

EE499 Cheat Sheet - Midterm 1

Vector space: A vector space V is a set of mathematical objects along with a field of scalars F that is closed under vector addition & scalar multiplication:

$$\forall x, y \in V, x + y \in V$$

$$\forall x \in V, \forall \alpha \in F, \alpha x \in V$$

Here these operations satisfy

- Additive commutativity
- Additive associativity
- Existence of additive identity
- Scalar multiplication associativity
- Distribution laws
- $0x = \underline{0}$ & $1x = x$ (which implies additive inverses)

Subspace: A subset of a vector space that is itself a vector space is called a subspace.

- The same operations are concerned as the parent vector space.
- A subspace must always contain the null vector.

Operations on subspaces: Let S & T be subspaces of a vector space V .

- $S \cap T$ is also a subspace
- $S \cup T$ is not necessarily a subspace

Span: For some vectors $x_i \in V$, $\text{span}(\{x_i\})$ is the set of all linear combinations of x_i .

Algebraic Sum: Let $S, T \subset V$. Then the algebraic sum of S & T is

$$S + T = \{s + t \mid s \in S \text{ & } t \in T\}$$

Direct Sum: The direct sum of two subspaces S, T of V is their algebraic sum if and only if S & T are disjoint subspaces, i.e. $S \cap T = \{\underline{0}\}$.

Linear combination: Consider a vector space V and its subset $S \subset V$. A vector $x \in V$ is a linear combination of elements of S if there exists a finite number of elements $\{s_1, \dots, s_n\} \subset S$ and $\{\alpha_1, \dots, \alpha_n\} \subset F$ such that

$$\alpha_1 s_1 + \dots + \alpha_n s_n = x$$

Linear Independence: x_1, \dots, x_n are linearly independent vectors if

$$\alpha_1 x_1 + \dots + \alpha_n x_n = 0 \iff \forall k \alpha_k = 0$$

Linearly independent vectors enable unique representation of other vectors, i.e.

$$\sum_{k=1}^n \alpha_k x_k = \sum_{k=1}^n \beta_k x_k \text{ if and only if } \forall k \alpha_k = \beta_k$$

Basis: A first definition: For a vector space V , the set $\{x_1, \dots, x_n\}$ are called a basis for V if

1. x_1, \dots, x_n are linearly independent
2. $\text{Span}(x_1, \dots, x_n) = V$.

A more general definition: The set of vectors $\Phi = \{\varphi_k\}_{k \in K} \subset V$ where K is either finite or countably infinite is called a basis for the normed vector space V when

1. $\forall x \in V \quad x = \sum_{k \in K} \alpha_k \varphi_k, \alpha_k \in \mathbb{R} \implies$ convergence dependent on the norm
2. Any sequence α_k that represents some $x \in V$ is unique.

We are mainly looking for existence & uniqueness of a representation. Linear independence of the elements is equivalent to the uniqueness of the representation.

Change of Basis: Say you have two different bases, e & v , for the same vector space. $[x]_e$ and $[x]_v$ are then the representations of some vector x in those bases. Then, one can pass between these representations through a simple matrix multiplication:

$$e = \{e_1, e_2, \dots, e_n\}$$

$$v = \{v_1, v_2, \dots, v_n\}$$

$$[x]_v = T_{e \rightarrow v} [x]_e \text{ where}$$

$$T_{e \rightarrow v} = \begin{bmatrix} [e_1]_v & [e_2]_v & \dots & [e_n]_v \end{bmatrix}$$

$$T_{v \rightarrow e} = (T_{e \rightarrow v})^{-1} = \begin{bmatrix} [v_1]_e & [v_2]_e & \dots & [v_n]_e \end{bmatrix}$$

Functional: A mapping T from a vector space V to the real numbers \mathbb{R} or the complex numbers \mathbb{C} .

$$T: V \rightarrow \mathbb{R}, \mathbb{C}$$

i.e. onto the field.

Norm: A functional is a norm if it satisfies the following axioms:

- Non-negativity: $\|x\| \geq 0 \quad \forall x \in V, \|x\| = 0 \Leftrightarrow x = \underline{0}$
- Homogeneity: $\|\alpha x\| = |\alpha| \cdot \|x\| \quad \forall x \in V, \forall \alpha \in F$
- Triangle inequality: $\|x+y\| \leq \|x\| + \|y\| \quad \forall x, y \in V$

Distance: A norm induces a metric $d(x, y) = \|x - y\|$.

Normed Vector Space: A normed vector space is a pair $(V, \|\cdot\|)$ where V is a vector space and $\|\cdot\|$ is a norm defined on V .

Some typical norms:

- $l_p, p \geq 1: (\sum |x_i|^p)^{1/p}, l_\infty: \sup |x_i|$
- $L_p[a, b], p \geq 1: (\int_a^b |x|^p dx)^{1/p}, L_\infty[a, b]: \max x(t)$
- On $\mathbb{C}^{N \times N}, \|A\| = \max |A_{ij}|$
- On $\mathbb{C}^{N \times N}, \|A\|_F = (\sum_i \sum_j |A_{ij}|^2)^{1/2} = \sqrt{\text{trace}(A^* A)}$

For l_p norms, we have

- $(\text{Unit ball})_p \subset (\text{Unit ball})_q$ if $p < q$
- $\|x\|_p \geq \|x\|_q$ if $p < q$
- $l_p \subset l_q$ if $p < q$

For function spaces, the opposite of the last item is true.

$$l_1 \subset \dots \subset l_\infty$$

$$L_1[a, b] \supset \dots \supset L_\infty[a, b]$$

Inner product: Let V be a vector space, with its field F . An inner product is a mapping $\langle \cdot, \cdot \rangle: V \times V \rightarrow F$ with

1. $\langle x, y \rangle = \langle y, x \rangle^* \quad \forall x, y \in V$
2. $\langle x+y, z \rangle = \langle x, z \rangle + \langle y, z \rangle \quad \forall x, y, z \in V$
3. $\langle \alpha x, z \rangle = \alpha \langle x, z \rangle$
4. $\langle x, x \rangle \geq 0 \quad \forall x \in V \text{ \& } x \neq \underline{0}, \langle x, x \rangle = 0 \text{ iff } x = \underline{0}$

Inner product space: An inner product space is a vector space equipped with an inner product.

