Summarizing the Essentials

January 12, 2024

Contents

1	Structure of Vector Space	2
2	Linear Transformations and Linear Operators	4
3	Matrix as "Coordinates"	5
4	Similarity and Diagonalization	7
5	Matrix Seen In Itself	10
6	Optimization by Least Square Method	12
7	SVD and the "Big Picture"	15

Summarizing the Essentials

1 Structure of Vector Space

- The fundamental notion that emerges naturally from the *axioms of vector spaces* is the concept of **linear combinations** (of vectors), from which the following fundamental relations among vectors naturally follow
 - Linearly independent: "linear combinations = 0" has only the trivial solution.
 - Linearly dependent: *"linear combinations* = 0" has non trivial solution.
- A set of vectors in linearly dependent if and only if one of them can be expressed in terms of the others.
 - Basis: A set of linearly independent vectors such that all vectors can be written as linear combinations of them.
 - Dimension: The number of vectors in any basis.
- **Dimension Formula:** For U, W two subspaces of V, then U + W and $U \cap W$ are subspaces and we have

$$\dim\left(U+W\right)=\dim U+\dim W-\dim\left(U\cap W\right)$$

 This is compatible with the general Inclusion-Exclusion Principal (容斥 原理) and can be generalized accordingly, for example

$$\dim (U_1 + U_2 + U_3) = \dim U_1 + \dim U_2 + \dim U_3 - \dim U_3$$

$$-dim\,(U_{1}\cap U_{2})-dim\,(U_{1}\cap U_{3})-dim\,(U_{2}\cap U_{3})+dim\,(U_{1}\cap U_{2}\cap U_{3})$$

- Any n + 1 vectors in a *n*-dimensional vector space must be linear independent.
- **Basis Extension:** Any set of linearly independent vectors in V can be extended to a basis of V.
 - If V is endowed with an **inner product**, so that the notion of **orthogonality** makes sense, then we have the concept of **orthonormal basis**. Then any

linearly independent set can be transformed into an orthogonal set by by applying the **Gram-Schmidt process**.

- In particularly, any inner product space has orthogonal basis. And any orthogonal set can be extended to a orthonormal basis.
- A set of orthogonal vectors must be linearly independent.
- Gram-Schmidt process is based on the notion of orthogonal projection:

$$Proj_{\mathbf{w}}\mathbf{v} := \frac{\langle \mathbf{v}, \, \mathbf{w} \rangle}{||\mathbf{w}||^2} \mathbf{w}$$

- In general, if W is a subspace of V that has orthogonal basis $\{\mathbf{w}_1, \cdots, \mathbf{w}_k\}$, then for any $\mathbf{v} \in V$, its orthogonal projection onto subspace W is given by

$$Proj_W \mathbf{v} = \frac{\langle \mathbf{v}, \mathbf{w}_1 \rangle}{||\mathbf{w}_1||^2} \mathbf{w}_1 + \dots + \frac{\langle \mathbf{v}, \mathbf{w}_k \rangle}{||\mathbf{w}_k||^2} \mathbf{w}_k$$

– In particularly, if $\{\mathbf{w}_1,\cdots,\mathbf{w}_k\}$ is **orthonormal**, then we have

$$Proj_W \mathbf{v} = \langle \mathbf{v}, \, \mathbf{w}_1 \rangle \mathbf{w}_1 + \dots + \langle \mathbf{v}, \, \mathbf{w}_k \rangle \mathbf{w}_k$$

- Inner product structure induces normal (distance) function: $||\mathbf{v}||^2 = \langle \mathbf{v}, \mathbf{v} \rangle$. Conversely, any normal induces an inner product via polarization identity.
 - Real case:

$$\langle \mathbf{u}, \mathbf{v} \rangle = \frac{||\mathbf{u} + \mathbf{v}||^2 - ||\mathbf{u} - \mathbf{v}||^2}{4}$$

- Complex case:

$$\langle \mathbf{u}, \mathbf{v} \rangle = \frac{||\mathbf{u} + \mathbf{v}||^2 - ||\mathbf{u} - \mathbf{v}||^2 + i||\mathbf{u} + i\mathbf{v}||^2 - i||\mathbf{u} - i\mathbf{v}||^2}{4}$$

