec{q} = \vec{y}$

Considering a matrix U , which is orthogonal and contains a basis for $C(X)$.
 $\vec{q} = X \vec{w} = U \vec{z}$ $\rightarrow (C(X)) = (C(U))$
 charge of
 leads to an orthonormal basis

$$\tilde{\vec{y}} = \underset{\vec{v}}{\operatorname{arg\!min}} \| \vec{q} - U\vec{v} \|^2 \Rightarrow \tilde{\vec{y}} = U(U^T U)^{-1} U^T \vec{q} = U(I)U^T \vec{q} = UU^T \vec{q}$$

$$\begin{aligned}
 X &= \begin{bmatrix} 1 & 1 \\ 1 & 3 \\ 0 & 0 \end{bmatrix}, \quad \vec{x}_1 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \quad \vec{x}_2 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \\
 U_1^T &= \frac{\vec{x}_1}{\|\vec{x}_1\|} = \frac{\begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}}{\sqrt{2}}, \quad U_2^T = \frac{\vec{x}_2}{\|\vec{x}_2\|} = \frac{\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}}{1} \\
 \vec{x}_2^D &= \vec{x}_2 - (\vec{x}_2^T U_1^T) U_1^T \quad \text{remove } \vec{x}_1 \\
 &= \vec{x}_2 - (\vec{x}_2^T U_1^T) U_1^T \vec{x}_2^D \\
 U_2 &= \frac{\vec{x}_2^D}{\|\vec{x}_2^D\|} = \frac{\begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix}}{\sqrt{2}} \quad \sim \sim \sim \rightarrow U = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & 0 \\ 0 & 0 \end{bmatrix} \rightarrow U^T U = I \\
 V &= \vec{v}_1 \vec{v}_2^T + \vec{v}_2 \vec{v}_1^T
 \end{aligned}$$

Drawbacks

$$\begin{aligned}
 X &= U \Sigma V^T \\
 \vec{\omega} &= (X^T X)^{-1} X^T \vec{y} = \\
 &= (\underbrace{V \Sigma^T U^T}_A \Sigma V^T)^{-1} V \Sigma^T U^T \vec{y} = \\
 &= (\underbrace{V \Sigma^T}_A \underbrace{V^T V}_I)^{-1} V \Sigma^T U^T \vec{y} = \\
 &= V (\underbrace{U^T \Sigma}_A)^{-1} V \Sigma^T U^T \vec{y} = \quad \text{using } V^T V = I \Rightarrow V^{-1} = V^T \\
 &= V (\Sigma^T \Sigma)^{-1} \underbrace{V^T \Sigma^T U^T \vec{y}}_B = \quad \text{using } (\Sigma^T)^{-1} = \Sigma \\
 &= V (\Sigma^T \Sigma)^{-1} \Sigma^T \underbrace{U^T \vec{y}}_B \quad \Sigma \rightarrow \text{pseudo-inverse of } \Sigma \sim \vec{y} = X (X^T X)^{-1} X^T \vec{y} \\
 &= V (\Sigma^T \Sigma)^{-1} \Sigma^T \vec{y} \\
 &\boxed{\vec{\omega} = V \Sigma^+ U^T \vec{y}}
 \end{aligned}$$

Assuming test $p \geq n$ in $X^{(n \times p)}$, X has n L.I. rows

$$X = \boxed{\underline{\underline{w}}}^T \rightarrow \text{no solutions for } \vec{w}$$

$$\vec{y} = X \vec{w}$$

$$\vec{z} = V \vec{z}^T \vec{U}^T \vec{q} \rightarrow \vec{\Sigma}^T = \vec{z}^T (\vec{z} \vec{z}^T)^{-1}$$

$$\hookrightarrow \vec{\Sigma}^T = \begin{bmatrix} \vec{v}_1 & 0 \\ 0 & \vec{v}_2 \end{bmatrix}$$

$$\vec{w}, \vec{w}^T, \vec{y} = X \vec{w} = X \vec{w}^T$$

$$\|\vec{w}\|_2^2 = \|\vec{w} - \vec{v}\|_2^2 = \|\vec{v} - \vec{w}\|_2^2 \geq \|\vec{v}\|_2^2 + \|\vec{w}\|_2^2 \Rightarrow \|\vec{w}\|_2^2 \geq \|\vec{v}\|_2^2$$

for any two vectors \vec{w} and \vec{v} ,
 \vec{w} will always have the minimum norm

$$(\vec{w} - \vec{v})^T \vec{v} = (\vec{v} - \vec{w})^T (\vec{v} - \vec{w}) =$$

$$= (X(\vec{w} - \vec{v}))^T (X(\vec{v}))^{-1} \vec{v} =$$

$$= (X(\vec{v} - \vec{w}))^T (X(X)^{-1})^T \vec{v} = \vec{B}_F$$

$$X = \begin{bmatrix} 1 & 0 & 0.02 \\ 0 & 1 & 0 \end{bmatrix} \quad \vec{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\vec{w}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \vec{w}_2 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \text{complex errors}$$

shifting 1st column
 shifting 2nd column
 improbable 3rd column

Matrix Completion

$X \in \mathbb{R}^{n \times p}$ rank(X) = $r \ll \min(n, p)$ \rightarrow there exists a low-rank representation of rank r

$$X_{ij} \text{ for } (i, j) \in \Omega \quad X = \begin{bmatrix} \vec{z}_1 & | & \vec{z}_2 & | & \vec{z}_3 \\ \vec{z}_1 & | & \vec{z}_2 & | & \vec{z}_3 \\ \vec{z}_1 & | & \vec{z}_2 & | & \vec{z}_3 \end{bmatrix} \quad \vec{z}_k \rightarrow \text{known values}$$

$$\left\{ \begin{array}{l} \text{"Ideal" estimator} \quad \text{non-convex problem} \\ \hat{X} = \underset{Z \in \mathbb{R}^{n \times p}}{\operatorname{argmin}} \operatorname{rank}(Z), \text{ subject to } \hat{X}_{ij} = X_{ij}, (i, j) \in \Omega \\ \hat{X} = X_i \text{ with } p = r \\ \text{"Frobenius" estimator} \quad \text{non-convex problem} \\ \hat{X} = \underset{Z \in \mathbb{R}^{n \times p}}{\operatorname{argmin}} \|Z\|_F^2, \text{ subject to } \hat{X}_{ij} = X_{ij}, (i, j) \in \Omega \end{array} \right.$$

\hookrightarrow the result is very close to the ideal one.