Induced norm: An inner product induces a unique norm as

$$\|x\| = \sqrt{\langle x, x \rangle}$$

Cauchy-Schwarz Inequality: For all vectors $x, y \in V$,

$$|\langle x, y \rangle| \leq \|x\| \cdot \|y\|$$

Angle: Using the Cauchy-Schwarz Inequality, we can define the concept of an angle between vectors:

$$\cos \theta = \frac{\langle x, y \rangle}{\|x\| \cdot \|y\|} \quad \text{for all } x, y \in V.$$

Various angle definitions exist for complex-valued inner products.

Parallelogram law: If V is an inner product space and $\|\cdot\|$ is the induced norm on V , then

$$\|x+y\|^2 + \|x-y\|^2 = 2\|x\|^2 + 2\|y\|^2 \quad \forall x, y \in V$$

If the norm of a normed space satisfies the parallelogram law, then it is induced by an inner product.

$\Rightarrow \ell_p$ or l_p spaces are IPS's iff $p=2$.

Orthogonality: Two vectors $x, y \in V$ are said to be orthogonal if

$$\langle x, y \rangle = 0$$

which is denoted by $x \perp y$.

Similarly, x & y are called parallel if $x = \alpha y$ for some $\alpha \in \mathbb{F}$.

Orthogonality of two sets is defined by the orthogonality of their elements:

$$S \perp T \iff s \perp t \quad \forall s \in S, \forall t \in T.$$

Orthogonal complement: The orthogonal complement of a set $S \subset V$ is defined as

$$S^\perp = \{x \in V \mid x \perp s \quad \forall s \in S\}$$

Orthonormal vectors: A set of vectors $\{x_1, \dots, x_n\}$ is orthonormal if

- $x_i \perp x_j$ for all $i \neq j$
- $\|x_i\| = 1$ for all i .

Gram - Schmidt Orthonormalisation Process: let x_1, \dots, x_n be a set of linearly independent vectors. Then one can find an orthogonal sequence inductively as follows:

$$e_1 = x_1$$

$$e_i = x_i - \sum_{j=1}^{i-1} \langle x_i, e_j \rangle e_j \quad \text{for all } 1 < j \leq n$$

If one divides each of these vectors by their norms, one obtains an orthonormal set.

$$u_i = \frac{e_i}{\|e_i\|} \quad \text{for all } 1 \leq i \leq n$$

The benefit of an orthonormal basis is that representation coefficients can be found by a simple inner product operation:

$$x = \sum_{k=1}^n \beta_k u_k, \quad \rightarrow \quad \beta_k = \langle x, u_k \rangle$$

and thus

$$x = \sum_{k=1}^n \langle x, u_k \rangle u_k$$

Operator: An operator is a mapping between vector spaces.

Linear operator: An operator $A: X \rightarrow Y$ is linear if

$$A(\alpha x + \beta y) = \alpha A(x) + \beta A(y)$$

or equivalently the two following conditions:

- Homogeneity: $A(\alpha x) = \alpha A(x)$
- Additivity: $A(x+y) = A(x) + A(y)$

Range space: The range space of an operator $A: X \rightarrow Y$ is defined as

$$R(A) = \{ y \in Y \mid \exists x \in X \ A(x) = y \} \subset Y$$

If $R(A) = Y$ for some operator A , then A is said to be a surjection, or onto/surjective.

Null space: The null space of an operator $A: X \rightarrow Y$ is defined as

$$N(A) = \{ x \in X \mid A(x) = 0 \} \subset X$$

If $N(A) = \{0\}$ for some operator A , then A is said to be an injection, or one-to-one/injective.

Both $N(A)$ & $R(A)$ are subspaces.

Inverse problem: Find x , given the operator A & $y = A(x)$.

The inverse problem has at least one solution if A is surjective.

It has at most one solution if A is injective.

It has a unique solution if A is surjective and injective, i.e. bijective.

Linear operators in finite dimensional spaces: An operator A between finite dimensional spaces always has a matrix representation.

Say e is a basis for X , v a basis for Y , and $A: X \rightarrow Y$ a linear operator. Then

$$[A(x)]_v = \underline{A} [x]_e$$

where

$$\underline{A} = \begin{bmatrix} [A(e_1)]_v & [A(e_2)]_v & \dots & [A(e_n)]_v \end{bmatrix}$$

If we were to change the basis, from e to e' and v to v' , then we can have a new matrix representation \underline{A}' as

$$\underline{A}' e' \rightarrow v' = T_{v \rightarrow v'} \underline{A} e \rightarrow v T_{e' \rightarrow e}$$

so that

$$[A(x)]_{v'} = \underline{A}' [x]_{e'}$$

Matrix Rank: The rank for a matrix $A \in \mathbb{C}^{M \times N}$ is the number of linearly independent columns of A .

We always have

$$\text{rank}(A) \leq \min\{M, N\}$$

Rank is also defined as the dimensionality of the range space of some matrix A .