- Norm and inner product enjoy the following properties
 - Triangle inequality: $||\mathbf{u} + \mathbf{v}|| \le ||\mathbf{u}|| + ||\mathbf{v}||$.
 - Pythagoras' Law: If $\mathbf{u} \perp \mathbf{v}$, then $||\mathbf{u} + \mathbf{v}|| = ||\mathbf{u}|| + ||\mathbf{v}||$.
 - Cauchy-Schwarz inequality: $|\langle u, \mathbf{v} \rangle| \le ||\mathbf{u}||\mathbf{v}||$.
- For complex inner product, it is **conjugate linear** in the 2^{nd} slot, i.e.,

$$\langle \mathbf{u}, \, k\mathbf{v} + l\mathbf{w} \rangle = k \langle \mathbf{u}, \, \mathbf{v} \rangle + l \langle \mathbf{u}, \, \mathbf{w} \rangle$$

which is equivalent to $\langle \mathbf{u}, \mathbf{v} \rangle = \overline{\langle \mathbf{v}, \mathbf{u} \rangle}$. In particularly, $\langle \mathbf{u}, \mathbf{u} \rangle \in \mathbb{R}$. So that $||\mathbf{u}||^2 = \langle \mathbf{u}, \mathbf{u} \rangle$ holds in both cases.

• Let U^{\perp} the **orthogonal complement** of U in V, then $\dim U + \dim U^{\perp} = \dim V$.

2 Linear Transformations and Linear Operators

A linear transformation T between V and W is a map T : V → W that is linear, i.e., keeps the linear combinations of vectors in V and W. That is

$$T(k\mathbf{v} + l\mathbf{w}) = kT(\mathbf{v}) + lT(\mathbf{w})$$
 for $\mathbf{v} \in V, \mathbf{w} \in W$ and $\forall k, l$

- A linear transformation form V to V is called a **linear operator**.
- Linear transformations can be composed, $T_1 \circ T_2(\mathbf{v}) := T_1(T_2(\mathbf{v}))$. This is possible if $Rang(T_2) \subset Domain(T_1)$.
- Linear transformation can be **inverted** only if it is one-to-one.
- Linear transformation $T: V \to W$ gives a relation between V and W.
 - If T is both one to one and onto, then we call T establishes an **isomorphism** between V and W, and denoted by $V \cong W$.
 - Isomorphic vector spaces have the same dimension. And conversely, if $\dim V = \dim W$, then $V \cong W$. Indeed, suppose V and W have basis given respectively by $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ and $\{\mathbf{w}_1, \dots, \mathbf{w}_n\}$. Then the linear transformation given by $T(k_1\mathbf{v}_1 + \dots + k_n\mathbf{v}_n) := k_1\mathbf{w}_1 + \dots + k_n\mathbf{w}_n$ establishes an isomorphism $V \cong W$.
- The extent that a general linear transformation $T: V \to W$ fails to be an isomorphism is measured by its kernel and range.
 - $ker(T) := {\mathbf{v} \in V : T(\mathbf{v}) = \mathbf{0}} \subset V.$ T is one-to-one if and only if $ker(T) = {\mathbf{0}}$. So, the larger the size of ker, the further for T from being an **injection** (a.k.a one to one).
 - $Rang(T) := \{ \mathbf{w} \in W : \mathbf{w} = T(\mathbf{v}) \text{ for some } \mathbf{v} \in V \} \subset W$. T is onto if and only if Range(T) = W. So, the smaller the size of range, the further for T from being an surjection (a.k.a onto).
- Fundamental identity: $\dim \ker (T) + \dim \operatorname{Range} (T) = \dim (V)$.

3 Matrix as "Coordinates"

 We can *coordinalize* a vector space V by choosing a basis B := {v₁,..., v_n}, then we have the **coordinate map** which gives the following isomporphism

$$[\]_B: V \xrightarrow{\cong} \mathbb{R}^n \qquad \mathbf{v} = k_1 \mathbf{v}_1 + \dots + k_n \mathbf{v}_n \longmapsto \ [\mathbf{v}]_B := \begin{bmatrix} k_1 \\ \vdots \\ k_n \end{bmatrix}$$