SVT (Singular Value Thresholding) non-monotone algorithm

$$\begin{aligned} &\bullet \text{ Initialize } \hat{X} = \text{zeros}(n, p) \\ &\bullet \text{ Set } \hat{X}_{ij} = X_{ij} \text{ for } (i, j) \in \Omega \\ &\bullet \text{ For } k = 1, 2, \dots, N \\ &\quad \hat{X}_{\text{old}} = \hat{X} \\ &\quad [U, \Sigma, V^T] = \text{SVD}(\hat{X}) \\ &\quad \Sigma \rightarrow \hat{\Sigma} \quad \begin{cases} \sigma_{ii} & \text{if } i < T \\ 0 & \text{otherwise} \end{cases} \quad \left(\text{reduces the rank of the matrix} \right) \\ &\quad \hat{X} = U \hat{\Sigma} V^T \\ &\quad \hat{X}_{ij} = X_{ij}, (i, j) \in \Omega \\ &\bullet \text{ If } \|\hat{X} - \hat{X}_{\text{old}}\|^2 \geq \epsilon \rightarrow \text{Iteration return} \end{aligned}$$

$$\begin{aligned}
 & (x_i, y_i), i=1 \dots n \\
 & x_i \in \mathbb{R}^p, y_i \in \mathbb{R}, \vec{y} \in \mathbb{R}^n \\
 & X \in \mathbb{R}^{n \times p} \\
 & \text{Hypothesis: } X \text{ has } p \text{ linearly independent columns} \\
 & \text{Thus } \epsilon_1, \dots, \epsilon_p > 0
 \end{aligned}
 \quad \left\{ \begin{array}{l}
 \vec{\hat{y}} = X \vec{\hat{w}}_{LS} \quad \vec{y} \approx X \vec{w}_{LS} \quad \vec{y} = X \vec{w}^* + \vec{\epsilon} \\
 \vec{\hat{w}}_{LS} = (X^T X)^{-1} X^T \vec{y} = \\
 = (X^T X)^{-1} X^T (X \vec{w}^* + \vec{\epsilon}) = \\
 = (X^T X)^{-1} (X^T X \vec{w}^* + X^T \vec{\epsilon}) = \\
 = (X^T X)^{-1} (X^T X \vec{w}^* + (X^T X)^{-1} X^T \vec{\epsilon}) = \\
 = (X^T X)^{-1} (\vec{w}^* + (X^T X)^{-1} X^T \vec{\epsilon}) = \\
 = \vec{w}^* + (X^T X)^{-1} X^T \vec{\epsilon}
 \end{array} \right. \quad \begin{array}{l}
 \text{weight vector of the underlying reality} \\
 \text{choose } X = U \Sigma V^T \\
 (X^T X)^{-1} X^T = V \Sigma^+ U^T \\
 \epsilon_0 \approx 10^{-9} \quad \text{if singular value decomposition is used, then } \epsilon_0 \text{ is not zero} \\
 \text{recognifies the error by factor } 10^3, \vec{w}_{LS} \neq \vec{w}^* \text{ very different}
 \end{array} \quad \left[\begin{array}{c|c}
 \frac{1}{\epsilon_1} & \\
 \frac{1}{\epsilon_2} & \\
 \vdots & \\
 \frac{1}{\epsilon_p} & \\
 \hline & 0
 \end{array} \right]$$

Ridge Regression

$$\vec{w}_R = \underset{\vec{w}}{\operatorname{arg\,min}} \| \vec{y} - X\vec{w} \|_2^2 + \lambda \| \vec{w} \|_2^2$$

$$\nabla f(\vec{w}) = \|\vec{y} - X\vec{w}\|_2^2 + \lambda \|\vec{w}\|_2^2 = \vec{y}^T \vec{y} - 2\vec{w}^T X^T \vec{y} + \lambda \vec{w}^T \vec{w} = 0$$

$$\vec{w}_p = (X^T X + \lambda I)^{-1} X^T \vec{y}$$

$$\begin{aligned}\vec{y}^* &= (\vec{X}^T \vec{X} + \lambda I)^{-1} \vec{X}^T \vec{y} = \\ &= (\vec{X}^T \vec{X} + \lambda I)^{-1} \vec{X}^T (\vec{X} \vec{w} + \vec{\epsilon}) = \\ &= (\vec{X}^T \vec{X} + \lambda I)^{-1} \vec{X}^T \vec{X} \vec{w} + (\vec{X}^T \vec{X} + \lambda I)^{-1} \vec{X}^T \vec{\epsilon}\end{aligned}$$

if $\tilde{E}' \approx 0$, \tilde{w}_L is a little bit worse than w_L^*
 However, if not, \tilde{w}_L will regularize small singular values, which w_L^* does not

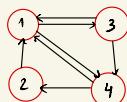
$$\begin{aligned}
 X &= UZV^T \\
 \overline{U}^T &= (VZ^TU^TVZV^T + \lambda VV^T)^{-1}VZ^TU^T\overline{V}^T \\
 &= [V\left(Z^TUVZ + \lambda I\right)V^T]^{-1}VZ^TU^T\overline{V}^T \\
 &= V\left(\cancel{Z^TUVZ} + \lambda I\right)^{-1}VZ^TU^T\overline{V}^T \\
 &= V\left(\cancel{Z^T} + \lambda I\right)^{-1}Z^TU^T\overline{V}^T
 \end{aligned}$$

$$E^+ = \begin{bmatrix} 1/\epsilon_2 & 0/\epsilon_2 & \dots & | & 0 \end{bmatrix}$$

$$(\Sigma^T \Sigma + \lambda I) \Sigma^T = \begin{bmatrix} \frac{\sigma_1^2}{\sigma_1^2 + \lambda} & & \\ & \ddots & \\ & & \frac{\sigma_p^2}{\sigma_p^2 + \lambda} \end{bmatrix}$$