Null-space Dimensionality: For a matrix $A \in \mathbb{C}^{M \times N}$

$$\dim(N(A)) = N - \text{rank}(A)$$

Rank-Nullity Theorem: Let $A \in \mathbb{C}^{M \times N}$

$$\dim(X) = \text{rank}(A) + \dim(N(A))$$

Projectors: $\hat{x} \in S \subset V$ is the projection of some $x \in V$ if

$$\hat{x} = \operatorname{argmin}_{s \in S} \|x - s\|^2$$

Orthonormal bases to Projections: let b_1, \dots, b_n be an orthonormal basis for a subspace $S \subset V$. Then the projection matrix \underline{P} is given by

$$\underline{P} = \underline{B} \underline{B}^H$$

where

$$\underline{B} = [b_1 \quad b_2 \quad \dots \quad b_n]$$

If b_i 's are a basis for S , although not orthonormal, we have

$$\underline{P} = \underline{B} (\underline{B}^H \underline{B})^{-1} \underline{B}^H$$

Additional content:

- Hölder's Inequality: let $p, q \in (1, \infty)$, $1/p + 1/q = 1$
$$\int |f(t)g(t)| dt = \left(\int |f(t)|^p dt \right)^{1/p} \left(\int |g(t)|^q dt \right)^{1/q}$$
- A matrix (or a transform) U is called unitary if the columns of U form an orthonormal basis. In this case,

$$U^H U = I$$

A unitary transformation/matrix preserves the inner product, and hence the norm:

$$\langle Ux, Uy \rangle = \langle x, y \rangle \quad \forall x, y \in V$$

$$\|Ux\| = \|x\|$$

EE499 Cheat Sheet - Midterm 2

Projectors in \mathbb{C}^N : Let $x \in \mathbb{C}^N$, $S \subseteq \mathbb{C}^N$, S a subspace of \mathbb{C}^N . Then we call \hat{x} a projection of x onto S if

$$\hat{x} = \underset{s \in S}{\operatorname{argmin}} \|x - s\|^2$$

so \hat{x} is the closest vector in S to x .

If b_1, \dots, b_r is a basis (hence linearly independent) for S , we can represent the projection operation as a matrix multiplication:

$$\hat{x} = Px \quad \text{where}$$

$$P = B(B^H B)^{-1} B^H \quad \text{for } B = [b_1 \dots b_r]$$

If further the basis is orthonormal, then

$$P = \underbrace{B(B^H B)^{-1} B^H}_{I} = BB^H$$

Projection theorem in IPS: Let X be an IPS, S a subspace of X , $x \in X$. If there is a vector $\hat{x} \in S$ such that

$$\|x - \hat{x}\| \leq \|x - s\| \quad \forall s \in S$$

then such an \hat{x} is unique. A necessary & sufficient condition that $\hat{x} \in S$ be the unique minimizing vector in S is that the error vector $e = x - \hat{x}$ be orthogonal to S . This is called as the orthogonality principle.

In an arbitrary inner product space, a projection may not exist but if it does, it is unique & fully characterized by the orthogonality principle.

Completeness: A space X is said to be complete if every Cauchy sequence is convergent in X . A vector sequence $\{x_n\}_{n=1}^{\infty}$ in a normed vector space is Cauchy if

$$\lim_{m, n \rightarrow \infty} \|x_n - x_m\| = 0$$

Every convergent sequence is Cauchy.

A complete IPS is called a Hilbert space. All finite dimensional spaces, ℓ_2 & $L_2[a, b]$ are Hilbert spaces.

Projection theorem in Hilbert spaces: Let X be a Hilbert space. Then the projection operator P_S (or equivalently a projection \hat{x}) exists for a subspace S if and only if S is closed or equivalently complete.

Because \mathbb{C}^N and a subspace $S \subseteq \mathbb{C}^N$ are finite dimensional, we are sure that a projection always exists & is unique and complies with the orthogonality principle. Furthermore, knowing that a projector is a linear operator, it has a matrix representation, as derived earlier.

3 methods of computing projections:

1) Optimization based arguments: Define the cost function J as

$$J: S \rightarrow \mathbb{R}$$

$$s \rightarrow \|x - s\| \quad \text{or} \quad \|x - s\|^2$$

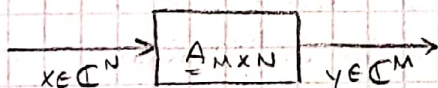
and minimize it by solving

$$\hat{x} = \underset{s \in S}{\operatorname{argmin}} J(s) \iff \nabla J(\hat{x}) = 0$$

2) Find \hat{x} using the orthogonality principle, which provides a system of equations.

3) Find the projection matrix P for S & compute $\hat{x} = Px$.

Solving linear systems of Equations:



Given y & A , find x such that $y = Ax$.

• Existence of a solution: For every $y \in \mathbb{C}^M$, there exists at least one solution iff

- $\iff R(A) = \mathbb{C}^M$
- $\iff \operatorname{rank}(A) = \dim(\underbrace{R(A)}_{\mathbb{C}^M}) = M$ (# of rows)
- \iff rows of A are linearly independent (A has full row rank)

One necessary condition is to have $M \leq N$, i.e. A be a square or fat matrix.

• Uniqueness of a solution: If there exists a solution, it is unique iff

- $\iff N(A) = \{0\}$
- $\iff \operatorname{rank}(A) = N$
- \iff columns of A are linearly independent (A has full column rank)

One necessary condition for uniqueness is to have $N \leq M$, i.e. A be a square or a tall matrix.

- Existence & Uniqueness: For every $y \in \mathbb{C}^M$, there exists a unique solution iff

C
A
S
E
O

$$\Leftrightarrow R(A) = \mathbb{C}^M \text{ \& } N(A) = \{0\}$$

$$\Leftrightarrow \text{rank}(A) = M = N$$

\Leftrightarrow rows of A are linearly independent & columns of A are linearly independent

$\Leftrightarrow A$ is invertible.