 Different basis of V are related by invertible matrix. If B = {v₁, ..., v_n} and B' = {v'₁, ..., v'_n} are two basis of V. Then the transition matrix from B to B', which is denoted by P_{B'←B}, is defined to be

$$P_{B' \leftarrow B} := \left[[\mathbf{v}_1]_{B'} \mid \cdots \mid [\mathbf{v}_n]_{B'} \right] \quad \Longleftrightarrow \quad B = B' P_{B' \leftarrow B}$$

Then for $\mathbf{v} \in V$, we have the coordinate transformation formula

$$[\mathbf{v}]_{B'} = P_{B' \leftarrow B}[\mathbf{v}]_B$$

• Transition matrix has the following properties

$$\begin{array}{l} - \ P_{B\leftarrow B'} = P_{B'\leftarrow B}^{-1} \\ - \ P_{B\leftarrow B'}P_{B'\leftarrow B''} = P_{B\leftarrow B''}. \end{array}$$

- If S is the standard basis for \mathbb{R}^n , then for any basis $B = {\mathbf{v}_1, \dots, \mathbf{v}_n}$, we have

$$P_{S\leftarrow B} = [\mathbf{v}_1 \,| \cdots \,| \, \mathbf{v}_n] =: B$$

- For any two basis B and B' of \mathbb{R}^n , then above properties implies

$$P_{B'\leftarrow B}=P_{B'\leftarrow S}P_{S\leftarrow B}=P_{S\leftarrow B'}^{-1}P_{S\leftarrow B}=(B')^{-1}\cdot B$$

- Given $T: V \to W$, we can coordinalize T as follows
 - Choosing basis $B = {\mathbf{v}_1, \dots, \mathbf{v}_n}$ and $B' = {\mathbf{w}_1, \dots, \mathbf{w}_m}$ of V and W respectively. Then the matrix (representation) for $T : V \longrightarrow W$ relative to the basis B and B' is defined to be

$$[T]_{B'B} = \left[\begin{array}{c|c} [T(\mathbf{v}_1)]_{B'} \end{array} \middle| [T(\mathbf{v}_2)]_{B'} \end{array} \middle| \cdots \bigg| [T(\mathbf{v}_n)]_{B'} \end{array} \right]$$

• In terms of coordinate representations of vectors in V and W via the basis B and B' respectively, the linear transformation $T: V \longrightarrow W$ can be represented as a matrix transformation

$$[T(\mathbf{x})]_{B'} = [T]_{B'B} \cdot [\mathbf{x}]_B$$

In other words, at the level of coordinate representation, T : V → W is represented as T_A : ℝⁿ → ℝ^m, where A = [T]_{B'B} is the matrix of T relative to the basis B and B'. Graphically, we have the following

Thus, if we know the coordinate of $T(\mathbf{x})$, namely $[T(\mathbf{x})]_{B'}$, we can "recover" $T(\mathbf{x})$ from the inverse of the *coordinate map* $[\]_{B'}$ as follows

$$T = [\]_{B'}^{-1} \circ T_A \circ [\]_B$$

- If V = W, $T : V \to V$ is represented by a square matrix relative to a basis B of V. In this case, we write $[T]_B$ instead of $[T]_{BB}$.
- If $T: V \to V$ is invertible, then for any basis B of V, the matrix of T relative to it is also invertible, and $[T^{-1}]_B = [T]_B^{-1}$.
- The composition of linear transformations corresponds to the multiplication of the corresponding matrices.



• For a linear operator $T: V \to V$. Two matrix representations $[T]_B$ and $[T]_{B'}$ are related by similarity relation.

$$[T]_{B'} = P_{B \leftarrow B'}^{-1} [T]_B P_{B \leftarrow B'} = P_{B' \leftarrow B} [T]_B P_{B \leftarrow B'}$$

• If a real vector space V is endowed with an inner product structure, and we consider orthonormal basis with respect to it. Then the transition matrix between two orthonormal basis is **orthogonal**. i.e.,

$$A^T = A^{-1}$$
 or $AA^T = A^T A = I$.

