

Matrix Methods

MIT 18.065

in Data Analysis, Signal Processing, and Machine Learning

2020. 11. 29

Sungchan Yi 

Matrix Multiplication.

$$A = \begin{bmatrix} 2 & 3 \\ 4 & 7 \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & 0 \\ 2 & 1 \end{bmatrix}}_L \underbrace{\begin{bmatrix} 2 & 3 \\ 0 & 1 \end{bmatrix}}_U^T = l_1 u_1^T + l_2 u_2^T$$

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

$$= \begin{bmatrix} 2 & 3 \\ 4 & 6 \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 3 \\ 4 & 7 \end{bmatrix}$$

Eigenvalue of Rotations

$$R_\theta = \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} \quad \det R_\theta = 1, \quad R_\theta \in O(2)$$

$$R_{\frac{\pi}{2}} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} = A \quad p_A(t) = t^2 + 1 = 0$$

$$\therefore \lambda = \pm i \quad (\text{Imaginary})$$

Reflections in \mathbb{R}^n

• Reflection always has an axis

(Some hyperplane of $n-1$ dimension)

Take $n-1$ orthonormal vectors of this hyperplane

: u_1, u_2, \dots, u_{n-1}

Then the reflection does not change the vector.

$\therefore S u_i = u_i$ for $i=1, 2, \dots, n-1$

$\Rightarrow 1$ is an eigenvalue. (multiplicity $n-1$)

$\exists v \in (\text{hyperplane})^\perp$ then $Sv = -v$.

$\Rightarrow -1$ is an eigenvalue

$\therefore S$ has a full set of independent eigenvectors.

S is diagonalizable

$$S \sim \text{diag}(-1, 1, \dots, 1), \quad \det S = -1.$$

Orthogonal Matrices have Orthogonal Eigenvectors.

Proof. Let $Q \in O(n)$, λ, μ be eigenvalues of Q , where $\lambda \neq \mu$.

$$\exists v, w \in \mathbb{C}^n \text{ s.t. } Qv = \lambda v, \quad Qw = \mu w$$

We show that $v^* w = 0$.

$$Qv = \lambda v \Rightarrow v^* Q^* = \bar{\lambda} v^*$$

$$\Rightarrow v^* Q^* \bar{\lambda} w = \bar{\lambda} v^* Qw$$

$$\Rightarrow v^* w = \bar{\lambda} \mu v^* w$$

Since eigenvalues of orthogonal matrices have length 1, $\bar{\lambda} \mu \neq 0$ and $v^* w = 0$ \blacksquare

$$(\det Q = \pm 1)$$

Singular Value Decomposition.

$$A \in M_{m,n}(\mathbb{F})$$

$$\left\{ \begin{array}{l} Av_1 = \sigma_1 u_1 \\ Av_2 = \sigma_2 u_2 \\ \vdots \\ Av_r = \sigma_r u_r \end{array} \right. \quad \begin{array}{l} (v_i \in \mathbb{F}^m, u_i \in \mathbb{F}^n) \\ (\sigma_1 \geq \sigma_2 \geq \dots \geq 0) \end{array}$$

$$A \begin{bmatrix} v_1 & \dots & v_r \end{bmatrix} = \begin{bmatrix} u_1 & \dots & u_r \end{bmatrix} \Sigma$$

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

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

Because V is orthogonal, $A = U \Sigma V^T$

Such U, V exist because $A^T A$, $A A^T$ are positive semi-definite

\rightarrow Has orthogonal eigenvectors. (basis)

Singular values of A are $\sqrt{\text{eigenvalues of } A^T A}$

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

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

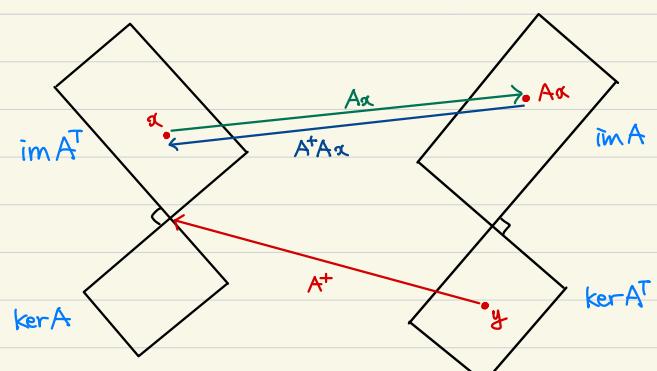
$$r = \text{rk } A \text{ and } \sigma_{r+1}, \dots, \sigma_{\min(m,n)} = 0$$

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

For positive semi-definite matrices,

SVD is equal to diagonalization.

Four Subspaces.



$$A^T A x = x \text{ for all } x \in \text{im } A^T$$

$$\text{but } A^T y = 0 \text{ for all } y \in \ker A^T$$

If A is invertible, $A^+ = A^{-1}$.

If A is not invertible, ... SVD!

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

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

$$= U \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r & \\ & & & 0 & \\ & & & & \ddots & \\ & & & & & 0 \end{bmatrix} V^T$$

But this only works when A has independent columns.

$\rightarrow A^T A$ is invertible.

Proof. A has independent columns $\Rightarrow \ker A = 0$.

Suppose $A^T A \alpha = 0$ then $A\alpha \in \ker A^T$.

since $\ker A^T \perp \text{im } A$, and $A\alpha \in \text{im } A$,

$A\alpha \in \ker A^T \cap \text{im } A = \{0\}$. $\therefore A\alpha = 0$.

But $\ker A = 0$ thus $\alpha = 0$.

$\therefore \ker A^T A = 0$ and $A^T A$ is invertible.

If $\ker A = 0$, (or $\text{rk } A = r = n$)

$$A^+ = (A^T A)^{-1} A^T$$

(Note that $A A^T$ is not invertible)

If A has independent rows, $A A^T$ is invertible.

$$A^+ = A^T (A A^T)^{-1}$$

Low Rank Approximation.

Minimize $\|A - A_k\|_2$ over A_k subject to $\text{rk } A_k = k$.

Eckart-Young-Mirsky Theorem.

Let $A = U \Sigma V^T \in M_{m,n}(\mathbb{R})$

Define $U_k \in M_{m,k}(\mathbb{R})$, $\Sigma_k \in M_{k,k}(\mathbb{R})$,

$V_k \in M_{n,k}(\mathbb{R})$ as

$$U = \begin{bmatrix} U_1 & * \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_1 & 0 \\ 0 & * \end{bmatrix}, \quad V = \begin{bmatrix} V_1 & * \end{bmatrix}$$

Then, $A_k = U_k \Sigma_k V_k^T$.

Proof. WLOG. $m \geq n$.

Let $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_k, \dots, \sigma_m)$

$$A = \sum_{i=1}^n \sigma_i u_i v_i^T, \quad A_k = \sum_{i=1}^k \sigma_i u_i v_i^T$$

$$\text{Thus } \|A - A_k\|_2 = \left\| \sum_{i=k+1}^m \sigma_i u_i v_i^T \right\|_2 = \sigma_{k+1}.$$

Suppose $B = X Y^T$ where X, Y have k columns.

$$\Rightarrow \text{rk } B = k.$$

We show that $\|A - A_k\|_2 = \sigma_{k+1} \leq \|A - B\|_2$

Since Y has k columns, there exists a non-trivial linear combination of the first $k+1$ columns of V ,

$$w = c_1 v_1 + \dots + c_{k+1} v_{k+1} \text{ such that } Y^T w = 0.$$

WLOG. scale $\|w\|_2 = 1$ ($c_1^2 + \dots + c_{k+1}^2 = 1$)

$$\text{Then, } \|A - B\|_2^2 \geq \|(A - B)w\|_2^2$$

$$= \|Aw\|_2^2$$

$$= c_1^2 \sigma_1^2 + \dots + c_{k+1}^2 \sigma_{k+1}^2$$

$$= \sigma_{k+1}^2 + c_1^2 (\sigma_1^2 - \sigma_{k+1}^2) + \dots$$

$$\geq \sigma_{k+1}^2 \quad \blacksquare$$

Matrix Norms.