A necessary condition is to have $M=N$, so A be a square matrix.

- If conditions of both Case 1 & 2 are violated, existence & if so uniqueness are not guaranteed. (C A S E 3)

Basics of matrix calculus:

- $\nabla_x (a^T x) = \nabla_x (x^T a) = a$ for any vector or ^(square) matrix a

- $\nabla_x (x^T B x) = Bx + B^T x = (B + B^T)x$

Case 0: Least squares solution (derived for real valued case)

- $J(x) = \|y - Ax\|^2 = (y - Ax)^T (y - Ax)$

$$= y^T y - \underbrace{y^T A}_{a^T} x - x^T A^T y + x^T A^T A x$$

$$\nabla_x (J(x)) = -(y^T A)^T - A^T y + (A^T A + A^T A)x = 0$$

$$-\cancel{A^T y} + \cancel{A^T A} \hat{x}_{ls} = 0$$

$$\underbrace{A^T A}_{\text{invertible as } A \text{ has linearly independent cols (Case 0)}} x = A^T y \Rightarrow \hat{x}_{ls} = (A^T A)^{-1} A^T y \rightarrow \text{NORMAL EQUATION}$$

For complex-valued case:

$$A^H A \hat{x}_{ls} = A^H y \Rightarrow \hat{x}_{ls} = (A^H A)^{-1} A^H y$$

- One can derive the normal equation using projection based arguments as well.

Equivalently define

$$\hat{b}_{LS} = \underset{b \in R(A)}{\operatorname{argmin}} \|y - b\|^2$$

↳ looking for a projection of y onto $R(A)$

$$\hat{b}_{LS} = A \hat{x}_{LS}$$

By the orthogonality principle, $e \perp R(A) = \operatorname{span}\{a_1, \dots, a_n\}$
 $\hookrightarrow A = [a_1 \dots a_n]$

$$e \perp a_k \quad \forall k$$

$$\langle e, a_k \rangle = a_k^H e = 0 \quad \forall k$$

$$\begin{bmatrix} a_1^H e \\ a_2^H e \\ \vdots \\ a_n^H e \end{bmatrix} = \begin{bmatrix} a_1^H \\ a_2^H \\ \vdots \\ a_n^H \end{bmatrix} e = 0$$

$$\rightarrow A^H e = 0 \quad \leftarrow \text{insert } e = y - \hat{b}_{LS} = y - A \hat{x}_{LS}$$

$$A^H (y - A \hat{x}_{LS}) = 0$$

$$A^H A \hat{x}_{LS} = A^H y \Rightarrow \hat{x}_{LS} = (A^H A)^{-1} A^H y$$

$$\hat{b}_{LS} = \underbrace{A (A^H A)^{-1} A^H}_{\text{Projection matrix onto } R(A)} y$$

Projection matrix onto $R(A)$!

Case 2: At least one solution exists, the problem is to pick one. If x_p is a solution, elements of $x_p + N(A)$ are all solutions. So, we constrain our solution further. One possible choice is to have a minimum-norm solution.

$$\hat{x}_{MN} = A^H (A A^H)^{-1} y$$

Case 3: We deal with both issues. We can find a least squares solution through the normal equation, but that solution is not unique either, so we may enforce a minimum norm solution as well. So the problem becomes

$$\hat{x}_{MNLS} = \underset{\substack{x \text{ s.t.} \\ A^H A x = A^H y}}{\operatorname{argmin}} \|x\|_2$$

which turns out to have a closed form solution, which is

$$\hat{x}_{MNLS} = A^+ y$$

where A^+ denotes the pseudoinverse of A .

Spectral Decomposition Theorem: Let $A \in \mathbb{C}^{n \times n}$ be a Hermitian symmetric matrix. Then there exist a unitary matrix $U \in \mathbb{C}^{n \times n}$ and a real diagonal matrix Λ such that

$$A = U \Lambda U^H = [u_1 \dots u_n] \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} \begin{bmatrix} u_1^H \\ \vdots \\ u_n^H \end{bmatrix}$$

$$= \sum_{k=1}^n \lambda_k u_k u_k^H$$

λ_k 's are the eigenvalues of A & u_k 's are the orthonormal eigenvectors of A . We then conclude the following facts:

- $A u_k = \lambda_k u_k$
- $U^H U = I$
- $A U = U \Lambda$
- $\Lambda = U^H A U$
- $\text{rank}(A) = \#$ of non-zero eigenvalues.

The importance is that when multiplying with such a matrix, U^H makes a transition into some other basis of \mathbb{C}^n and the matrix multiplication in that domain becomes elementwise multiplication. We return back to the original basis after the elementwise multiplication using $U = (U^H)^{-1}$.

Singular Value Decomposition Theorem: Let $A \in \mathbb{C}^{m \times n}$ and $\text{rank}(A) = r$, $r \leq m$ & $r \leq n$. Then there exist unitary matrices $U \in \mathbb{C}^{m \times m}$ & $V \in \mathbb{C}^{n \times n}$ and a diagonal matrix $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r)$ with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ such that

$$A = U \Sigma V^H = [u_1 \dots u_r] \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{bmatrix} \begin{bmatrix} v_1^H \\ \vdots \\ v_r^H \end{bmatrix}$$

$$= \sum_{i=1}^r \sigma_i u_i v_i^H$$

Unitariness here for matrices of arbitrary size is

$$U^H U = I$$

$$V^H V = I$$

So we may indeed have $U U^H \neq I$ & $V V^H \neq I$.