- The rows and columns of an orthogonal matrix form an orthonormal basis.
- The product of two orthogonal matrix is orthogonal.
- The inverse of an orthogonal matrix is orthogonal.
- If a complex vector space V is endowed with an inner product structure, and we consider orthonormal basis with respect to it. Then the transition matrix between two orthonormal basis is **unitary**. i.e.,

$$A^* := \overline{A}^T = A^{-1} \quad \text{or} \quad AA^* = A^*A = I.$$

- The rows and columns of an unitary matrix form an orthonormal basis.
- The product of two unitary matrix is orthogonal.
- The inverse of an unitary matrix is unitary.
- A is orthogonal (unitary) if and only if A keeps the dot product, i.e.,

$$A\mathbf{v} \bullet A\mathbf{w} = \mathbf{v} \bullet \mathbf{w}$$

if and only if A keeps the norm, i.e.,

$$||A\mathbf{v}|| = ||\mathbf{v}||$$

4 Similarity and Diagonalization

- For a linear operator $T: V \to V$, which is represented by $A = [T]_B$. If we have another basis $B' = \{\mathbf{w}_1, \cdots, \mathbf{w}_n\}$ such that $[T]_{B'} = D = diag\{\lambda_1, \cdots, \lambda_n\}$. Then we call the linear operator T (or the corresponding matrix A) is **diagonalizable**.
 - This means $T\mathbf{w}_i = \lambda_i \mathbf{w}_i \quad \forall i = 1, \cdots, n$. And we have

$$D = [T]_{B'} = P_{B' \leftarrow B}[T]_B P_{B \leftarrow B'} = P_{B \leftarrow B'}^{-1} A P_{B \leftarrow B'}$$

Denote by $P = P_{B \leftarrow B'}$ the transition matrix, then the above becomes

$$D = P^{-1}AP$$

- In particularly, given $A \in M_{n \times n}$, we view it as matrix transformation $T_A : \mathbb{R}^n \to \mathbb{R}^n$. Then we see that
 - A is itself the matrix representation of T_A relative to the standard basis $S = \{\mathbf{e}_1, \cdots, \mathbf{e}_n\}.$
 - If A is diagonalizable, then \exists new basis $P = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ such that $A\mathbf{v}_i = \lambda_i \mathbf{v}_i$ for some scalars λ_i , $i = 1, \dots, n$.
 - Denote by $P = [\mathbf{v}_1 | \cdots | \mathbf{v}_n]$, then $A\mathbf{v}_i = \lambda_i \mathbf{v}_i$ means AP = PD for $D = diag\{\lambda_1, \cdots, \lambda_n\}$. That is $D = P^{-1}AP$ where $P = P_{S \leftarrow B}$.
- Conditions of diagonalizability (for both real and complex matrices):
 - − $A \in \mathbb{M}_{n \times n}$ is diagonalizable if and only if A has n linearly independent eigenvectors.
 - -A is diagonalizable if and only if the sum of dimensions of distinct eigenspaces (a.k.a geometric multiplicities) is equal to n.
 - -A is diagonalizable if and only if for each eigenvalue, its algebraic multiplicity equals its geometric multiplicity.
 - * In particularly, if A has n distinct eigenvalues, i.e., if all eigenvalues of A are distinct (all eigenvalues having algebraic multiplicity 1), then A is diagonalizable.
 - * Geometric multiplicity \leq Algebraic multiplicity.

• Methods of diaonalizability:

- First, find the eigenvalues of A by computing the *characteristic polynomial*

$$P_A(\lambda) = \det\left(\lambda I - A\right) = 0$$

- Then for each eigenvalue λ , find a basis for the corresponding eigenspace

$$E_{\lambda} = \{ \mathbf{v} : A\mathbf{v} = \lambda \mathbf{v} \} = null \, (\lambda I - A).$$

by solving the corresponding linear system $(\lambda I - A)\mathbf{x} = \mathbf{0}$.

- Putting all n basis vectors for all eigenspaces together to form the transition matrix P (from the basis formed by these basis vectors to the standard basis).
- And P can diagonalize A in the sense $P^{-1}AP = D$ where D is diagonal with diagonal entries being the eigenvalues (counted with multiplicity) that are so arranged in order of the corresponding eigenvectors in P.