Spectral Norm. (Maximum singular value)

$$\|A\|_2 = \sigma_1 = \max_{\alpha \in V} \frac{\|A\alpha\|_2}{\|\alpha\|_2}$$

Frobenius Norm

$$\|A\|_F = \sqrt{\sum_{i,j} |a_{ij}|^2} = \sqrt{\text{tr}(A^* A)} = \sqrt{\sum_i \sigma_i^2}$$

Nuclear Norm (Sum of singular values)

$$\|A\|_N = \sum_i |\sigma_i|$$

Least Squares

$A : m \times n$ and $\text{rk } A = r$

Solve $Ax = b$!

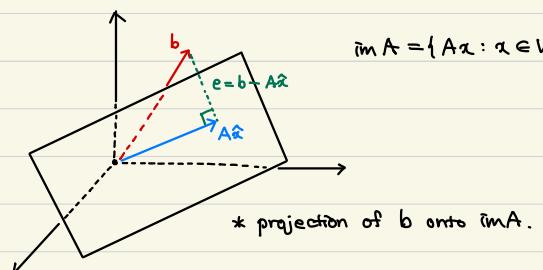
If $m = n = r$, $x = A^{-1}b$.

Otherwise, we minimize $\|Ax - b\|_2^2$

$$\|Ax - b\|_2^2 = (Ax - b)^T (Ax - b)$$

$$= x^T A^T A x - 2b^T A x + b^T b$$

\Rightarrow has minimum at $A^T A \hat{x} = A^T b$



10. Survey of Difficulties with $Ax = b$

- ① Always works. $x = A^+b$.
- ② Good normal case : Square matrix
Condition number σ_1/σ_n is OK. $x = A^{-1}b$
- ③ $m > n = \text{rank } A$ $A^T A \hat{x} = A^T b$
 $A^T A$ is not too big (reasonable to compute)
- ④ $m < n$ (Not enough equations - Deep Learning)
 - minimum norm A^+b (ℓ^2 norm)
Good algorithm
 - min ℓ^1 norm → pick some algorithm and solve.
Good results in test data.
 - under determined / many solutions
- ⑤ Columns are nearly almost linearly dependent
(In bad condition) → Gram Schmidt
 $A = QR$ (Q : orthogonal, R : triangular)
 - Column pivoting: reordering columns
 - If pivot is too slow, then reorder the rows. $PA = LU$
- ⑥ Nearly singular / Inverse problems \Rightarrow Add penalty.
 - Minimize $\|Ax - b\|^2 + \delta^2 \|x\|^2$
try to make it invertible
- ⑦ Way too big → Iterative methods (conjugate gradient)
- ⑧ Randomized numerical linear algebra
Sample the columns/rows of A ...

$$\text{SVD: } A = U \Sigma V^T \rightarrow A^+ = V \Sigma^+ U^T$$

Nearly singular → small σ_i 's
 $\rightarrow \Sigma^{-1}$ has large singular values.

ℓ^2 norm : minimize $\|Ax - b\|^2 + \delta^2 \|x\|^2$

Choose $\delta > 0 \rightarrow$ Solvable problem.

$$\begin{bmatrix} A \\ \delta I \end{bmatrix} x = \begin{bmatrix} b \\ 0 \end{bmatrix} \quad (\text{Augmented system})$$

$A^* x = b^*$ ← Apply least squares here.

- What should δ be?
- What happens as $\delta \rightarrow 0$?

$$(A^*)^T A^* \hat{x} = (A^*)^T b^*$$

$$\Rightarrow [A^T \delta I] \hat{x} = [A^T \delta I] \begin{bmatrix} b \\ 0 \end{bmatrix}$$

$$\Rightarrow (A^T A + \delta^2 I) \hat{x} = A^T b \quad \dots (x)$$

positive definite (invertible)

Suppose A is 1×1 $A = [\sigma]$ $b = [b]$

$$\text{Then } (x) \Rightarrow (\sigma^2 + \delta^2)x = \sigma b \quad x = \frac{\sigma}{\sigma^2 + \delta^2} b$$

$$(i) \sigma > 0 : \text{as } \delta \rightarrow 0, x \rightarrow \frac{b}{\sigma} \text{ OK.}$$

$$(ii) \sigma = 0 : x = 0$$

... Similar to pseudo inverses. $\delta \rightarrow 0, x \rightarrow A^+ b$!!!

For any A . $(A^T A + \delta^2 I)^{-1} A^T \rightarrow A^+$ as $\delta \rightarrow 0$

(Note that we are staying with ℓ^2 norms)

$$p(t) A = U \Sigma V^T \quad A^T A = V \Sigma^T V^T = V \Sigma^T \Sigma V^T$$

$$A^T A + \delta^2 I = V \Sigma^T \Sigma V^T + \delta^2 V V^T = V (\Sigma^T \Sigma + \delta^2 I) V^T$$

Let $D = \Sigma^T \Sigma + \delta^2 I$ (diagonal, D^{-1} exists).

$$\therefore (A^T A + \delta^2 I)^{-1} A^T = V D^{-1} V^T / \Sigma^T V^T = V D^{-1} \Sigma^T V^T$$

Observe that $D = (d_{ij})_{n \times n}$ where

$$d_{ij} = \begin{cases} \sigma_i^2 + \delta^2 & \text{if } i=j \leq r \\ 0 & \text{otherwise} \end{cases}$$

$$\text{therefore } D^{-1} \Sigma^T = (c_{ij})_{n \times m} \text{ where}$$

$$c_{ij} = \begin{cases} \sigma_i / (\sigma_i^2 + \delta^2) & \text{if } i=j \leq r \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Since } \lim_{\delta \rightarrow 0} \frac{\sigma_i}{\sigma_i^2 + \delta^2} = \frac{1}{\sigma_i}, \quad D^{-1} \Sigma^T \rightarrow \Sigma^T$$

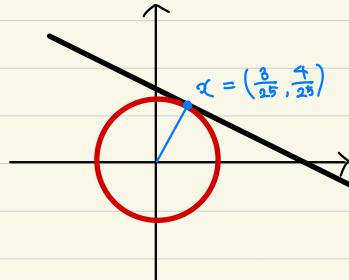
$$\therefore (A^T A + \delta^2 I)^{-1} A^T \rightarrow V \Sigma^T V^T = A^+ \text{ as } \delta \rightarrow 0 \blacksquare$$

11. Minimizing $\|x\|$ Subject to $Ax = b$

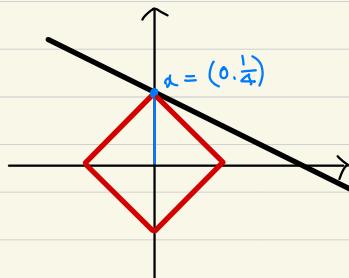
Suppose we want to minimize $\|x\|$,

$$\text{with } 3x_1 + 4x_2 = 1 \quad (\text{Geometrically})$$

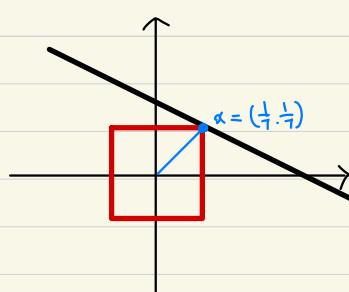
ℓ_2 norm



ℓ_1 norm



ℓ_∞ norm



* Gram - Schmidt

1. Standard way $A \rightarrow QR$

2. Column exchanges

3. Krylov - Arnoldi - Lanczos

Let $A = [a_1 \dots a_n]$.

Gram-Schmidt $A = QR$ where Q is orthogonal

What is R ?

→ Combinations of q_i will give columns of A .

→ Since Q is orthogonal, $Q^{-1} = Q^T$.

$$\therefore R = Q^T A = \begin{bmatrix} -q_1^T \\ \vdots \\ -q_n^T \end{bmatrix} \begin{bmatrix} a_1 \dots a_n \end{bmatrix}$$

Then (i,j) -component of R should be $q_i^T a_j$

(dot product of q_i and a_j)

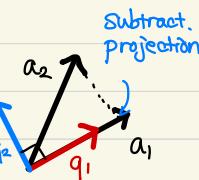
Gram-Schmidt Process.