Comparing SVD & Spectral Decomposition:

Spectral Decomposition

- $A \in \mathbb{C}^{n \times n}$, $A^H = A$
- $A = U \Lambda U^H$
 $= \sum_{k=1}^n \lambda_k u_k u_k^H$
- $U = [u_1 \dots u_n] \in \mathbb{C}^{n \times n}$
 \hookrightarrow eigenvectors of A

- $U^H U = I_n = U U^H$
- $A = U \Lambda U^H \Leftrightarrow \Lambda = U^H A U$
 \hookrightarrow Applying $A \Leftrightarrow$ elt. wise multiplication in U domain

Singular Value Decomposition

- $A \in \mathbb{C}^{m \times n}$, no restriction
- $A = U_{m \times r} \Sigma_{r \times r} (V_{n \times r})^H$
 $= \sum_{k=1}^r \sigma_k u_k v_k^H$
- $U = [u_1 \dots u_r] \in \mathbb{C}^{m \times r}$
 \hookrightarrow left singular vectors of A ,
eigenvectors of AA^H
- $V = [v_1 \dots v_r] \in \mathbb{C}^{n \times r}$
 \hookrightarrow right singular vectors of A ,
eigenvectors of $A^H A$.

- $U^H U = I$, possibly $U U^H \neq I$
 $V^H V = I$, possibly $V V^H \neq I$.

- $R(U) = R(A)$, $P_{R(A)} = U U^H$
 $R(V) = R(A^H)$, $P_{R(A^H)} = V V^H$

Extended SVD: In this formulation $A \in \mathbb{C}^{m \times n}$ is expanded such that 0 singular values are also allowed and the singular vector matrices are square and fit the regular unitary matrix property. So,

$$A = \tilde{U}_{m \times m} \tilde{\Sigma}_{m \times n} (\tilde{V}_{n \times n})^H$$

$$= \underbrace{\begin{bmatrix} U_{m \times r} & u_{r+1} & \dots & u_m \\ \hline 0 & & & \end{bmatrix}}_{\tilde{U}_{m \times m}} \underbrace{\begin{bmatrix} \Sigma_{r \times r} & 0 \\ 0 & 0 \end{bmatrix}}_{\tilde{\Sigma}_{m \times n}} \underbrace{\begin{bmatrix} V_{n \times r} & v_{r+1} & \dots & v_n \\ \hline 0 & & & \end{bmatrix}}_{\tilde{V}_{n \times n}}^H$$

$$= \sum_{k=1}^r \sigma_k u_k v_k^H$$

with

$$\tilde{U}^H \tilde{U} = \tilde{U} \tilde{U}^H = I_m$$

$$\tilde{V}^H \tilde{V} = \tilde{V} \tilde{V}^H = I_n$$

Pseudo-inverse (Moore-Penrose): let $A \in \mathbb{C}^{m \times n}$ and $\text{rank}(A) = r$. Then the pseudo inverse is constructed using SVD as

$$A = U \Sigma V^H \rightarrow A^+ = V \Sigma^{-1} U^H$$

where

$$\Sigma^{-1} = \begin{bmatrix} 1/\sigma_1 & & & 0 \\ & \ddots & & \\ & & \ddots & \\ 0 & & & 1/\sigma_r \\ & & & & & 0 \end{bmatrix}$$

As a closed sum, we get

$$A = \sum_{k=1}^r \sigma_k u_k v_k^H \rightarrow A^+ = \sum_{k=1}^r \frac{1}{\sigma_k} v_k u_k^H$$

A corollary of this definition is

$$A^+ = A^H (A A^H)^+ = (A^H A)^+ A^H$$

In case A is square and has an inverse, A^+ is an inverse to A . As multiplicative inverses are unique, we obtain

$$A \in \mathbb{C}^{n \times n} \text{ \& \; } \det(A) \neq 0 \rightarrow A^{-1} = A^+$$

Case-by-case pseudoinverse:

- Case 0: A is square and nonsingular

$$A^+ = A^{-1} \Rightarrow \hat{x} = x = A^+ y = A^{-1} y$$

- Case 2: A has full column rank

$$A^+ = (A^H A)^{-1} A^H = A_{LS} \Rightarrow \hat{x}_{LS} = A^+ y$$

- Case 2: A has full row rank

$$A^+ = A^H (A A^H)^{-1} = A_{MN} \Rightarrow \hat{x}_{MN} = A^+ y$$

- Case 3: A has neither full row rank nor full column rank

$$A^+ = V \Sigma^{-1} U^H = A_{MNLs} \Rightarrow \hat{x}_{MNLs} = A^+ y$$

\Rightarrow Therefore no matter in which case we are, one can use the pseudoinverse to find the solution to the inverse problem!

Sensitivity & Conditioning of Linear Inverse Problems: There may be errors on y and on A , due to

- physical measurement errors (noise)
- modelling errors
- numerical implementation / storage errors (finite precision)

How can we quantify the error in our solution to the inverse problem?

Define the 2-norm condition number as

$$K(A) = \frac{\sigma_{\max}}{\sigma_{\min}}$$

where σ_{\max} & σ_{\min} are the maximum & minimum singular values of A , respectively.

We can observe that $K(A) \geq 1$, and $K(A) = 1$ iff all the singular values of A are the same.

- If $K(A) \approx 1$, then A is a well-conditioned matrix
- If $K(A) \gg 1$, then A is an ill-conditioned matrix, meaning it may amplify the error upon inverse problem solution.

• Case 0: We have non-singular & square A , & $x = A^{-1}y$. Then

$$\frac{\|\delta x\|_2}{\|x\|_2} \leq K(A) \frac{\|\delta y\|_2}{\|y\|_2}$$

$$\frac{\|\delta x\|_2}{\|x\|_2} \leq K(A) \left[\frac{\|\delta y\|_2}{\|y\|_2} + \frac{\|\delta A\|_2}{\|A\|_2} \right]$$

relative error in the solution x

• Cases 1-3: $\hat{x} = A^+y$. Similar bounds exist for these cases, and they again increase with $K(A)$.