{ if $\sigma_i \gg \lambda$, diagonal terms are $\propto \frac{1}{\sigma_i}$
 if $\sigma_i \approx 0$, diagonal terms are ≈ 0
 2 presents exploding terms
 regularizes

PAGE - RANK



Idea: Surf the web
randomly (random walk
on the graph)

$$\tilde{A} = \begin{bmatrix} 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 \end{bmatrix}$$

\Downarrow probability matrix

$$= \begin{bmatrix} 0 & 1 & \frac{1}{3} & \frac{1}{2} \\ 0 & 0 & \frac{1}{3} & \frac{1}{2} \\ \frac{1}{2} & 0 & 0 & 0 \\ \frac{1}{2} & 0 & \frac{1}{3} & 0 \end{bmatrix}$$

$$A\vec{e}_3 = \begin{bmatrix} 1/3 \\ 4/3 \\ 0 \\ 1/3 \end{bmatrix}$$

$$\vec{\pi} = \begin{bmatrix} \pi_1 \\ \pi_L \\ \pi_3 \\ \pi_4 \end{bmatrix}$$

if look for

$$A\vec{\pi} = \vec{\pi}$$

$A > 0$, Perron-Frobenius
 \Downarrow
 $\lambda_1 > \lambda_2 > \lambda_3 > \dots > \lambda_n$

Power - Method

$$\vec{\pi}^{(0)} \quad \| \vec{\pi}^{(0)} \| = 1$$

for $k = 1, 2, \dots$

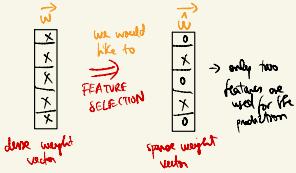
$$\vec{\pi}^{(k)} = \frac{A \vec{\pi}^{(k-1)}}{\| A \vec{\pi}^{(k-1)} \|}$$

stop condition

$$\| \vec{\pi}^{(k)} - \vec{\pi}^{(k-1)} \| < \epsilon$$

$$\begin{aligned}
 & \text{A can be diagonalized, } \sqrt{\lambda} \\
 & \hookrightarrow A = V \Lambda V^{-1} \\
 & \vec{v}_1^{(1)} = \frac{\vec{v}_1^{(1)}}{\|\vec{v}_1^{(1)}\|} = \frac{\vec{v}_1^{(1)}}{\|\Lambda^{\frac{1}{2}} \vec{v}_1^{(1)}\|} \\
 & \vec{v}_1^{(2)} = \frac{\vec{v}_1^{(2)}}{\|\vec{v}_1^{(2)}\|} = \frac{\vec{v}_1^{(2)}}{\|\Lambda^{\frac{1}{2}} \vec{v}_1^{(2)}\|}, \text{ then} \\
 & \vec{v}_1^{(K)} = \frac{\vec{v}_1^{(K)}}{\|\vec{v}_1^{(K)}\|} = \frac{\vec{v}_1^{(K)}}{\|\Lambda^{\frac{1}{2}} \vec{v}_1^{(K)}\|} \\
 & \vec{v}_2^{(1)} = \alpha_2 \vec{v}_2^{(1)} = \alpha_2 (\vec{v}_2^{(1)}) + \sum_{i=2}^K \alpha_i \vec{v}_i^{(1)} \\
 & \vec{v}_2^{(K)} = \alpha_2 \vec{v}_2^{(K)} + \dots + \alpha_K \vec{v}_2^{(K)} = V \left(\alpha_2 \lambda_2^{\frac{1}{2}} \vec{v}_2^{(1)} + \dots + \alpha_K \lambda_K^{\frac{1}{2}} \vec{v}_2^{(K)} \right) \\
 & = \alpha_2 \lambda_2^{\frac{1}{2}} \vec{v}_2^{(1)} + \dots + \alpha_K \lambda_K^{\frac{1}{2}} \vec{v}_2^{(K)} = \alpha_2 \lambda_2^{\frac{1}{2}} \left(v_2 + \frac{\alpha_3}{\lambda_2} \lambda_2^{\frac{1}{2}} \vec{v}_2^{(2)} + \dots + \frac{\alpha_K}{\lambda_2} \lambda_2^{\frac{1}{2}} \vec{v}_2^{(K)} \right) = \\
 & \xrightarrow{\text{similarly } \lambda_2 > \lambda_3 \left(\frac{\lambda_2}{\lambda_3} \right)^K \rightarrow 0 \text{ when } K \rightarrow \infty} \\
 & \Rightarrow \text{converges}
 \end{aligned}$$

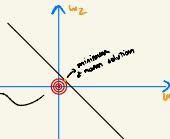
$\vec{y} = X \vec{w}$ For LS and RR, $\vec{w} \in \mathbb{R}^p$ and is dense



$$\vec{y} = X \vec{w} \quad X \in \mathbb{R}^{n \times p}, p > n$$

Example:
2 features, 1 sample ($n=2, p=2$)