For given n vectors a_1, \dots, a_n .
 a_1

Let $q_1 = \frac{a_1}{\|a_1\|}$

$$u_2 = a_2 - \text{proj}_{q_1} a_2 \Rightarrow q_2 = \frac{u_2}{\|u_2\|}$$

$$= a_2 - (a_2^T q_1) q_1$$



Repeat similarly, ...

$$u_3 = a_3 - (a_3^T q_1) q_1 - (a_3^T q_2) q_2$$

$$\Rightarrow q_3 = \frac{u_3}{\|u_3\|}$$

$$\therefore u_i = a_i - (a_i^T q_1) q_1 - \dots - (a_i^T q_{i-1}) q_{i-1}$$

$$\Rightarrow q_i = \frac{u_i}{\|u_i\|} \quad \text{for } i=1, \dots, n.$$

Check that $q_i^T q_j = \delta_{ij}$ (Kronecker Delta)

2. Column pivoting possible

: How to decide q_2

→ Compute $a_i - (a_i^T q_1) q_1$ for all $i=2, \dots, n$
 and set u_2 as the maximum..

$$\text{Now } q_2 = \frac{u_2}{\|u_2\|}.$$

Computing all $a_i - (a_i^T q_1) q_1$ isn't costly since
 $(a_i^T q_1) q_1$ was going to be subtracted anyway.

3. Krylov. : Large, sparse A .

Solve $Ax=b$...

Sparse → matrix \times vector is easy.

Compute b , Ab , A^2b , ..., $A^{n-1}b$.
 iterate $Ax = (A^k b)$
 dimension

Their combination gives the Krylov space K_J

x_j = best vector in K_j (closest least squares solution)

Why is orthonormal basis so good?

Let $\alpha = c_1 q_1 + c_2 q_2 + \dots + c_n q_n = Q \begin{bmatrix} c_1 \\ \vdots \\ c_n \end{bmatrix} = QC$.

Then $C = Q^{-1}\alpha = Q^T\alpha$!!!

$$q_i^T \alpha = c_i q_i^T q_i = c_i \quad (\text{since } q_i^T q_j = 0 \text{ if } i \neq j)$$

Any matrix Symmetric matrix
 Aronldi - Lanczos found a great way to orthogonalize the vectors in K_j . (Using standard Gram-Schmidt)

* Trefethen - Bau * Golub - Van Loan

12. Computing Eigenvalues and Singular Values

Eigenvalues of A . - QR Method.

Let $A_i = Q_i R_i$ (orthogonalize, R : upper triangular)

Set $A_{i+1} = R_i Q_i$ for $i=0, 1, \dots$

Claim. $A_i \sim A_{i+1}$

Proof. $A_{i+1} = R_i Q_i = R_i A_i R_i^{-1}$ ↗ Can't I use A_i^{-1} ?

After each iteration, components below the diagonals tend to get smaller (in absolute value)

⇒ Then A_k is nearly upper triangular, and has λ_i 's on the main diagonal.

* Introducing shifts (shifting eigenvalues)

Shift by εI .: $A_i - \varepsilon I = Q_i R_i$, $A_{i+1} = R_i Q_i + \varepsilon I$

→ Eigenvectors do not change, eigenvalues are shifted.

Does $A_i \sim A_{i+1}$ still hold?

$$A_{i+1} = R_i Q_i + \varepsilon I = R_i (A_i - \varepsilon I) R_i^{-1} + \varepsilon I$$

$$= R_i A_i R_i^{-1} - R_i (\varepsilon I) R_i^{-1} + \varepsilon I = R_i A_i R_i^{-1}$$

Since λ_m is usually computed first, ε can be set to λ_m .

Non-zeros directly under the main diagonal : Hessenburg Matrix.

* If A_i is nearly upper-triangular → easy to compute.

Computing Eigenvalues.

① Reduce A (by similarity transformation) to Hessenburg form.

② Apply QR decomposition with shifts

LAPACK

If A is symmetric, its Hessenburg form will be tridiagonal.

→ reduced order of operations.

Symmetric S . then find $Q = Q_0 Q_1 \dots$ such that

$Q S Q^{-1}$ is tridiagonal (and has same eigenvalues)

Multiplying orthogonal matrices does not change σ_i 's.

σ_i 's are eigenvalues of $A^T A$.

Then for QA , $(QA)^T QA = A^T Q^T QA = A^T A$

thus the singular values of QA are eigenvalues of $A^T A$

Then if $A = U \Sigma V^T$, multiply left/right hand sides by orthogonal matrices Q_1, Q_2 so that

$$A' = (Q_1 U) \Sigma (V^T Q_2) \rightarrow \text{bidiagonal} !!$$

Now $(A')^T A$ is tridiagonal, now σ_i can be calculated.

using the QR decomposition.

But these methods are $O(n^3)$

13. Randomized Matrix Multiplication

For matrices that are **HUGE!!!**

$$AB = \begin{bmatrix} 1 & \dots & n \end{bmatrix} \begin{bmatrix} -b_1^T \\ \vdots \\ -b_n^T \end{bmatrix} = a_1 b_1^T + \dots + a_n b_n^T$$

IDEA: Choose columns randomly with probability $\sum p_i = 1$.

+ Choose the probability that minimize the variance

↑ Normalize.

$$p_i = \|a_i\| \|b_i^T\| / c \quad \text{where } c = \sum_{i=1}^n \|a_i\| \|b_i^T\|$$

↑ probability of choosing $a_i \times b_i^T$ (to multiply)

the probability is proportional to the norm.

s: number of samples

With probability p_i , choose $\frac{a_i b_i^T}{s p_i}$

Then the approximate AB is $\sum_{i=1}^s \frac{a_i b_i^T}{s p_i}$

Mean of 1 sample is $\sum_{i=1}^n p_i \frac{a_i b_i^T}{s p_i} = \frac{\sum a_i b_i^T}{s}$

→ s samples give $\sum_{i=1}^s \frac{a_i b_i^T}{s} = \frac{\sum a_i b_i^T}{s} = AB$

Variance (will depend on p_i)

$$\sum_{i=1}^n \frac{\|a_i\| \|b_i^T\|^2}{s p_i} - \frac{1}{s} \|AB\|_F^2 = \sum_{i=1}^n \frac{1}{s} \|a_i\| \|b_i^T\|^2 - \frac{1}{s} \|AB\|_F^2$$

variance.

$$= \frac{1}{s} (c^2 - \|AB\|_F^2)$$

→ s samples : $c^2 - \|AB\|_F^2$ (with optimal p_i)

* How to choose p_i such that the variance is minimized?

Minimize $\sum_{i=1}^n \frac{\|a_i\| \|b_i^T\|^2}{p_i}$ with the constraint $\sum p_i = 1$

Lagrange Multipliers.

$$\text{Let } F = \sum_{i=1}^n \frac{\|a_i\| \|b_i^T\|^2}{p_i} + \lambda (\sum p_i - 1)$$

$$\frac{\partial F}{\partial p_i} = -\frac{\|a_i\| \|b_i^T\|^2}{p_i^2} + \lambda = 0 \Rightarrow p_i^2 = \frac{\|a_i\| \|b_i^T\|^2}{\lambda} \quad \text{for each } i=1, \dots, n$$

Since p_i 's add up to 1, $\sum p_i = \frac{\sum \|a_i\| \|b_i^T\|^2}{\sqrt{\lambda}} = 1$.

$$\therefore \sqrt{\lambda} = \sum_{i=1}^n \|a_i\| \|b_i^T\| = c. \quad \therefore p_i = \frac{\|a_i\| \|b_i^T\|^2}{c} = \frac{\|a_i\| \|b_i^T\|}{\sqrt{\lambda}}$$

14. Low Rank Changes in A and its Inverse

Start with I and perturb it by uv^T (rank 1 matrix).

What is the inverse of $I - uv^T$?

→ Matrix Inversion Formula

$$(I - uv^T)^{-1} = I + \frac{uv^T}{1 - v^T u} \rightarrow |x|.$$

→ The inverse also has rank 1 perturbation.

$$\text{Proof. } (I - uv^T)(I + \frac{uv^T}{1 - v^T u})$$

$$= I + \frac{uv^T}{1 - v^T u} - uv^T - \frac{uv^T uv^T}{1 - v^T u}$$

$$= I + \frac{uv^T}{1 - v^T u} - uv^T - \frac{u(v^T u)v^T}{1 - v^T u}$$

$$= I + \frac{(1 - v^T u)uv^T}{1 - v^T u} - uv^T = I + uv^T - uv^T = I$$

How about for rank k matrix?