Orthogonality principle in IPS (revisited): Say that

$$y = \sum_{i=1}^n c_i a_i + e$$

for some y . Then if we minimize e , $\sum c_i a_i$ becomes the projection of y onto $\text{span}\{a_1, \dots, a_n\}$. As we know in such a case, we require e to be perpendicular to $\text{span}\{a_i\}$, as per the orthogonality principle dictates. This gives us a certain set of equations we may solve to obtain the optimal c_i values, i.e. the coordinates of the projection \hat{y} in the $\{a_i\}_{i=1}^n$ basis.

$$\langle y - \sum_{i=1}^n c_i a_i, a_k \rangle = 0 \quad \forall k$$

$$\langle y, a_k \rangle = \sum_{i=1}^n c_i \langle a_i, a_k \rangle \quad \forall k$$

$$\langle y, a_k \rangle = \sum_{i=1}^n c_i a_k^H a_i \quad \forall k$$

$$\begin{bmatrix} \langle y, a_1 \rangle \\ \langle y, a_2 \rangle \\ \dots \\ \langle y, a_n \rangle \end{bmatrix} = \begin{bmatrix} G_{1j} = \langle a_j, a_1 \rangle \\ \dots \\ G_{ij} = \langle a_j, a_i \rangle \\ \dots \\ G_{nj} = \langle a_j, a_n \rangle \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ \dots \\ c_n \end{bmatrix}$$

→ optimal c_i 's

$$d = G \cdot c_{opt}$$

→ "Gram Matrix"

If a_i 's are linearly independent, G is invertible, and so

$$c_{opt} = G^{-1}d$$

Gram Matrix: For a set of vectors $\{a_i\}_{i=1}^n$, the Gram matrix is defined as

$$G_{ij} = \langle a_j, a_i \rangle$$

G is Hermitian symmetric ($G^H = G$) and is positive semi-definite. If further $\{a_i\}_{i=1}^n$ is linearly independent, G is positive definite and hence invertible. Furthermore if $\{a_i\}$'s are orthogonal, G is a diagonal matrix.

The Gram matrix method can be used to derive the projection matrix as well. Note that

$$d = A^H y \quad \& \quad G = A^H A$$

and a_i cpts are the coordinates in $\{a_i\}_{i=1}^n$ basis, (supposing they are),

$$\hat{y} = A c_{opt} = A G^{-1} d = \underbrace{A (A^H A)^{-1} A^H}_P y = P y$$

VECTOR SPACE OF RANDOM VARIABLES

IPS of random variables: let $\{y_i\}_{i=1}^m$ be a finite collection of ^{real valued} random variables, with $E[y_i^2] < \infty$ $\forall i$. Then $V = \text{span}\{y_i\}_{i=1}^m$ is an inner product space under the following inner product

$$\langle x, y \rangle = E[x y]$$

$$\|x\| = \sqrt{E[x^2]}$$

We treat $x=y$ iff $P(X(\omega) \neq Y(\omega), \omega: \text{outcome}) = 0$, so $x=y$ a.e.

The usual IPS facts look like the following

$$\bullet \frac{|E[x y]|^2}{E[x^2] E[y^2]} \Rightarrow \text{Cauchy-Schwarz Ineq.}$$

$$\bullet \sqrt{E[(x+y)^2]} \leq \sqrt{E[x^2]} + \sqrt{E[y^2]}$$

$$\bullet \rho_{x,y} = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y} = \frac{E[(x-\mu_x)(y-\mu_y)]}{\sqrt{E[(x-\mu_x)^2] E[(y-\mu_y)^2]}} = \cos \theta$$

θ is the angle between $x-\mu_x$ & $y-\mu_y$

$$\bullet x \perp y \Leftrightarrow E[x y] = 0$$

$$\bullet \text{Uncorrelated: } E[(x-\mu_x)(y-\mu_y)] = 0 \text{ i.e. } x-\mu_x \perp y-\mu_y$$

For zero mean r.v.'s, uncorrelated \Leftrightarrow orthogonal.

Decorrelation / Whitening. It is the process of obtaining rv's that are uncorrelated.

- 1) Let x, y be random variables. Then u, v given by the following equation are orthogonal

$$\begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 1/\|x\| & 0 \\ 0 & 1/\|y\| \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

- 2) Gram-Schmidt orthogonalization is another way of decorrelating a given set of rv's. If they are also zero mean, the process is also whitening.

Gram matrix: For random variables x_1, \dots, x_N , define $x = [x_1, \dots, x_N]^T$

$$G = \begin{bmatrix} \langle x_1, x_1 \rangle & \langle x_2, x_1 \rangle & \dots & \langle x_N, x_1 \rangle \\ \langle x_1, x_2 \rangle & \langle x_2, x_2 \rangle & \dots & \langle x_N, x_2 \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle x_1, x_N \rangle & \langle x_2, x_N \rangle & \dots & \langle x_N, x_N \rangle \end{bmatrix}$$

$$= \begin{bmatrix} E[x_1 x_1] & E[x_2 x_1] & \dots & E[x_N x_1] \\ E[x_1 x_2] & E[x_2 x_2] & \dots & E[x_N x_2] \\ \vdots & \vdots & \ddots & \vdots \\ E[x_1 x_N] & E[x_2 x_N] & \dots & E[x_N x_N] \end{bmatrix} = R_{xx}$$

Covariance matrix of x

$$= E[x x^T]$$

The Gram matrix is again Hermitian symmetric, & positive semi-definite. If further all vectors are linearly independent, G is positive definite & hence invertible.