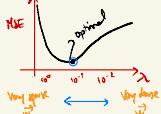
$$X = [2, 3] \quad y = [1] \quad \sim z = [2 \ 3] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \Rightarrow z = 2w_1 + 3w_2 \quad \text{line (infinite solutions)}$$



$$\|z\|^2 = \sum_{i=1}^2 (z_i)^2, z \in \mathbb{R}^n$$

$$\|z\|^2 = \|X\vec{w}\|^2 \rightarrow \text{more norm}$$

$\hookrightarrow F(\vec{w}) = \|X\vec{w} - \vec{y}\|^2 + \lambda \|\vec{w}\|^2$
LASSO (Least Absolute Shrinkage and Selection Operator)
and **SVR** (Support Vector Regression)



$$F(\vec{w}) = \|X\vec{w} - \vec{y}\|^2_L + \lambda_1 \|\vec{w}\|_1 + \lambda_2 \|\vec{w}\|^2$$



$$\vec{w}_R = V(\Sigma^T \Sigma + \lambda I)^{-1} \Sigma^T \vec{y} \quad \text{Ridge regression}$$

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

$$\vec{w}_R = V\Sigma^T (\Sigma^T \Sigma + \lambda I)^{-1} \Sigma^T \vec{y} = V\Sigma^T \vec{U} \vec{U}^T (\Sigma^T \Sigma + \lambda I)^{-1} \Sigma^T \vec{y} =$$

$$= \vec{X} \vec{\alpha} \quad \vec{\alpha} \in \mathbb{R}^n$$

$$\vec{w}_R = \vec{X}^T \vec{\alpha} = \sum_{i=1}^n w_i x_i \vec{x}^T \quad \hat{y}_i = \vec{x}_i^T \vec{w} = w_1 x_{1i} + w_2 x_{2i} + \dots + w_n x_{ni}$$

Add features

$$\hat{y}_i = w_1 x_{1i} + w_2 x_{2i} + w_3 x_{3i}^2 + w_4 x_{4i}^2 + w_5 x_{5i} x_{2i}$$

Using $\varphi(\vec{x})$ in feature map

$$\varphi(\vec{x}) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} \quad \hat{y}_i = \varphi(\vec{x})^T \vec{w}$$

$\vec{w} \in \mathbb{R}^{d+1}$

Kernel method: aim of this method is to avoid the necessity of computing high orders

(in case of computing high orders)

$$\varphi(\vec{x}_1)^T \varphi(\vec{x}_2) = \vec{x}_1^T \vec{x}_2 + 2x_1 x_2 + x_1^2 + x_2^2$$

$$\vec{w}_R = \varphi^T \vec{\alpha} \quad \vec{\alpha} = V(\Sigma \varphi^T + \lambda I)^{-1} \Sigma^T \vec{y} =$$

$$= (\vec{X} \varphi^T + \lambda I)^{-1} \vec{y} \rightarrow \text{how convex feature map}$$

KER

$$K_{ij} = \varphi(\vec{x}_i)^T \varphi(\vec{x}_j) = K(x_i, x_j)$$

$$\vec{\alpha} = (\vec{X} \varphi^T)^{-1} \vec{y} \rightarrow \text{How to predict a label for a new feature vector}$$

$$\hat{y} = \varphi(\vec{x})^T \vec{\alpha} = \vec{\alpha}^T \varphi(\vec{x}) =$$

$$= (\vec{X} \varphi^T) \vec{\alpha}^T = \vec{\alpha}^T \varphi(\vec{x}) = \sum_{i=1}^n \alpha_i \varphi(x_i)^T \varphi(\vec{x}) = \sum_{i=1}^n \alpha_i K(x_i, \vec{x})$$

Kernel trick allows to create non-linear models by generating the feature map without having to extract it by computing all features

Kernel Function

$$K(\vec{x}_1, \vec{x}_2) = \vec{x}_1^T \vec{x}_2$$

$$\vec{x}_1^2 = \begin{bmatrix} x_{11} \\ x_{12} \\ x_{13} \\ x_{14} \\ x_{15} \end{bmatrix} \quad \varphi(\vec{x}_1) = \begin{bmatrix} 1 \\ x_{11} \\ x_{12} \\ x_{13} \\ x_{14} \\ x_{15} \end{bmatrix}$$

$$\varphi(\vec{x}_1)^T \varphi(\vec{x}_2) = x_{11}^2 + 2x_{11}x_{21} + x_{21}^2 + x_{12}^2 + 2x_{12}x_{22} + x_{22}^2 + \dots$$

$$K(\vec{x}_1, \vec{x}_2) = (\vec{x}_1^T \vec{x}_2)^2 = (x_{11}x_{21} + x_{12}x_{22})^2 = \text{(4)}$$

i) Polynomial of degree q

$$K(\vec{x}_1, \vec{x}_2) = (\vec{x}_1^T \vec{x}_2)^q$$

ii) Polynomial of degree $< q$

$$K(\vec{x}_1, \vec{x}_2) = (\vec{x}_1^T \vec{x}_2 + 1)^q$$

iii) Gaussian Function

$$K(\vec{x}_1, \vec{x}_2) = \exp\left(-\frac{||\vec{x}_1 - \vec{x}_2||^2}{2\sigma^2}\right)$$

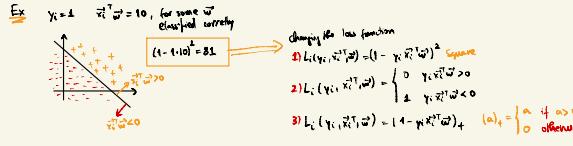
$X \in \mathbb{R}^{n \times n}$ (cloudy entries)
 $\Phi = X^T V$ factors
 $V = U S V^T$

$$\begin{aligned} C &= \sum_{i=1}^n V_i^T X \rightarrow \text{exp. factorize as } C = V A V^T \\ &\quad \text{where } A = \sum_{i=1}^n X_i^T \text{ (diagonal elements)} \\ &\quad \Phi = X V \\ &\quad \Phi = X V = U S V^T = \sum_{i=1}^n U_i S_i V_i^T \\ &\quad \Phi = X V = U S V^T = \sum_{i=1}^n U_i S_i V_i^T \end{aligned}$$

$$\vec{w}_k^T = \underset{\vec{w}}{\operatorname{arg\,min}} \|y - \vec{w}^T \vec{w}\|_2^2 + \lambda \|\vec{w}\|_2^2$$