• Orthogonal and Unitary Diagonalization

- If we can find an orthogonal P such that $P^T A P = D$, we say that A is orthogonally diagonalizable. This amounts saying we can choose an orthonormal basis $B = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ for \mathbb{R}^n such that $A\mathbf{v}_i = \lambda_i \mathbf{v}_i$ for some scalars λ_i .
 - * A real matrix A is orthogonally diagonalizable if and only if A is symmetric, i.e., $A^T = A$.
 - * For symmetric A, all its eigenvalues are real, and the eigenvectors corresponding to different eigenvalues are orthogonal.
- If we can find an unitary P such that $P^*AP = D$, we say that A is **unitarily diagonalizable**. This amounts saying we can choose an orthonormal basis $B = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ for \mathbb{C}^n such that $A\mathbf{v}_i = \lambda_i \mathbf{v}_i$ for some scalars λ_i .
 - * A complex matrix A is unitarily diagonalizable if and only if A is **Herim**itian, i.e., $A^* = A$.
 - * For Hermitian A, all its eigenvalues are real, and the eigenvectors corresponding to different eigenvalues are orthogonal.

• Methods of Orthogonal (Unitary) Diagonalization:

- First, find all eigenvalues of A by computing the characteristic polynomial.
- Find a basis for each eigenspace. It must be an orthogonal basis.
- Putting the basis vectors for all distinct eigenspaces together, and apply the Gram-Schmidt process to get an orthonormal basis.
- Using these basis as columns of an invertible P (which must be orthogonal in the real case, and unitary in the complex case).
- -P computed as above will orthogonally (unitarily) diagonalize A.

5 Matrix Seen In Itself

- For $A \in \mathbb{M}_{m \times n}$, the space spanned by its rows (columns) is called the *row (column)* space of A, denoted by *row*(A) and *col*(A) respectively.
- $\dim row(A) = \dim col(A) =: rank(A)$.
- A is invertible if and only if it is of full rank.
- Elementary row and column operations do not alter the rank of a matrix.
 - If R is a reduced row echelon form of A, then

rank(A) = rank(R) = number of nonzero rows in R = number of leading 1's in R

- If
$$rank(A) = r$$
 then \exists invertible P and Q such that $PAQ = \begin{bmatrix} I_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$

- $A\mathbf{x} = \mathbf{b}$ is solvable if and only if $ranl(A) = rank(A|\mathbf{b})$; if and only if $\mathbf{b} \in col(A)$.
- For any matrix A, $rank(A) = rank(A^T) = rank(AA^T) = rank(A^TA)$.
- $row(A)^{\perp} = null(A)$ and $col(A)^{\perp} = null(A^T)$
- $nullity(A) := dim \, null(A) = dim \{ \mathbf{x} : A\mathbf{x} = \mathbf{0} \}.$
- Fundamental relation: For $A \in \mathbb{M}_{m \times n}$, we have

$$nullity(A) + rank(A) = n$$

- Viewing A as matrix transformation $T_A : \mathbb{R}^n \to \mathbb{R}^m$, the above relation is the same as $\dim \ker (T_A) + \dim \operatorname{Range} (T_A) = n$.
- QR-decomposition: If A is an m×n matrix with linearly independent column vectors, i.e., full column rank (列满秩) then A can be factored as A = QR where Q is an m×n matrix with orthonormal column vectors, and R is an n×n invertible upper triangular matrix. In particularly, every invertible matrix has a QR-decomposition.
- If $A \in \mathbb{M}_{m \times n}$, and if $\lambda_1 \ge \lambda_2 \ge \dots \ge \lambda_n \ge 0$ are the eigenvalues of $A^T A$, then the numbers

$$\sigma_1=\sqrt{\lambda_1}, \ \sigma_2=\sqrt{\lambda_2}, \ \cdots, \ \sigma_n=\sqrt{\lambda_n}$$

are called the (singular values) 奇异值 of A.