LMMSE: Suppose you make $\{y_i\}_{i=1}^N$ observations of a random variable x and want to estimate it, by a linear sum $\alpha_1 y_1 + \dots + \alpha_N y_N$. Minimize the error by minimizing $\|x - \hat{x}\|^2 = \|e\|^2 = E[(x - \hat{x})^2]$. This amounts to projecting x onto $\text{span}\{y_i\}$.

$$e \perp \text{span}\{y_i\} \rightarrow e \perp y_i \quad \forall i$$

By the same derivation with Gram matrix previously,

$$\begin{matrix} G \alpha = d \\ \downarrow \\ R_{yy} \end{matrix} = \begin{matrix} \langle x, y_1 \rangle \\ \langle x, y_2 \rangle \\ \vdots \\ \langle x, y_N \rangle \end{matrix} = r_{xy}$$

$$\alpha_{\text{opt}} = R_{yy}^{-1} r_{xy} \Rightarrow \hat{x}_{\text{opt}} = \alpha_{\text{opt}}^T y$$

LMMSE estimator

Matrix inverse: $A^{-1} = \frac{1}{\det(A)} ((-1)^{ij} M_{ij})^T$
↪ cofactor

EE499 Cheat Sheet - Final

MSE of LMMSE Estimator:

$$\begin{aligned}
 E[(x-\hat{x})^2] &= E[e^2] = \|e\|^2 = \|x\|^2 - \|\hat{x}\|^2 = E[x^2] - E[\hat{x}^2] \\
 &= E[x^2] - E[\alpha^T y y^T \alpha] \\
 &= E[x^2] - \alpha^T E[yy^T] \alpha = E[x^2] - \alpha^T R_y \alpha \\
 &= E[x^2] - r_{xy}^T R_y^{-1} r_y \\
 &= E[x^2] - r_{xy}^T R_y^{-1} r_{xy} \\
 &= E[x^2] - r_{xy}^T \alpha = E[x^2] - \alpha^T r_{xy}
 \end{aligned}$$

because $e \perp \hat{x}$

$R_y^{-1} r_{xy}$

Vector LMMSE estimation: Say we have x_1, \dots, x_m many r.v.'s, which we observe through y_1, \dots, y_n . We know that y_i 's may depend on possibly all x_i 's. We will estimate x_i 's through y_j 's, using a linear model as follows:

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}_{M \times 1}, \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_{N \times 1}, \quad k_i = \begin{bmatrix} k_{i1} \\ \vdots \\ k_{in} \end{bmatrix}_{N \times 1}$$

$$x_i = k_i^T y \Rightarrow \hat{x} = \begin{bmatrix} -k_1^T \\ \vdots \\ -k_m^T \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = Ky$$

K y

Minimize MSE: $E[e^T e] = E\left[\sum_{i=1}^m (x_i - k_i^T y)^2\right] = \sum_{i=1}^m E[(x_i - k_i^T y)^2]$

\Rightarrow minimize each $E[(x_i - k_i^T y)^2]$ equivalently k_i

$\Rightarrow R_y k_i = r_{x_i y} \quad k_i$

Stack these equations horizontally

$$R_y \underbrace{\begin{bmatrix} k_1 & k_2 & \dots & k_m \end{bmatrix}}_{K^T} = \underbrace{\begin{bmatrix} r_{x_1 y} & r_{x_2 y} & \dots & r_{x_m y} \end{bmatrix}}_{R_{xy}^T, R_{xy} = E[xy^T]}$$

$R_y K^T = R_{xy}^T$

$\Rightarrow K R_y = R_{xy} \quad (\text{WIENER-HOPF EQUATION})$

If R_y is invertible,

$$K_{opt} = R_{xy} R_y^{-1}$$

What's the error for this estimator? (MSE = trace (Re))

$$\begin{aligned}
 R_e &= E[ee^T] = E[e(x - \hat{x})^T] = \\
 &= E[ex^T] - E[e\hat{x}^T] \quad \text{due to } \perp \text{ principle} \\
 &= E[(x - Ky)x^T] = E[xx^T] - E[Kyx^T] \\
 &= R_x - K R_{yx} = R_x - R_{xy} R_y^{-1} R_{yx}
 \end{aligned}$$

Now, say our observations y are of the form $y = Ax + n$

$$\begin{aligned}
 R_y &= E[yy^T] = E[(Ax + n)(Ax + n)^T] \\
 &= A E[xx^T] A^T + E[nn^T] + \dots \\
 &= A R_x A^T + R_n
 \end{aligned}$$

$$\begin{aligned}
 R_{xy} &= E[xy^T] = E[x(Ax + n)^T] = E[xx^T A^T] + E[xn^T] \\
 &= R_x A^T
 \end{aligned}$$

So the K estimator is, by Wiener-Hopf eqn.

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

Then our estimations are

$$\hat{x} = R_x A^T (A R_x A^T + R_n)^{-1} y$$

By the matrix inversion lemma, this is equivalent to

$$\hat{x} = (A^T R_n^{-1} A + R_x^{-1})^{-1} A^T R_n^{-1} y$$

Maximum Likelihood Estimation: Treat the unknown x as deterministic, observations as random.

$$\hat{x}_{ML} = \underset{x}{\operatorname{argmax}} \underbrace{p(y|x)}_{\text{likelihood function}} = \underset{x}{\operatorname{argmax}} \underbrace{\log p(y|x)}_{\text{log-likelihood function}}$$

ML estimators with different noise statistics:

- iid Gaussian: $\hat{x}_{ML} = \underset{x}{\operatorname{argmin}} \|y - Ax\|_2^2 = \hat{x}_{LS}$
- correlated Gaussian: $\hat{x}_{ML} = \underset{x}{\operatorname{argmin}} (y - Ax)^T \Sigma^{-1} (y - Ax)$
 $= \underset{x}{\operatorname{argmin}} \|y - Ax\|_{\Sigma^{-1}}$
"weighted norm"