($y_i^T \vec{w}_k$) = binary classification problem $y_i \in \{-1, 1\}, x_i \in \mathbb{R}^p$

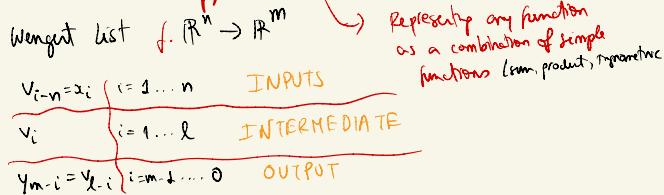
 $\vec{w}_{k+1}^T = \underset{\vec{w}}{\operatorname{arg\,min}} \sum_i (y_i - \vec{w}^T \vec{w})^2 = \underset{\vec{w}}{\operatorname{arg\,min}} \sum_i (1 - x_i^T \vec{w})^2 \quad \text{if } k=2$
 $\vec{w}_{k+1}^T = \underset{\vec{w}}{\operatorname{arg\,min}} \sum_i (1 - x_i^T \vec{w})^2 = \underset{\vec{w}}{\operatorname{arg\,min}} \sum_i (1 + x_i^T \vec{w})^2 \quad \text{if } k=2$



Derivative calculation

	Pros	Cons
Manual	Exact, Good for proofs	Prone to error Expensive for complex functions
Numerical Differentiation	Easy to program	FP precision Computational cost
Automatic Differentiation	Exact Fast	Implementation

Automatic Differentiation



$$y = f(x_1, x_2) = \sin((x_1 + x_2) x_2^2)$$

$v_{-1} = x_1 = 1$ (suppose) INPUTS

$v_0 = x_2 = 2$ INPUTS

$v_1 = x_1 + x_2 = 3$ INTERMEDIATE

$v_2 = x_2^2 = 4$ INTERMEDIATE

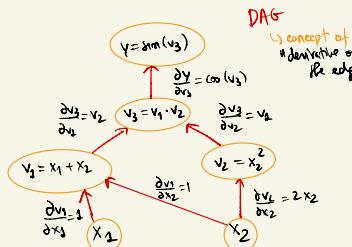
$v_3 = v_1 \cdot v_2 = 12$ INTERMEDIATE

$v_4 = \sin(v_3) \approx 0.536$ INTERMEDIATE

$y = v_4$ OUTPUTS

$$\begin{aligned} \dot{v}_2 &= \dot{x}_2 = 1 && \text{derivative w.r.t } x_2 \\ \dot{v}_0 &= \dot{x}_2 = 0 && \text{derivative w.r.t } x_2 \\ \dot{v}_1 &= \dot{x}_1 + \dot{x}_2 = 1 \\ \dot{v}_2 &= 0 \\ \dot{v}_3 &= \dot{v}_1 \dot{v}_2 + v_1 \dot{v}_2 = 4 \\ \dot{v}_4 &= \cos(v_3) \dot{v}_3 \approx 3.37 \\ \dot{y} &= \dot{v}_4 = 3.37 \end{aligned}$$

Forward Mode of AD



$y = f(x_1, x_2) = \sin((x_1 + x_2) \times x_2)$

two forward passes
for NN issue with FM

$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$, $m \gg n$

DAG

$y = v_4 = \sin(v_3)$

$\frac{\partial y}{\partial x_1} = \frac{\partial y}{\partial v_1} \cdot \frac{\partial v_1}{\partial x_1} = \cos(v_3) \cdot 1 = \cos(v_3)$

$J = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_n} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$

$\vec{x} = \vec{e}_i$

$y_j = \frac{\partial y_j}{\partial x_i}, i=1 \dots m$

similarly to Newton's method
 \Rightarrow able to avoid the explicit construction of the whole matrix through more tricks
 \hookrightarrow save memory

$V_1 = x_1 = 1$
 $V_0 = x_2 = 2$
 $V_2 = x_1 + x_2 = 3$
 $V_3 = x_2^2 = 4$
 $V_4 = v_1 = \sin(v_3)$

$v_1' = x_1' = 1$
 $v_0' = x_2' = 0$
 $v_2' = x_1' + x_2' = 1$
 $v_3' = 2x_2' = 0$
 $v_4' = v_1' v_2' + v_1 v_2' = 4$
 $v_4' = v_2^2 = 16$
 $v_4 = \cos(v_3) = \cos(\sin(v_3)) = \cos(16)$

$\frac{\partial y}{\partial x_1}$
 $\frac{\partial y}{\partial x_2}$

Dual numbers

$a+b\varepsilon: a, b \in \mathbb{R}, \varepsilon \neq 0, \varepsilon^2=0$

Operations

$(a+b\varepsilon) + (c+d\varepsilon) = (a+c) + (b+d)\varepsilon$

$(a+b\varepsilon) \cdot (c+d\varepsilon) = ac + ad\varepsilon + bc\varepsilon + bd\varepsilon^2 = ac + (ad+bc)\varepsilon$

$f(a+b\varepsilon) = f(a) + f'(a)b\varepsilon + \frac{f''(a)}{2}b^2\varepsilon^2 \dots$

$+ b=1 \quad f(a+\varepsilon) = f(a) + \varepsilon f'(a)$

Ex.: $f(x) = g(x) h(x)$

$f(a+\varepsilon) = g(a+\varepsilon) h(a+\varepsilon) =$
 $= (g(a) + g'(a)\varepsilon)(h(a) + h'(a)\varepsilon) =$
 $= g(a)h(a) + g'(a)h(a)\varepsilon + g(a)h'(a)\varepsilon + g'(a)h'(a)\varepsilon^2 =$
 $= g(a)h(a) + (g'(a)h(a) + g(a)h'(a))\varepsilon$

$f(x) = g(h(x))$

$f(a+\varepsilon) = g(h(a+\varepsilon)) =$
 $= g(h(a) + \varepsilon h'(a)) =$
 $= g(h(a)) + g'(h(a))h'(a)\varepsilon$