• Singular Value Decomposition (SVD) Expanded Form: If A is an $m \times n$ matrix of rank k, then A can be factored as

$$A = U\Sigma V^T = \begin{bmatrix} \mathbf{u}_1 \ \mathbf{u}_2 \ \cdots \ \mathbf{u}_k \ | \ \mathbf{u}_{k+1} \ \cdots \ \mathbf{u}_m \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 \ \cdots & 0 & \\ 0 & \sigma_2 \ \cdots & 0 & \\ \vdots & \vdots & \vdots & \\ 0 & 0 \ \cdots & \sigma_k & \\ \hline \mathbf{0}_{(m-k) \times k} & \mathbf{0}_{(m-k) \times (n-k)} \end{bmatrix} \begin{bmatrix} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \vdots \\ \mathbf{v}_k^T \\ \hline \mathbf{v}_{k+1}^T \\ \vdots \\ \mathbf{v}_n^T \end{bmatrix}$$

in which U, Σ and V have sizes $m \times m, m \times n$ and $n \times n$ respectively, and in which:

- a). $V = [\mathbf{v}_1 \ \mathbf{v}_2 \cdots \mathbf{v}_n]$ orthogonally diagonalizes $A^T A$.
- b). The nonzero diagonal entries of Σ are $\sigma_1 = \sqrt{\lambda_1}$, $\sigma_2 = \sqrt{\lambda_2}$, \cdots , $\sigma_k = \sqrt{\lambda_k}$, where $\lambda_1, \lambda_2, \cdots, \lambda_k$ are the nonzero eigenvalues of $A^T A$ corresponding to the column vectors of V.
- c). The column vectors of V are ordered so that $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_k > 0.$

d).
$$\mathbf{u}_i = \frac{A\mathbf{v}_i}{||A\mathbf{v}_i||} = \frac{1}{\sigma_i} \, A\mathbf{v}_i \quad i=1,2,\cdots,k$$

- e). $\{\mathbf{u}_{1},\mathbf{u}_{2},\cdots,\mathbf{u}_{k}\}$ is an orthonormal basis for $col\left(A\right).$
- f). { $\mathbf{u}_1, \mathbf{u}_2, \cdots, \mathbf{u}_k, \mathbf{u}_{k+1}, \cdots, \mathbf{u}_m$ } is an extension of { $\mathbf{u}_1, \mathbf{u}_2, \cdots, \mathbf{u}_k$ } to an orthonormal basis for \mathbb{R}^m .

- The reduced singular value decompositon of A reads

$$A = \underbrace{\left[\mathbf{u}_1 \ \mathbf{u}_2 \cdots \mathbf{u}_k\right]}_{U_1} \underbrace{\left[\begin{array}{cccc} \sigma_1 & 0 & \cdots & 0\\ 0 & \sigma_2 & \cdots & 0\\ \vdots & \vdots & & \vdots\\ 0 & 0 & \cdots & \sigma_k\end{array}\right]}_{\Sigma_1} \underbrace{\left[\begin{array}{cccc} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \vdots \\ \mathbf{v}_k^T \end{array}\right]}_{V_1^T}$$

In this form, the sizes of U_1 , Σ_1 and V_1^T are $m \times k$, $k \times k$ and $k \times n$ respectively. And Σ_1 is invertible. When expanding, the above can be further written as

$$A = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \dots + \sigma_k \mathbf{u}_k \mathbf{v}_k^T$$

which is called a **reduced singular value expansion** of *A*. This applies to **all** matrices, whereas the *spectral decomposition* applies only to symmetric matrices.

* Spectral decomposition for symmetric matrix: Given $A \in \mathbb{M}_{n \times n}(\mathbb{R})$ that is symmetric, we know that $\exists orthogonal P = [\mathbf{u}_1 \cdots \mathbf{u}_n]$ that diagonalizes A, i.e.,

6 Optimization by Least Square Method

The least squares problem: $\min_{\forall \mathbf{x} \in \mathbb{R}^n} \{ ||\mathbf{b} - A\mathbf{x}|| \}$ admits solutions that are given by $\hat{\mathbf{x}}$ such that

$$A\hat{\mathbf{x}} = proj_{col\,(A)}\mathbf{b}. \qquad \Longrightarrow \quad ||\mathbf{b} - A\hat{\mathbf{x}}|| = \min_{\forall \, \mathbf{x} \in \mathbb{R}^n} \left\{ ||\mathbf{b} - A\mathbf{x}|| \right\}$$

Notice that such $\hat{\mathbf{x}}$, when exists, may not be unique, and any such solution is called a **least squares solution** of $A\mathbf{x} = \mathbf{b}$. Each such solution $\hat{\mathbf{x}}$ has the same error vector and thus the same least squares error.