Let $U \in M_{n \times k}$, $V \in M_{k \times m}$. UV^T is rank k.

$$(I_n - UV^T)^{-1} = I_n + U(I_k - V^T U)^{-1} V^T$$

K x K inverse

Proof) $(I_n - UV^T)(I_n + U(I_k - V^T U)^{-1} V^T)$

$$= I_n - UV^T + U(I_k - V^T U)^{-1} V^T - UV^T U(I_k - V^T U)^{-1} V^T$$

$$= I_n - UV^T + U(I_k - V^T U)^{-1} V^T - V^T U(I_k - V^T U)^{-1} V^T$$

$$= I_n - UV^T + U(I_k - V^T U)^{-1} V^T$$

$$= I_n - UV^T + UV^T = I_n$$

Finally, (Woodbury Matrix Identity)

$$(A - UV^T)^{-1} = A^{-1} - A^{-1} U (I_k - V^T A^{-1} U)^{-1} V^T A^{-1}$$

To use the formula.

① Solve $(A - UV^T)x = b$

Suppose $Aw = b$ is solved for w A = LU and back-substitute. (once)

Also solve $Az = w$ for z

Now solve $(A - uv^T)x = b$ for x, quickly!

$$\therefore x = w + \frac{zu^T w}{1 - v^T z}$$

$$\text{Proof). } (A - uv^T)(w + \frac{zu^T w}{1 - v^T z})$$

scalar

$$= Aw - uw^T w + A \frac{zu^T w}{1 - v^T z} - uv^T \frac{zu^T w}{1 - v^T z}$$

$$= b - uw^T w + \frac{zu^T w}{1 - v^T z} - (v^T z) \frac{zu^T w}{1 - v^T z}$$

$$= b - uw^T w + uw^T w = b$$

② New measurement in least squares!

(OLD) $Ax = b \xrightarrow{\text{normal}} A^T A \hat{x} = A^T b$ (Recursive LSA)

Make 1 more measurement!

$$\left[\begin{array}{c} A \\ v^T \end{array} \right] x = \left[\begin{array}{c} b \\ b_{\text{new}} \end{array} \right] \leftarrow j \text{ new row.}$$

$$\xrightarrow{\text{normal}} [A^T \ v] \left[\begin{array}{c} A \\ v^T \end{array} \right] \hat{x}_{\text{new}} = [A^T \ v] \left[\begin{array}{c} b \\ b_{\text{new}} \end{array} \right]$$

This is actually $(A^T A + vv^T) \hat{x}_{\text{new}} = A^T b + v b_{\text{new}}$
doesn't change rank 1 change in ATA.

* Kalman Filter : Dynamic least squares.

- Idea of using the covariance matrix. (correlation)
- (Control Theory) State equation.

Least Squares : Standard.

Data not correlated, covariance matrix is I.

15. Matrices $A(t)$ and Derivative dA/dt

Setup $A(t)$ is a matrix depending on t.

A is invertible and we know $\frac{dA}{dt}$. What is $\frac{dA^{-1}}{dt}$?

Notice that $B^{-1} - A^{-1} = B^{-1}(A - B)A^{-1}$

$$\Rightarrow \Delta A^{-1} = (A + \Delta A)^{-1}(-\Delta A)A^{-1}$$

divide by Δt and let $\Delta t \rightarrow 0$. then

$$\frac{dA^{-1}}{dt} = -A^{-1} \frac{dA}{dt} A^{-1} \quad \text{Compare with } |x| \text{ case.}$$

$A(t) = t$

$A - B$ and $B^{-1} - A^{-1}$
have the same rank
multiplying by invertible matrix
doesn't change the rank.

* Changes in Eigenvalues

Let $A(t)x(t) = \lambda(t)x(t)$ (eigenvector of A)
 $y^T(t)A(t) = \lambda(t)y^T(t)$ (eigenvector of A^T)
 Norm given by $y^T x = 1$.

$$\Rightarrow AX = X\lambda, Y^TA = \lambda Y^T, Y^TX = I$$

Also, $y^T Ax = y^T(\lambda x) = \lambda$. Take the derivative!

$$\begin{aligned} \frac{dy}{dt} &= \frac{dy^T}{dt} A x + y^T \frac{dA}{dt} x + y^T A \frac{dx}{dt} \quad (\text{Product rule}) \\ &= \frac{dy^T}{dt} \lambda x + y^T \frac{dA}{dt} x + \lambda y^T \frac{dx}{dt} \\ &= \lambda \left(\frac{dy^T}{dt} x + y^T \frac{dx}{dt} \right) + y^T \frac{dA}{dt} x \\ &= \lambda \frac{d(y^T x)}{dt} + y^T \frac{dA}{dt} x = y^T \frac{dA}{dt} x \\ &\rightarrow \text{Does not involve } \frac{dy^T}{dt}, \frac{dx}{dt}! \end{aligned}$$

For symmetric matrix S. Note that uu^T is positive semi-definite.
 Eigenvalues of $S + uu^T$?

Let λ_i, μ_i be eigenvalues of S, $S + uu^T$.

Then the eigenvalues interlace!

$$\mu_1 \geq \lambda_1 \geq \mu_2 \geq \lambda_2 \geq \dots \geq \dots$$

Repeat the process and make rank 1 changes.

Let v_i be the eigenvalues of $S + uu^T + ww^T$, then

$$v_1 \geq \mu_1 \geq v_2 \geq \mu_2 \geq \dots$$

thus we know that $\lambda_1 \geq v_3 !!!$

Suppose $Su = \lambda_2 u$. (2nd unit eigenvector)

Look at $S + 20uu^T$.

$$(S + 20uu^T)u = Su + 20uu^Tu = (\lambda_2 + 20)u$$

How does it not contradict the interlacing?

16. Derivatives of Inverse and Singular Values

* Changes in Singular Values.

From $AV = U\Sigma$, $\sigma(t) = u^T(t)A(t)v(t)$.

$$\begin{aligned} \frac{du}{dt} &= \frac{du^T}{dt} Av + u^T \frac{dA}{dt} v + u^T A \frac{dv}{dt} \\ &= \frac{du^T}{dt} \sigma u + u^T \frac{dA}{dt} v + \sigma v^T \frac{du}{dt} = u^T \frac{dA}{dt} v \\ (\because u^T u = 1) \Rightarrow \frac{du^T}{dt} u + u^T \frac{du}{dt} &= 0 \Rightarrow \frac{du^T}{dt} u = 0 \end{aligned}$$

We can't get exact formulas but we can get bounds.

* Interlacing.

Eigenvalues of S : $\lambda_1 \geq \lambda_2 \geq \dots$

Eigenvalues of $S + kuu^T$: $\mu_1 \geq \mu_2 \geq \dots$

$$\mu_1 \geq \lambda_1 \geq \mu_2 \geq \lambda_2 \geq \dots$$

If $\mu_2 = \lambda_2 + k$, take k as large as μ_2 passes λ_1 .

\Rightarrow As k increases, μ_2 is now μ_1 and the inequality still holds.

* Weyl's Inequality

For symmetric S.T., $\lambda_{i+j-1}(S+T) \leq \lambda_i(S) + \lambda_j(T)$

\nearrow i-th eigenvalue of S

Compressed Sensing

Nuclear norm $\|A\|_N = \sigma_1 + \sigma_2 + \dots + \sigma_r$

Minimizing in L_1 norm makes the answer sparse.

17. Rapidly Decreasing Singular Values

What makes low rank matrices appear in problems?

Let $X \in M_{m,n}(\mathbb{R})$, singular values $\sigma_1, \dots, \sigma_n$.