$f(x) = \frac{1}{x}$
 $f(x+\varepsilon) = \frac{1}{x+\varepsilon} = \frac{x-\varepsilon}{(x+\varepsilon)(x-\varepsilon)} = \frac{x-\varepsilon}{x^2} - \frac{x}{x^2} - \frac{\varepsilon}{x^2} = \frac{1}{x} - \frac{1}{x^2}\varepsilon$

Backward (reverse) mode
 \hookrightarrow sensitivity of the output w.r.t. the inputs

$\bar{v}_i = \frac{\partial y}{\partial v_i}$

$\frac{\partial v_1}{\partial v_1} = 1$
 $\bar{v}_3 = \frac{\partial v_1}{\partial v_3} = \cos(v_3)$
 $\bar{v}_2 = \frac{\partial v_1}{\partial v_2} = v_2$
 $\bar{v}_1 = \frac{\partial v_1}{\partial v_1} = \frac{\partial v_1}{\partial v_2} \frac{\partial v_2}{\partial v_1} = v_2 \cos(v_3)$
 $\bar{v}_1 = \frac{\partial v_1}{\partial v_1} = \frac{\partial v_1}{\partial v_3} \frac{\partial v_3}{\partial v_1} = v_3 \cos(v_3)$
 $\bar{v}_0 = \frac{\partial v_1}{\partial v_0} = \frac{\partial v_1}{\partial v_3} \frac{\partial v_3}{\partial v_0} = v_3 \cos(v_3)$
 $\bar{v}_1 = \frac{\partial v_2}{\partial v_1} = \frac{\partial v_2}{\partial v_3} \frac{\partial v_3}{\partial v_1} = v_3 \cos(v_3)$
 $\bar{v}_0 = \frac{\partial v_2}{\partial v_0} = \frac{\partial v_2}{\partial v_3} \frac{\partial v_3}{\partial v_0} = v_3 \cos(v_3)$
 $\bar{v}_1 = \frac{\partial v_3}{\partial v_1} = \frac{\partial v_3}{\partial v_2} \frac{\partial v_2}{\partial v_1} = v_2 \cos(v_3)$
 $\bar{v}_0 = \frac{\partial v_3}{\partial v_0} = \frac{\partial v_3}{\partial v_2} \frac{\partial v_2}{\partial v_0} = v_2 \cos(v_3)$
 $\bar{v}_1 = \frac{\partial v_3}{\partial v_1} = \frac{\partial v_3}{\partial v_0} = 2v_0$
 $\bar{v}_0 = \frac{\partial v_3}{\partial v_0} = \frac{\partial v_3}{\partial v_0} = 2v_0$
 $\bar{v}_1 = \frac{\partial v_4}{\partial v_1} = \frac{\partial v_4}{\partial v_3} \frac{\partial v_3}{\partial v_1} = 2v_0 \cos(v_3)$
 $\bar{v}_0 = \frac{\partial v_4}{\partial v_0} = \frac{\partial v_4}{\partial v_3} \frac{\partial v_3}{\partial v_0} = 2v_0 \cos(v_3)$
 $\bar{v}_1 = \frac{\partial v_4}{\partial v_2} = \frac{\partial v_4}{\partial v_3} \frac{\partial v_3}{\partial v_2} = v_3 \cos(v_3)$
 $\bar{v}_0 = \frac{\partial v_4}{\partial v_0} = \frac{\partial v_4}{\partial v_3} \frac{\partial v_3}{\partial v_0} = v_3 \cos(v_3)$
 $\bar{v}_1 = \frac{\partial v_4}{\partial v_1} = \frac{\partial v_4}{\partial v_0} = 2v_0 \cos(v_3)$
 $\bar{v}_0 = \frac{\partial v_4}{\partial v_0} = \frac{\partial v_4}{\partial v_0} = 2v_0 \cos(v_3)$

$f: \mathbb{R}^n \rightarrow \mathbb{R}^m, m \ll n$

\hookrightarrow evaluation of f costs $ops(f)$

• FM $\Rightarrow \sim m \cdot c \cdot ops(f)$ operations

• BM $\Rightarrow \sim m \cdot c \cdot ops(f)$ operations

c ≈ 2.5
 Jacobian computation
 clear FM or BM according to sign of m/n

autodiff.org

Griewank book

• FM
 $\text{If } \vec{x}^2 = \vec{r}^2 \Rightarrow \vec{x}^T \vec{r}$
 - BM
 If $\vec{y} = \vec{r} \Rightarrow \vec{y}^T \vec{r}$
 ↗ sensitivity vector

Exercise: Write the tangent line for BM $f(x_1, x_2) = \sin((x_1+x_2)x_2)$

$$\begin{aligned} v_1 &= x_2 \\ v_2 &= x_2 \\ v_3 &= x_1 + x_2 \\ v_4 &= x_1 \\ v_5 &= x_1 + x_2 \\ v_6 &= x_1 + x_2 \end{aligned}$$

$$\begin{aligned} \vec{v}_1 &= \frac{\partial \vec{v}}{\partial v_1} = 1 \\ \vec{v}_2 &= \frac{\partial \vec{v}}{\partial v_2} = \cos(v_3) \\ \vec{v}_3 &= \frac{\partial \vec{v}}{\partial v_3} = \sin(v_3) \cdot v_2 \\ \vec{v}_4 &= \frac{\partial \vec{v}}{\partial v_4} = \frac{\partial \vec{v}}{\partial v_5} = \cos(v_3) \cdot v_2 \\ \vec{v}_5 &= \frac{\partial \vec{v}}{\partial v_5} = \frac{\partial \vec{v}}{\partial v_6} = \cos(v_3) \cdot v_2 \\ \vec{v}_6 &= \frac{\partial \vec{v}}{\partial v_6} = \frac{\partial \vec{v}}{\partial v_1} = \cos(v_3) \cdot v_2 \end{aligned}$$

$$\begin{aligned} \vec{v}_6 &= \frac{\partial \vec{v}}{\partial v_6} = \frac{\partial v_6}{\partial v_1} \frac{\partial v_3}{\partial v_1} \frac{\partial v_2}{\partial v_3} + \frac{\partial v_6}{\partial v_2} \frac{\partial v_3}{\partial v_2} = \\ &= (\cos(v_3))' \cdot v_2 + \cos(v_3) \cdot v_2 \cdot 2 \cdot v_0 = \cos(v_3) (v_2 + 2v_1 v_0) \\ \vec{v}_3 &= \frac{\partial \vec{v}}{\partial v_3} = \frac{\partial v_3}{\partial v_1} \frac{\partial v_2}{\partial v_1} = (\cos(v_3))' v_2 \end{aligned}$$