Solving the above least squares problem is equivalent to solving the so called *normal equation:* $A^T A \mathbf{x} = A^T \mathbf{b}$. Notice that the normal equation is always consistent even if $A \mathbf{x} = \mathbf{b}$ is inconsistent. Indeed, this establishes an *one to one* correspondence between the following solution sets

$$\left\{\begin{array}{c} Solutions \ to \\ A\mathbf{x} = proj_{col\,(A)}\mathbf{b} \end{array}\right\} \quad \stackrel{1-1}{\longleftrightarrow} \quad \left\{\begin{array}{c} Solutions \ to \\ A^TA\mathbf{x} = A^T\mathbf{b} \end{array}\right\}$$

And the **uniqueness** of least square solutions is equivalent to the **uniqueness** of solutions to the associated normal system.

- When A is of full column rank, i.e., the column vectors of A are linearly independent, then $A^T A$ is invertible, and the normal equation $A^T A \mathbf{x} = A^T \mathbf{b}$ has **unique solution** $\hat{\mathbf{x}} = (A^T A)^{-1} A^T \mathbf{b}$, thus the least squares problem admits unique solution. In particular, the orthogonal projection $proj_{col(A)}\mathbf{b}$ is given by $A\hat{\mathbf{x}} = A(A^T A)^{-1}A^T\mathbf{b}$. In general, if $W \subseteq V$ is a subspace that can be described by col(A) for some matrix A, then orthogonal projection $proj_W$ is given by $proj_W = A(A^T A)^{-1}A^T$
- When A is of full column rank, by applying Gram-Schmidt process to the set of column vectors followed by normalization, we will get the QR decomposition of A, i.e., A = QR, where Q is orthogonal and R is upper triangular, then for each $\mathbf{b} \in \mathbb{R}^m$ the system $A\mathbf{x} = \mathbf{b}$ has a unique least squares solution given by $\hat{\mathbf{x}} = R^{-1}Q^T\mathbf{b}$

Least Squares Lines of Best Fit: Suppose we want to fit a straight line y = a + bx to the experimentally determined points

$$(x_1,\,y_1),\,\,(x_2,\,y_2),\,\cdots,(x_n,\,y_n)$$

Optimization Problem: $\min_{a,b} \left\{ \sum_{i=1}^{n} [y_i - (a + bx_i)]^2 = \sum_{i=1}^{n} d_i^2 \right\}$ where $d_i = |y_i - (a + bx_i)|$ measures the *error in* y_i *at the point* x_i , which are called the *residuals* (残余).

To solve this problem, we first turn it into a linear algebra problem. If there's no error, we would have

$$\begin{cases} y_1 = a + bx_1 \\ y_2 = a + bx_2 \\ \vdots \\ y_n = a + bx_n \end{cases}$$

Denote by $\mathbf{y} = [y_1 \ y_2 \ \cdots \ y_n]^T$, $\mathbf{v} = [a \ b]^T$ and $M = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$ then the above

system can be written as $M\mathbf{v} = \mathbf{y}$.

When the data $\{(x_i, y_i)\}$ are not co-linear (共线), then there is **NO** $\mathbf{v} = [a \ b]^T$ that solves $M\mathbf{v} = \mathbf{y}$.

However, we can seek the **least squares solutions** that always exist. The *least squares error* in this case is given by

$$e = ||\mathbf{y} - M\mathbf{v}|| = \left\| \begin{bmatrix} y_1 - (a + bx_1) \\ y_2 - (a + bx_2) \\ \vdots \\ y_n - (a + bx_n) \end{bmatrix} \right\| = \sqrt{sum \text{ of the squares of the data errors}}$$

Thus, the least squares solution of $M\mathbf{v} = \mathbf{y}$ will yield the least squares line of best fit. We consider the normal equation associated to $M\mathbf{v} = \mathbf{y}$.

$$M^T M \mathbf{v} = M^T \mathbf{y}$$

$$M^{T}M = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_{1} & x_{2} & \cdots & x_{n} \end{bmatrix} \begin{bmatrix} 1 & x_{1} \\ 1 & x_{2} \\ \vdots & \vdots \\ 1 & x_{n} \end{bmatrix} = \begin{bmatrix} n & \sum_{i=1}^{n} x_{i} \\ \sum_{i=1}^{n} x_{i} & \sum_{i=1}^{n} x_{i}^{2} \end{bmatrix}$$