Fact: If k is the largest integer s.t. $\sigma_{k+1} \leq 0 < \sigma_k$,

$$\textcircled{1} \quad \text{rk } X = k$$

$$\textcircled{2} \quad X = \sigma_1 u_1 v_1^T + \dots + \sigma_k u_k v_k^T$$

$$\textcircled{3} \quad \text{Rank Theorem.}$$

Solving Images
 $\frac{2km}{n} < n^2$
 \uparrow
 SVD
 U, V
 All entries

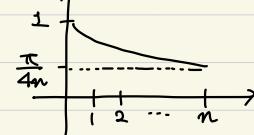
Def X is a low rank matrix if $\text{rk } X < \frac{n}{2}$.

Often, we demand $k \ll \frac{n}{2}$

Rank 1 matrix \rightarrow Highly aligned.

$$\begin{aligned} \text{Triangular Flag} \quad \square &= \begin{bmatrix} 1 & \dots & 0 \\ 1 & \ddots & \\ \vdots & \ddots & 1 \end{bmatrix} = X \\ X^{-1} &= \begin{bmatrix} 1 & \dots & 1 \\ -1 & \ddots & \\ \vdots & \ddots & -1 \end{bmatrix} \quad (X^T X)^{-1} = \begin{bmatrix} 1 & \dots & 1 \\ 1 & \ddots & \\ \vdots & \ddots & 1 \end{bmatrix} \\ \Rightarrow \text{Singular values } \sigma_k &= \sqrt{2 \sin \frac{\pi(2k-1)}{2(2n+1)}} \\ \sigma_1 \approx \frac{2n}{\pi}, \dots, \sigma_n \approx \frac{1}{2} & \end{aligned}$$

Plot of σ_k/σ_1



Japan Flag

$$\begin{aligned} \text{rk} \left(\begin{array}{c|c} \text{eye} & \text{all 0} \end{array} \right) &\leq \text{rk} \left(\begin{array}{c|c} \text{eye} & \text{all 0} \end{array} \right) + \text{rk} \left(\begin{array}{c|c} \text{all 0} & \text{eye} \end{array} \right) \\ &\leq \text{rk} \left(\begin{array}{c|c} \text{eye} & \text{all 0} \end{array} \right) + \text{rk} \left(\begin{array}{c|c} \text{all 0} & \text{eye} \end{array} \right) + 1 \\ &\leq 2 \left\lceil r \left(1 - \frac{1}{\sqrt{2}} \right) \right\rceil + 1 \end{aligned}$$

$$\text{rk } X = k$$

Numerical rank for $0 < \epsilon < 1$ (tolerance) is k if

k is the largest integer s.t. $\sigma_{k+1} \leq \epsilon \sigma_1 < \sigma_k$

By definition, $\text{rk}_0 X = \text{rk } X$

Eckart-Young : $\sigma_{k+1}, \|X - X_k\|_2$

(X_k : best rank k matrix)

* Matrices of low numerical rank.

All low rank matrices.

Hilbert matrix (full rank but low numerical rank)

$$H_{ij} = \frac{1}{i+j-1}$$

Vandermonde matrix $V = \begin{bmatrix} 1 & x_1 & \dots & x_1^{n-1} \\ 1 & x_2 & \dots & x_2^{n-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & \dots & x_n^{n-1} \end{bmatrix}$ ($x_i \in \mathbb{R}$)

Often we want V^{-1} , but hard to compute

The world is Smooth. (Reade, 1983)

$$\text{Ex. } p(x,y) := 1 + \alpha + xy$$

Define $\alpha_{ij} = p(i,j) = 1 + i + ij$.

$$X = \underbrace{\begin{bmatrix} 1 & \dots & 1 \\ 1 & \dots & 1 \\ \vdots & \vdots & \vdots \\ 1 & \dots & 1 \end{bmatrix}}_{\text{rank 1}} + \underbrace{\begin{bmatrix} 1 & \dots & 1 \\ 2 & \dots & 2 \\ \vdots & \vdots & \vdots \\ n & \dots & n \end{bmatrix}}_{\text{rank 1}} + \underbrace{\begin{bmatrix} 1 & \dots & 1 \\ 1 & \dots & 1 \\ \vdots & \vdots & \vdots \\ 1 & \dots & 1 \end{bmatrix}}_{\text{rank 1}}$$

$$\therefore \text{rk } X = 3.$$

$$\text{In general, } p(xy) = \sum_{s,t}^{\infty} \alpha_{st} x^s y^t$$

$$\alpha_{ij} = p(i,j) \text{ then } \text{rk } X \leq n^2$$

$$\text{For Hilbert matrix, } f(xy) = \frac{1}{|xy|-1}.$$

$$\text{Approximate } f \text{ by a polynomial. } \Rightarrow |f(xy) - p(xy)| \leq \frac{\varepsilon}{n} \|X\|_2$$

$$\text{Then matrix } Y = (y_{ij}), \text{ where } y_{ij} = p(i,j)$$

should have finite rank and $\|X - Y\|_2 \leq \varepsilon \|X\|_2$
(will have low numerical rank)

Doesn't work well ...

$$\text{rk } H_{1000} = 1000, \text{ rk}_\varepsilon H_{1000} = 28 \quad (\varepsilon = 10^{-15})$$

$$\text{But Reade's argument ... } \text{rk}_\varepsilon H_{1000} \leq 719$$

→ Not strong enough to show that H_{1000} is of low numerical rank.

* The world is Sylvester Sylvester Equation.

Many matrices satisfy $AX - XB = C$ for some A, B, C

Want to find A, B, C → X is of low numerical rank.

$$\begin{bmatrix} 1/2 & & & \\ & 3/2 & & \\ & & \ddots & \\ & & & n-1/2 \end{bmatrix} H - H \begin{bmatrix} -1/2 & & & \\ & -3/2 & & \\ & & \ddots & \\ & & & -n+1/2 \end{bmatrix} = \begin{bmatrix} 1 & \dots & 1 \\ 1 & \dots & 1 \\ \vdots & \vdots & \vdots \\ 1 & \dots & 1 \end{bmatrix}$$

$$\begin{bmatrix} x_1 & & & \\ & \ddots & & \\ & & x_m \end{bmatrix} V - V \begin{bmatrix} 0 & & & \\ & 1 & \dots & -1 \\ & & \ddots & \\ & & & 0 \end{bmatrix} = \begin{bmatrix} 0 & \dots & * \\ \vdots & \dots & * \\ 0 & \dots & * \end{bmatrix}$$

C happens to be low rank.

Ithm If X satisfies $AX - XB = C$ for normal matrices A, B, and let $r = \text{rk } C$, then

$$\sigma_{k+1} \leq Z_k(E, F)_\sigma,$$

* Z : Zolotarev number, E, F: Eigenvalue set of A, B

Key: E and F are separated.

$$\text{For } H: E \subseteq [\frac{1}{2}, n - \frac{1}{2}], F \subseteq [-n + \frac{1}{2}, -\frac{1}{2}]$$

$Z_k(E, F) \rightarrow 0$ as k increases.

$$\rightarrow \text{rk}_\varepsilon H_{1000} \leq 37. \quad (\text{Huge improvement})$$

18. Counting Parameters in SVD, LU, QR.

Counting free parameters in matrix factorization.

$$L: \frac{1}{2}n(n-1), \quad U: \frac{1}{2}n(n-1)$$

$$Q(\text{orthogonal}): \frac{1}{2}n(n-1)$$

$$\Lambda(\text{diagonal}): n$$

$$X(\text{eigenvectors}): n^2 - n \quad (\text{Normalize them})$$

$$S(\text{symmetric}): \frac{1}{2}n(n-1)$$

Let A be an $n \times m$ matrix with rank r.

$r \times r \Sigma$: has r parameters

$m \times r U$: has r orthonormal n -dimensional vectors.

$$(n-1) + (n-2) + \dots + (n-r) = rn - \frac{1}{2}r(r+1)$$

$r \times m V^T$: has r orthonormal m -dimensional vectors.

$$(m-1) + (m-2) + \dots + (m-r) = rm - \frac{1}{2}r(r+1)$$

$$\Rightarrow \text{Total: } (m+n)r - r^2 \quad \left(\begin{array}{c|cc} r \times & \vdots & \\ \hline ; & \ddots & \\ i: & & \times \end{array} \right) \quad \text{Automatically Desired.}$$

* Saddle Point from constraints.

Positive definite S.