- independent Poisson:

- iid Generalised Gaussian: $\hat{x}_{ML} = \operatorname{argmin}_x \|y - Ax\|_p$

Gaussian distribution reminder:

$$p_x(x) = \frac{1}{\sqrt{2\pi\sigma_x^2}} \exp\left(-\frac{1}{2} \frac{(x-\mu_x)^2}{\sigma_x^2}\right) = \frac{1}{\sqrt{2\pi\sigma_x^2}} \exp\left(-\frac{(x-\mu_x)^2}{2\sigma_x^2}\right)$$

$$p_x(x) = \frac{1}{\sqrt{(2\pi)^m |\Sigma|}} \exp\left(-\frac{1}{2} (\underline{x}-\underline{\mu})^T \Sigma^{-1} (\underline{x}-\underline{\mu})\right)$$

$x \in \mathbb{C}^m$

Bayesian Estimation: x & y are both random

$$\hat{x}_{\text{Bayesian}} = \operatorname{argmin}_{\hat{x}} E[\text{cost}(\hat{x}, x)]$$

↳ MMSE: $\text{cost}(\hat{x}-x) = (\hat{x}-x)^2$
 MAP: $\text{cost}(\hat{x}-x) = \begin{cases} 0, & |\hat{x}-x| \leq \Delta \\ 1, & |\hat{x}-x| > \Delta \end{cases}$

MMSE Estimation:

$$\hat{x}_{\text{MMSE}} = \operatorname{argmin}_{\hat{x}} E[(\hat{x}-x)^2] = E[x|y]$$

MAP Estimation:

$$\hat{x}_{\text{MAP}} = \operatorname{argmax}_x p(x|y) = \operatorname{argmax}_x \frac{p(y|x)p(x)}{p(y)}$$

$$= \operatorname{argmax}_x \log \frac{p(y|x)p(x)}{p(y)} \rightarrow \text{discard}$$

$$= \operatorname{argmax}_x \log p(y|x) + \log p(x)$$

$$= \operatorname{argmin}_x -\log p(y|x) - \log p(x)$$

General form of MAP estimators

$$\hat{x}_{\text{MAP}} = \operatorname{argmin}_x \underbrace{\Delta(y, Ax)}_{\text{model mismatch term / data fidelity term}} + \lambda R(x)$$

model mismatch term /
data fidelity term

↳ regularization term
↳ regularization parameter

Unconstrained optimization: Say $J: \mathbb{R}^N \rightarrow \mathbb{R}$ is an objective/cost function.
Then this problem is

$$\hat{x} = \underset{x}{\operatorname{argmin}} J(x)$$

Convexity: A cost function is called convex if

$$\alpha J(x_1) + (1-\alpha)J(x_2) \geq J(\alpha x_1 + (1-\alpha)x_2)$$

$$\forall x_1, x_2 \quad \alpha \in [0, 1]$$

If a cost function is convex, then it has no local minima other than the global minimum. So we can use any local optimization method at hand.

Local minimization methods: The algorithm is as follows:

$$n = 0$$

$$x^{(n)} = x^{(0)} \quad (\text{initial guess})$$

repeat

compute $d^{(n)}$, an update direction for x

compute $\tau^{(n)}$, a step size

$$x^{(n+1)} = x^{(n)} + \tau^{(n)} d^{(n)}$$

$$n = n + 1$$

until convergence

$$\hat{x} = x^{(n)}$$

Steepest Descent: Pick $d^{(n)}$ in the opposite direction of the gradient:

$$d^{(n)} = - \frac{\nabla J(x^{(n)})}{\|\nabla J(x^{(n)})\|_2}$$

Newton's Method: Use second order approximation of J to obtain faster convergence.

$$J(x + \tau d) \approx J(x) + \tau d^T \nabla J(x) + \frac{\tau^2}{2} d^T \nabla^2 J(x) d \quad \leftarrow \text{goal } J(x)$$

Minimize wrt. d & choose $\tau = 1$

$$\nabla_d (J(x) + \tau d^T \nabla J(x) + \frac{\tau^2}{2} d^T \nabla^2 J(x) d) = 0$$

$$\nabla J(x) + \frac{1}{2} d^T \nabla^2 J(x) d = 0$$

$$\Rightarrow d = - (\nabla^2 J(x^{(n)}))^{-1} \nabla J(x^{(n)})$$

$$\Rightarrow x^{(n+1)} = x^{(n)} - \underbrace{(\nabla^2 J(x^{(n)}))^{-1}} \nabla J(x^{(n)})$$

Quasi-Newton methods: Choose this ∇^2 matrix as something else, say B :

$$x^{(n+1)} = x^{(n)} - B \nabla J(x^{(n)})$$

Conjugate Gradient methods: Update the direction as well:

$$d^{(0)} = -\nabla J(x^{(0)})$$

$$x^{(n+1)} = x^{(n)} + \alpha d^{(n)}$$

$$d^{(n+1)} = -\nabla J(x^{(n+1)}) + B^{(n)} d^{(n)}$$

Choosing the step-size parameter α ideally

$$\alpha^{(n)} = \underset{\alpha}{\operatorname{argmin}} J(x^{(n)} + \alpha d^{(n)})$$

but this is computationally costly, unnecessary in guaranteeing convergence and there are approximate methods of choosing α

Global Optimization methods:

- Deterministic methods

- ↳ Exhaustive Search,
- ↳ Branch & Bound,
- ↳ Dynamic programming, ---

- Stochastic methods

- ↳ Simulated annealing,
- ↳ Genetic Algorithms, ---