Newton's Method
 $h(\vec{x}) \rightsquigarrow (H\vec{x})$
 ↗ useful calculate

f: $\mathbb{R}^n \rightarrow \mathbb{R}, \vec{x} \rightarrow y$

reverse-on-forward

$$i) \nabla f \cdot \vec{v} = \frac{d f}{d \vec{v}} \quad \vec{x} = \vec{v} \quad (\text{forward step})$$

ii) Apply the backward mode to the result of i)

$$\nabla^2 f \cdot \vec{v} = H\vec{v}$$

Convolution

Functions

$$(f * g)(x) = \int_{-\infty}^{\infty} f(t)g(x-t)dt$$

Vectors

$$(\vec{c} * \vec{d})_k = \sum_{i+j=k} c_i d_j$$

Cyclic Convolutions

$$(\vec{c} \circledast \vec{d})_k = \sum_{cyclic\;order} c_i d_j$$

$$\begin{aligned} c_0 + c_1 x + \dots + c_{n-1} x^{n-1} &= c(x) \\ d_0 + d_1 x + \dots + d_{n-1} x^{n-1} &= d(x) \\ c(x)d(x) \xrightarrow{*} &\Rightarrow c_0 d_0 + c_1 d_{n-1} + \dots \\ &\text{sum of indices are always } K \end{aligned}$$

Convolution

Toeplitze matrix
Elements are recursive.
Elements are given by a
(2n-3) length sequence
 $\{t_k : -n+1 \leq k \leq n-1\}$

$$\begin{matrix} T(1,1) = t_{-n+1} \\ T(1,2) = t_{-n+2} \\ \vdots \\ T(n,n) = t_n \end{matrix}$$

Cyclic Convolution

Convolution
For each rowvector, the elements are given by a
(2n-3) length sequence
 $\{t_k : -n+1 \leq k \leq n-1\}$

$$C(c, d) = C(c_0, d_0)$$

$$C(c, d) = \begin{bmatrix} c_0 & c_1 & \dots & c_{n-1} \\ c_0 & c_1 & \dots & c_{n-1} \\ \vdots & \vdots & \ddots & \vdots \\ c_0 & c_1 & \dots & c_{n-1} \end{bmatrix}$$

$$ex: \begin{bmatrix} 3 & 8 & 5 & 3 \\ 5 & 3 & 8 & 5 \\ 8 & 5 & 3 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ c_4 \end{bmatrix} = \begin{bmatrix} c_2 \\ c_3 \\ c_4 \\ c_1 \end{bmatrix}$$

$$\vec{c} = c_0 \mathbb{I} + c_1 \vec{x} + c_2 \vec{x}^2 + c_3 \vec{x}^3 \quad (\text{polynomial in } \vec{x})$$

$$D = d_0 \mathbb{I} + d_1 \vec{x} + d_2 \vec{x}^2 + d_3 \vec{x}^3$$

$$CD = ? \Rightarrow CD \text{ is also diagonal}$$

$$CD = (c_0 \mathbb{I} + c_1 \vec{x} + c_2 \vec{x}^2 + c_3 \vec{x}^3)(d_0 \mathbb{I} + d_1 \vec{x} + d_2 \vec{x}^2 + d_3 \vec{x}^3) =$$

$$\stackrel{\text{Gauss Rule}}{\sim} \vec{x}^n \Rightarrow \vec{x} = I$$

$$\vec{x} = I$$

SGD

Cost function : $J(w) = \frac{1}{N} \sum_{i=1}^N J_i(w)$

Algorithm : sample x_k randomly

$$w^{(k+1)} = w^{(k)} - \eta_k \nabla J_{i_k}(w^{(k)})$$

Simple convergence results for SGD

$$\|w^{(k+1)} - w^*\|^2 = \|J(w^{(k+1)}) - J(w^*)\|$$

measurements
to prove
convergence

Theorem: The following statements are equivalent to the condition that a diff function is μ -strongly convex

- i) $f(y) \geq f(x) + \nabla f(x)^T(y-x) + \frac{\mu}{2} \|y-x\|^2$
- ii) $g(x) = f(x) - \frac{\mu}{2} \|x\|^2$ is convex
- iii) $(\nabla f(x) - \nabla f(y))^T(y-x) \geq \mu \|y-x\|^2$

$$A) \# \|w^{(k+1)} - w^*\|^2 \leq (1-2\eta\mu)^K \|w^{(0)} - w^*\|^2 + \frac{\eta^2 c^2}{2\mu}$$

convergence
results in
expectation

$$B) \# [J(w^{(k)})] - J(w^*) \leq (1-2\mu\eta)^K [J(w^{(0)} - \nabla J(w^*)) + \frac{L\eta^2}{2\mu}]$$

more
details
in notes

Then: Assume y^* is the minimizer of J and

- i) ∇J is L -Lipschitz $\rightarrow J(y) \leq J(x) + \nabla J(x)^T(y-x) + \frac{L}{2} \|y-x\|^2$
- ii) J is μ -strongly convex $\rightarrow J(y) \geq J(x) + \nabla J(x)^T(y-x) + \frac{\mu}{2} \|y-x\|^2$
- iii) $\|\nabla J(x)\| \leq c, \forall x, \forall i$
- iv) $\sigma < \eta < \frac{1}{L}$
- v) $\# [g_K] = \nabla J(w^*)$