Minimize $\frac{1}{2}\alpha^T S \alpha$ with respect to $A\alpha = b$.

$$\text{Lagrangian. } L(\alpha, \lambda) = \frac{1}{2}\alpha^T S \alpha + \lambda^T (A\alpha - b)$$

$$\alpha \in \mathbb{R}^n, \lambda \in \mathbb{R}^m \quad \text{minus is also OK.}$$

$$\frac{\partial L}{\partial \alpha} = S\alpha - A^T \lambda = 0, \quad \frac{\partial L}{\partial \lambda} = A\alpha - b = 0.$$

$$\Rightarrow \begin{bmatrix} S & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} \alpha \\ \lambda \end{bmatrix} = \begin{bmatrix} 0 \\ b \end{bmatrix}$$

Not positive definite. The solution $\begin{bmatrix} \alpha \\ \lambda \end{bmatrix}$ is a saddle point of $L(\alpha, \lambda)$

• KKT matrix / conditions

Saddle points from Rayleigh quotient. $R(\alpha) = \frac{\alpha^T S \alpha}{\alpha^T \alpha}$

$$\lambda_{\max} = \max_{\alpha \in V} R(\alpha), \text{ at eigenvector } v_k \quad (Sv_k = \lambda_{\max} v_k)$$

$$\lambda_{\min} = \min_{\alpha \in V} R(\alpha), \text{ at eigenvector } y_k \quad (Sy_k = \lambda_{\min} y_k)$$

$$\therefore \lambda_{\min} \leq R(\alpha) \leq \lambda_{\max}$$

At $\alpha = v_k$ (k -th eigenvector), saddle point.

$$R(v_k) = \lambda_k, \quad \text{grad } R(v_k) = 0$$

$$\text{* Min-Max Principle: } \lambda_k = \max_{V \subseteq \mathbb{R}^n} \min_{\alpha \in V} \frac{\alpha^T S \alpha}{\alpha^T \alpha}$$

• Read more about this.

20. Definitions & Inequalities

Expectation $\mu = E[X]$

Variance $\sigma^2 = E[(X-\mu)^2]$

$$E[f(X)] = \sum_i p_i \cdot f(x_i)$$

$$\Rightarrow \sigma^2 = \sum_i p_i (x_i - \mu)^2 = \sum_i p_i x_i^2 - \mu^2 = E[X^2] - E[X]^2$$

* Markov Inequality.

Holds for nonnegative random variable X .

$$\text{For } a > 0, P(X \geq a) \leq \frac{E[X]}{a}$$

$$\begin{aligned} \text{Proof). } E[X] &= \int_0^\infty x f(x) dx = \int_0^a x f(x) dx + \int_a^\infty x f(x) dx \\ &\geq \int_a^\infty a f(x) dx \geq \int_a^\infty a f(a) dx = a P(X \geq a) \end{aligned}$$

* Chebyshew Inequality

$$\text{For random variable } X, P(|X-\mu| \geq a) \leq \frac{\text{Var}(X)}{a^2}$$

Proof). From Markov,

$$P(|X-\mu| \geq a) = P((X-\mu)^2 \geq a^2) \leq \frac{E[(X-\mu)^2]}{a^2} = \frac{\text{Var}(X)}{a^2}$$

* Covariance Matrix.

m experiments at once. Random variables X_i ($i=1, \dots, m$)

$$V = (v_{ij})_{m \times m} \text{ where } v_{ij} = \text{Cov}(X_i, X_j)$$

$$\text{And } \text{Cov}(X_i, X_j) = E[(X_i - \mu)(X_j - \mu)]$$

Alternatively, let $X = (X_1, X_2, \dots, X_m)$.

Then $\bar{X} = (\bar{X}_1, \bar{X}_2, \dots, \bar{X}_m)$, and

$$\begin{aligned} V &= E[(X - \bar{X})(X - \bar{X})^T] \text{ (always P.S.D.)} \\ &= E[XX^T] - \bar{X}\bar{X}^T \end{aligned}$$

21. Minimizing a Function Step by Step.

Taylor Expansion. (Single variable)

$$f(x + \Delta x) \approx f(x) + f'(x) \Delta x + \frac{1}{2} f''(x) \Delta x^2$$

Taylor Expansion (Multi-variable)

$$F(x + \Delta x) = F(x) + (\Delta x)^T \nabla F(x) + \frac{1}{2} (\Delta x)^T H \Delta x$$

where H is the Hessian matrix.

$$\text{Hessian matrix } (i, j) \text{- component : } \frac{\partial^2 F}{\partial x_i \partial x_j}$$

and also H is symmetric.

Let $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$, defined by

$$f(x_1, \dots, x_n) = f(\vec{x}) = (f_1(\vec{x}), \dots, f_n(\vec{x}))$$

where $f_i: \mathbb{R}^n \rightarrow \mathbb{R}$.

$$f(x + \Delta x) = f(x) + J \Delta x$$

where J is the Jacobian, which has $\frac{\partial f_i}{\partial x_j}$ as its (i, j) -component.

How would we solve $f = \vec{0}$ $\left\{ \begin{array}{l} n \text{ equations} \\ n \text{ unknowns} \end{array} \right.$

* Newton's Method.

$$x_{i+1} = x_i - J(x_i)^{-1} f(x_i) \quad \text{for } i = 0, 1, \dots$$

Can be used to solve $f(\vec{x}) = 0$.

* Minimizing Functions.

Minimizing $F(x) = 0$ is similar (!) to solving $\nabla F = 0$.

① Steepest Descent

$$x_{i+1} = x_i - s_i \nabla F \quad (\text{linear convergence})$$

(Gradient is the steepest direction)

② Newton's Method. (Solve $\nabla F = 0$)

$$x_{i+1} = x_i - H^{-1} \nabla F \quad (\text{quadratic convergence})$$

* Convexity (for convex optimization)

A set K is a **convex set** iff for any $x, y \in K$ and any $t \in [0, 1]$, $(1-t)x + ty \in K$.

For two convex sets A, B , $A \cup B$ is not convex, but $A \cap B$ is always convex.

Convex sets are usually the domain where we minimize the objective function.

A function f is **convex** iff for any $x, y \in \text{dom}(f)$ and any $t \in [0, 1]$, $f((1-t)x + ty) \leq (1-t)f(x) + tf(y)$

For two convex functions f, g , $\min(f, g)$ is not convex.

but $\max(f, g)$ is always convex.

Using derivatives,

$f: \mathbb{R} \rightarrow \mathbb{R}$, $f''(x) \geq 0$ then convex.

$f: \mathbb{R}^n \rightarrow \mathbb{R}$, $H \geq 0$ then convex. (P.S.D.)

22. Gradient Descent: Downhill to Minimum

Gradient Descent.

$$x_{i+1} = x_i - \alpha \nabla f(x_i)$$

α : learning rate. $\left\langle \begin{array}{l} \text{If } \alpha \text{ is too big} \rightarrow \text{oscillation.} \\ \text{If } \alpha \text{ is too small} \rightarrow \text{long training.} \end{array} \right.$

Must choose the optimal α (step size) : exact line search.

Backtracking line search: use $\alpha, \frac{1}{2}\alpha, \frac{1}{4}\alpha, \dots$

23. Accelerating Gradient Descent

With momentum, $\frac{1-b}{1+b} \rightarrow \frac{1-\sqrt{b}}{1+\sqrt{b}}$ (rate of decay)

Momentum

$$\left\{ \begin{array}{l} x_{i+1} = x_i - \alpha z_i \\ z_i = \nabla f(x_i) + \beta z_{i-1} \quad (\approx \text{memory term}) \end{array} \right.$$

Similar to second order differential equations

\rightarrow damping factors appear.

Try writing it with two first-order equations.

We try to minimize $f(\alpha) = \frac{1}{2} \alpha^T S \alpha \rightarrow \nabla f = S \alpha$

$$\begin{cases} \alpha_{k+1} = \alpha_k - d \alpha_k \\ z_{k+1} - S \alpha_{k+1} = \beta z_k \end{cases}$$

Let $S q = \lambda q$ (eigen) we track eigenvectors.

$\alpha_k = C_k q$, $z_k = d_k q$. Then $S \alpha_{k+1} = \lambda C_{k+1} q$

$$\Rightarrow \begin{bmatrix} 1 & 0 \\ -\lambda & 1 \end{bmatrix} \begin{bmatrix} C_{k+1} \\ d_{k+1} \end{bmatrix} = \begin{bmatrix} 1 & -d \\ 0 & \beta \end{bmatrix} \begin{bmatrix} C_k \\ d_k \end{bmatrix}$$

$$\begin{aligned} \therefore \begin{bmatrix} C_{k+1} \\ d_{k+1} \end{bmatrix} &= \begin{bmatrix} 1 & 0 \\ \lambda & 1 \end{bmatrix} \begin{bmatrix} 1 & -d \\ 0 & \beta \end{bmatrix} \begin{bmatrix} C_k \\ d_k \end{bmatrix} \\ &= \begin{bmatrix} 1 & -d \\ \lambda & \beta - \lambda d \end{bmatrix} \begin{bmatrix} C_k \\ d_k \end{bmatrix} \end{aligned}$$

Now we want to choose R to be as small as possible.

Suppose we know that $0 < m \leq \lambda \leq M$ (eigenvalues of S)

The condition number $K = \frac{M}{m} = \frac{\lambda_{\max}}{\lambda_{\min}}$

If the condition number is large, the problem is harder

Suppose R has eigenvalues e_1, e_2 . they depend on α and β . So we want to minimize the largest e_i , for all possible λ of S . ($0 < m \leq \lambda \leq M$)

$$\alpha_{\text{optimal}} = \left(\frac{2}{\sqrt{M} + \sqrt{m}} \right)^2 \quad \beta_{\text{optimal}} = \left(\frac{\sqrt{M} - \sqrt{m}}{\sqrt{M} + \sqrt{m}} \right)^2$$

So what ... ?

Nesterov Accelerated Gradient (NAG)

$$x_{i+1} = x_i + \beta(x_i - x_{i-1}) - \alpha \nabla f(x_i + \gamma(x_i - x_{i-1}))$$

Evaluate it somewhere on a line ↗

24. Linear Programming and Two-Person Games

Linear Programming

Want to minimize $c^T \alpha = c_1 \alpha_1 + \dots + c_m \alpha_m$.

Constraint on α . $A \alpha = b$, $A \in M_{m \times n}(\mathbb{R})$, $m < n$

and all $\alpha_i \geq 0$. "feasible set"

① Simplex Method. (Dantzig)

Start from 1 corner → next corner → ... that lowers the cost

② Karmarkar's Method. "interior"

Dual LP : maximum of $b^T y = b_1 y_1 + \dots + b_m y_m$

constraint: $A^T y \leq c$

Weak duality for any feasible α, y , $b^T y \leq c^T \alpha$.

Proof. $b^T y = \alpha^T A^T y \leq \alpha^T c = c^T \alpha$

$\alpha \geq 0$ keeps the direction of the inequality

Strong duality $b^T y^* = c^T \alpha^*$ y^* maximizes $b^T y$
 α^* minimizes $c^T \alpha$

Max flow - min cut

For a graph with each edge having a capacity,

what is the maximum flow from the source to sink?

Maximum flow is equal to the min cut.

A cut partitions the vertices into two sets. A and B.

the value of a cut is computed by the sum of all edges $e = (s, t)$ such that $s \in A$ and $t \in B$.

(its ends are in different sets)

2b. Structure of Neural Nets for Deep Learning

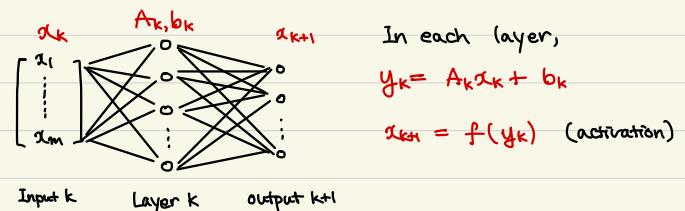
We have training data = feature vectors and the answer

$F: \alpha = (x_1, \dots, x_m) \rightarrow$ correct answer

We want to create a function $F(\alpha)$ that is correct on almost all of the data (don't want overfitting)

Also we have activation functions (usually nonlinear).

$$\text{ReLU}(x) = \max(0, x)$$



$r(N, m)$: Number of flat pieces in m dimensions with N folds

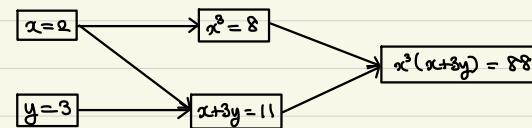
$$r(N, m) = \sum_{r=0}^m \binom{N}{r}$$

$$r(N, m) = r(N-1, m) + r(N-1, m-1)$$

27. Backpropagation

$$F(x, y) = x^3(x+3y)$$

Computational Graph of F.



Forward (not recommended)

$$\begin{aligned} \frac{\partial x^3}{\partial x} &= 3x^2 & \frac{\partial}{\partial x} (x^3(x+3y)) \\ x &\quad \frac{\partial x^3}{\partial y} = 0 &= 3x^2(x+3y) + x^3 = 4x^3 + 3x^2y \\ &&= 140 \\ y &\quad \frac{\partial (x+3y)}{\partial x} = 1 & \frac{\partial}{\partial y} (x^3(x+3y)) \\ &\quad \frac{\partial (x+3y)}{\partial y} = 3 &= 0 + 3x^3 = 3x^3 = 24 \end{aligned}$$

$$\begin{array}{c}
 \text{Simplifying: } \\
 \begin{aligned}
 x &= 2 & c &= x^3 = 8 \\
 y &= 3 & s &= x+3y = 8+3=11 \\
 && \frac{\partial F}{\partial x} &= 3x^2 = 3 \cdot 8 = 24 \\
 && \frac{\partial F}{\partial y} &= 1 \\
 && F &= x^3(x+3y) = 8 \cdot 11 = 88
 \end{aligned}
 \end{array}$$

31. Eigenvectors of Circulants. Fourier Matrix

But note that the eigenvalues of P are roots of unity!

Because $P^n = I$, characteristic polynomial is $t^n - 1$.

\therefore eigenvalues are ω^i for $i=0, 1, \dots, n-1$, $\omega = e^{\frac{2\pi i}{n}}$

Fourier matrix is the eigenvector matrix of P .

$$F_n = \frac{1}{\sqrt{n}} \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 1 & \omega & \cdots & \omega^{n-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \omega^{n-1} & \cdots & \omega^{(n-1)^2} \end{bmatrix}$$

(i,j)- component : $(\omega^{j-1})^{i-1}$

32. Convolution Rule

Convolution of functions $f(t), g(t)$. $h(t) = (f * g)(t)$

$$h(t) = \int_0^t f(u) g(t-u) du$$

Convolution Rule

Eigenvalues of P are $\omega^0, \omega^1, \dots, \omega^{n-1}$ ($\omega = e^{\frac{2\pi i}{n}}$)

thus eigenvalues of $C = c_0I + c_1P + \dots + c_{n-1}P^{n-1}$

should be $f(\omega^0), \dots, f(\omega^{n-1})$ where

$$f(t) = c_0 + c_1 t + \dots + c_{n-1} t^{n-1}$$

Since the i -th row of F is

$$[F]_i = (1, \omega^{i-1}, \omega^{2(i-1)}, \dots, \omega^{(n-1)(i-1)})$$

$$\begin{aligned}
 [F]_i \cdot (c_0, c_1, \dots, c_{n-1})^T &= c_0 + c_1 \omega^{i-1} + \dots + c_{n-1} \omega^{(n-1)(i-1)} \\
 &= f(\omega^{i-1})
 \end{aligned}$$

Thus $F \cdot (c_0, \dots, c_{n-1})^T$ gives the eigenvalues of C

Eigenvalues of $CD = (\text{Eigenvalues of } C) \cdot (\text{Eigenvalues of } D)$

$$F(c \otimes d) = F(c) \odot F(d) \quad (\text{pointwise multiplication})$$

→ By FFT, this computation is fast.

* Kronecker product : 2D matrix convolution.

33. Neural Nets and the Learning Function

Learning function $F(x, v)$

x : weights in the neural network. A_k, b_k ,

v_i : sample feature vectors.