$$(1) \|w^{(k+1)} - w^*\| = \|(\nabla J_{i_k}(w^{(k)})) - w^*\|^2 =$$

$$= \|w^{(k)} - w^*\|^2 + 2\eta_k \nabla J_{i_k}(w^{(k)})^T(w^{(k)} - w^*) + \eta_k^2 \|\nabla J_{i_k}(w^{(k)})\|^2$$

$$\leq \|w^{(k)} - w^*\|^2 + 2\eta_k \nabla J_{i_k}(w^{(k)})^T(w^{(k)} - w^*) + \eta_k^2 c^2$$

$$\# [\|w^{(k+1)} - w^*\|^2] \leq \# [\|w^{(k)} - w^*\|^2 + 2\eta_k \nabla J_{i_k}(w^{(k)})^T(w^{(k)} - w^*) + \eta_k^2 c^2]$$

$$\leq \# [\|w^{(k)} - w^*\|^2] + 2\eta_k \# [\|(\nabla J_{i_k}(w^{(k)})) - w^*\|^2] + \eta_k^2 c^2 =$$

$$= (1-2\eta\mu) \# [\|w^{(k)} - w^*\|^2] + \eta_k^2 c^2$$

$$\# [\|w^{(k+1)} - w^*\|^2] \leq (1-\eta\mu)^K \|w^{(0)} - w^*\|^2 + \sum_{i=0}^K (1-2\eta\mu) \eta_i^2 c^2$$

$$= (1-\eta\mu)^K \|w^{(0)} - w^*\|^2 + \frac{\eta^2 c^2}{2\mu} =$$

$$= (1-\eta\mu)^K \|w^{(0)} - w^*\|^2 + \frac{m^2}{2\mu}$$

$$(2) J(w^{(k+1)}) \leq J(w^{(k)}) + \nabla J(w^{(k)})^T(w^{(k+1)} - w^{(k)}) + \frac{L}{2} \|w^{(k+1)} - w^{(k)}\|^2$$

$$= J(w^{(k)}) + \nabla J(w^{(k)})^T(w^{(k)} - \eta_k \nabla J_{i_k}(w^{(k)}) - w^{(k)}) + \frac{L}{2} \|\eta_k \nabla J_{i_k}(w^{(k)})\|^2$$

$$\leq J(w^{(k)}) - \eta_k \nabla J(w^{(k)})^T \nabla J_{i_k}(w^{(k)}) + \frac{L\eta_k^2 m^2}{2}$$

$$\# [J(w^{(k+1)}) - J(w^*)] \leq J(w^{(k)}) - J(w^*) - \eta_k \|\nabla J(w^{(k)})\|^2 + \frac{L\eta_k^2 m^2}{2}$$

$$\leq J(w^{(k)}) - J(w^*) - 2\eta_k \mu (J(w^{(k)}) - J(w^*)) + \frac{L\eta_k^2 m^2}{2}$$

$$= (1-2\eta_k \mu) (J(w^{(k)}) - J(w^*)) + \frac{L\eta_k^2 m^2}{2}$$

(iterate K times)

$$\# [J(w^{(k)}) - J(w^*)] \leq (1-2\eta_k \mu)^K (J(w^{(0)}) - J(w^*)) + \sum_{i=0}^K (1-2\eta_k \mu) \frac{L\eta_k^2 m^2}{2}$$

$$\frac{1}{2\eta_k \mu} \cdot \frac{L\eta_k^2 m^2}{2} = \frac{Lc^2 m^2}{4\mu}$$

$$(f(y) \geq f(x) + \nabla f(x)^T(y-x) + \frac{\mu}{2} \|y-x\|^2)$$

LHS when $y = \Sigma - \frac{1}{\mu} \nabla J(x)$
when $y = \Sigma - \frac{1}{\mu} \nabla J(w^*)$

$$J(w^*) \geq J(w) - \frac{1}{2\mu} \|\nabla J(w)\|^2$$

Line search Procedure

$$\# w^{(k+1)} = w^{(k)} - \eta_k p_k$$

$$\# p_k \text{ s.t. } p_k^T \nabla J(w^{(k)}) < 0$$

$$p_k = -\frac{\nabla J(w^{(k)})}{\|\nabla J(w^{(k)})\|}$$

$$J(x) = x_1 - x_2 + 2x_1 x_2 + 2x_1^2 + x_2^2$$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \nabla J(x) = \begin{bmatrix} 1+2x_2 + 4x_1^2 \\ -2+2x_1 + 2x_2 \end{bmatrix}$$

$$m_k = \arg \min J(w^{(k)} - \eta_k \nabla J(w^{(k)}))$$

$$\begin{cases} x^{(0)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} & \nabla J(x^{(0)}) = \begin{bmatrix} 1 \\ -2 \end{bmatrix} \\ J(x^{(0)} - \eta_k \nabla J(x^{(0)})) = \end{cases}$$

$$= J(-\eta_1, \eta_2) = \eta_2^2 - 2\eta_2$$

$$\beta = \frac{dJ}{dm_2} = 2\eta_2 - 2 = 2(m_2 - 1) \Rightarrow \eta_2 = 1$$

ADA GRAD

η for rare features, η is big
for features appearing often, η is small

$$\eta_{i_k} \propto \frac{1}{\sum_i \nabla J_i}$$

problem no
the summation grows,
 η get small and
the optimization gets stuck

ADADELT A
RHS Prop

Averaging the
gradients

ADAM

Universal Approx. Theorem

add first theory on error bounds
is not available yet

Comparison PINN vs FEM
on background

low data (\hookrightarrow much physics) } comparison:
some data (\hookrightarrow some physics) } show that
large data (\hookrightarrow large physics) PINN can solve
any of them
(unlike PINN
in 1D)