Optimize! using SGD etc

Each layer: $v_{k+1} = \text{ReLU}(A_k v_k + b_k)$ $k=0, 1, \dots, L$ (layers)

(No activation function for the last layer)

Often, weights are underdetermined because $\#x \gg \#v$

30. Completing a Rank 1 Matrix, Circulants

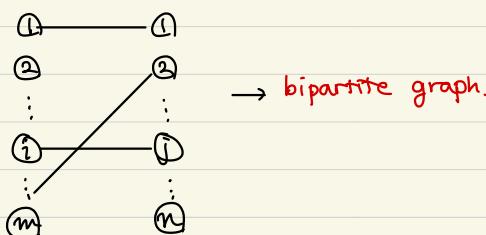
$$A = uv^T \text{ (mxn matrix). } a_{ij} = u_i v_j$$

We are given $m+n-1$ non-zero numbers (with position)

Can we complete the rank 1 matrix A ?

For which positions, is the problem solvable?

If element (i,j) is given, add an edge



If such bipartite graph has a cycle, the problem will be unsolvable.

Circulants (cyclic convolution matrix)

diagonals circle around. and each diagonal is constant.

$$C = \begin{bmatrix} 2 & 5 & 1 & 0 \\ 0 & 2 & 5 & 1 \\ 1 & 0 & 2 & 5 \\ 5 & 1 & 0 & 2 \end{bmatrix}$$

For two circulant matrices C, D .

① CD is also circulant.

$$\textcircled{2} \quad C = c_0I + c_1P + c_2P^2 + \dots \quad (\text{polynomial in } P)$$

where P is a permutation matrix of $\sigma = (1 \ 2 \ 3 \ \dots \ n)$

Note that $P^n = I$.

For $v = (v_1, \dots, v_n)$, $w = (w_1, \dots, w_n)$, convolution $v * w$

$$v * w = (c_1, c_2, \dots, c_{2n-1}) \text{ where}$$

$$c_k = \sum_{i=1}^k v_i \cdot w_{k+1-i} \quad (k=1, \dots, 2n-1)$$

A cyclic convolution $v \otimes w$ gives

$$v \otimes w = (d_1, d_2, \dots, d_n) \text{ where}$$

$$d_k = c_k + c_{k+n} \quad (c_{2n} = 0 \text{ for convenience})$$

Eigenvalues, Eigenvectors of Circulants

Eigenvectors of the permutation matrix P are also eigenvectors of a circulant.

If $C = f(P)$, and λ is an eigenvalue of P ,

$f(\lambda)$ is an eigenvalue of C

Loss Function.

Choose α to minimize L (loss function)

$$L = \frac{1}{N} \left[\sum_{i=1}^N F(\alpha_i, x_i) - \text{true}_{i,i} \right]$$

\downarrow Output of α_i \downarrow True answer of x_i

Popular Loss Functions

- ① Square loss (error is squared)
- ② L^1 loss (LASSO)
- ③ Hinge loss (-1.1 classification, regression)
- ④ Cross-Entropy loss

In SGD, we only use part of the data to calculate the loss and perform backpropagation

Distance Matrices

Given distances squared $\|x_i - x_j\|^2 = d_{ij}$ as matrix D

find positions $x_i \in \mathbb{R}^n$

$$d_{ij} = \|x_i - x_j\|^2 = \|x_i\|^2 + \|x_j\|^2 - 2x_i^T x_j$$

rank 1 matrix rank 1 matrix
cols repeated rows repeated

Define $G = X^T X$ and find X from G . (of course, X is not unique)

Also, let $d = (d_1, \dots, d_n)^T$ and $d = \|x\|^2$ \leftarrow How do we know this?

$$D = (d_{ij}) \quad \text{and} \quad d_{ij} = x_i^T x_j$$

From $D = 1d^T + d1^T - 2X^T X$, we see that

$$X^T X = G = -\frac{1}{2}D + \frac{1}{2}1d^T + \frac{1}{2}d1^T$$

Now, how to find X from $X^T X$?

$X^T X$ is positive semi-definite!

① Eigenvalues of $X^T X = Q \Lambda Q^T$ (by Spectral Thm)

then $X = Q \sqrt{\Lambda} Q^T$.

② Elimination on $X^T X \rightarrow X^T X = LDU$

but since $X^T X$ is symmetric, $L^T = U$.

Then $X = \sqrt{D} L^T$ (Cholesky) (Much faster)

34. Procrustes Problem

Given two sets of vectors $X = [x_1 \dots x_m]$, $Y = [y_1 \dots y_n]$

Find $\underset{Q}{\operatorname{argmin}} \|YQ - X\|_F^2$ subject to $Q^T Q = I$

Closest orthogonal transformation.

Trace is sum of eigenvalues

$$\text{Note that } \|A\|_F^2 = \operatorname{tr}(A^T A) = \operatorname{tr}(AA^T) = \sum \sigma_i^2$$

$$Y^T X = U \Sigma V^T \text{ then, best } Q = U V^T$$

Proof. $\underset{Q}{\operatorname{argmin}} \|YQ - X\|_F^2$

$$= \underset{Q}{\operatorname{argmin}} \operatorname{tr}((YQ - X)^T (YQ - X))$$

$$= \underset{Q}{\operatorname{argmin}} \operatorname{tr}(Q^T Y^T Y Q - X^T Y Q - Q^T Y^T X + X^T X)$$

$$= \underset{Q}{\operatorname{argmin}} (\|QY\|_F^2 + \|X\|_F^2 - 2\operatorname{tr}(Q^T Y^T X))$$

$$= \underset{Q}{\operatorname{argmin}} (\|X\|_F^2 + \|Y\|_F^2 - 2\operatorname{tr}(Q^T Y^T X))$$

$$= \underset{Q}{\operatorname{argmax}} \operatorname{tr}(Q^T Y^T X)$$

$$= \underset{Q}{\operatorname{argmax}} \operatorname{tr}(Q^T U \Sigma V^T) = \underset{Q}{\operatorname{argmax}} \operatorname{tr}(Q^T U \Sigma)$$

orthogonal

Since $Q^T U \Sigma V^T$ is orthogonal, $\operatorname{tr}(Q^T U \Sigma V^T)$ is maximized

when $I = Q^T U \Sigma V^T \therefore Q = U \Sigma V^T$

Each column of $Q^T U \Sigma V^T$ will have norm 1. when multiplied by Z , the columns are stretched.
To maximize the trace, the i -th component of the i -th column should be 1. Thus $Q^T U \Sigma V^T = I$.

35. Finding Clusters in Graphs

For partition A, B of X , find positions x, y such that

$$\text{minimizes : } \sum_{i \in A} \|x - x_i\|^2 + \sum_{i \in B} \|y - y_i\|^2$$

Given a_1, \dots, a_n , the best x that minimizes $\sum_i \|a_i - x\|^2$

is the **centroid** of a_i 's. $x = \frac{1}{n} \sum a_i$

K-means (here, $k=2$)

① Given A, B (clusters), find centroid x, y .

② Given x, y , form best clusters A, B .

- Each node (data) goes with the closer one of x, y .

Repeat! (until clusters are same)

Solution Method (spectral clustering)

"Spectrum" of a matrix \rightarrow Eigenvalues.

Start with graph **Laplacian matrix**

$$L = A^T A = D - G$$

Laplace's finite difference equation
connects with $u_{xx} + u_{yy} = 0$

size: $|E| \times |V|$
 A : incidence matrix. If $e_i = (u, v)$, $a_{iu} = -1$, $a_{iv} = 1$.

$|V| \times |V|$
 D : degree matrix. Diagonal, $d_{ii} = \deg(v_i)$

$|V| \times |V|$
 G : adjacency matrix. $g_{ij} = 1$ if $(i, j) \in E$, 0 otherwise.

L is PSD. and $\lambda_1 = 0$ with eigenvector $v_1 = c_1 \mathbf{1}$

Fiedler vector v_2 : eigenvector of smallest $\lambda > 0$

The cluster for data i is decided by the sign of the i -th component of v_2 .

Also, since L is PSD. $v_1 \perp v_2$, sum of components of v_2 is 0.

36. Alan Edelman and Julia Language.