So $det(M^T M) \neq 0$ when $\sigma^2 \neq 0$, i.e., there is error. In this case, the normal equation has unique solution, thus we have **unique** line of best fit given by

$$y = a^* + b^* x$$
 where $\mathbf{v}^* = \begin{bmatrix} a^* \\ b^* \end{bmatrix} = (M^T M)^{-1} M^T \mathbf{y}.$

The above can be generalized to the case when we need to find a polynomial $y = a_0 + a_1 x + a_2 x^2 + \dots + a_m x^m$ that best fit the data $\{(x_1, y_1), (x_2, y_n), \dots, (x_n, y_n)\}$. That is, find a_0, a_1, \dots, a_m such that the following quantity is minimized

$$erro = \left\| \left[\begin{array}{c} y_1 - a_0 - a_1 x_1 - \dots - a_m x_1^m \\ y_2 - a_0 - a_1 x_2 - \dots - a_m x_2^m \\ \vdots \\ y_n - a_0 - a_1 x_n - \dots - a_m x_n^m \end{array} \right] \right\|$$

That is we need to find the *least squares solutions* to the following system

$$\begin{cases} a_{0} + a_{1}x_{1} + \dots + a_{m}x_{1}^{m} = y_{1} \\ a_{0} + a_{1}x_{2} + \dots + a_{m}x_{2}^{m} = y_{2} \\ \vdots \\ a_{0} + a_{1}x_{n} + \dots + a_{m}x_{n}^{m} = y_{n} \end{cases} \iff \underbrace{ \begin{bmatrix} 1 & x_{1} & \dots & x_{1}^{m} \\ 1 & x_{2} & \dots & x_{2}^{m} \\ \vdots & \vdots & \vdots \\ 1 & x_{n} & \dots & x_{n}^{m} \end{bmatrix}}_{M} \underbrace{ \begin{bmatrix} a_{0} \\ a_{1} \\ \vdots \\ a_{m} \end{bmatrix}}_{\mathbf{x}} = \underbrace{ \begin{bmatrix} y_{1} \\ y_{2} \\ \vdots \\ y_{n} \end{bmatrix}}_{\mathbf{x}}$$

which is equivalent to solving the corresponding normal equation: $M^T M \mathbf{x} = M^T \mathbf{y}$.

7 SVD and the "Big Picture"

Given $A \in \mathbb{M}_{m \times n}(\mathbb{R})$. A can be viewed as a linear transformation from \mathbb{R}^n to \mathbb{R}^m by left multiplication; similarly A^T can be viewed as a linear transformation from \mathbb{R}^m to \mathbb{R}^n .



We know that the fundamental facts about the pair of maps T_A, T_{A^T} are given as follows:

$$\begin{split} \dim \ker \left(T_A \right) + \dim R(T_A) &= \dim \left(\mathbb{R}^n \right) = n \\ \dim \ker \left(T_{A^T} \right) + \dim R(T_{A^T}) &= \dim \left(\mathbb{R}^m \right) = m \end{split}$$

In terms of the language of matrix, the above become

$$nullity\left(A\right)+rank\left(A\right)=n;\quad nullity\left(A^{T}\right)+rank\left(A^{T}\right)=m\quad (*)$$

Of course we know that $rank(A) = rank(A^T)$ (i.e., the row rank equals the column rank), that the above two equations can be combined into one identity (recall that we have called it the "baby index theorem")

$$nullity(A) - nullity(A^T) = n - m$$

which gives the full information about the solvability of the linear system $A\mathbf{x} = \mathbf{b}$.

Next, we endow \mathbb{R}^n and \mathbb{R}^m with the standard inner product structure, then we have the notion of *orthogonal complement* of a subspace. And we know the following

important facts:

$$null(A)^{\perp} = col(A^T);$$
 $null(A^T)^{\perp} = col(A)$

As we know that if $U \subseteq V$, then $\dim U + \dim U^{\perp} = \dim V$, thus by applying "dim" on both sides of the above orthogonal relation, we get the fundamental relations (*).

The above discussion can be "encoded" into the following picture, which we call "big picture".



In this framework, SVD for A turns to be an algorithm for computing all ingredients contained in the above "big picture".

$$A = U\Sigma V^T = \begin{bmatrix} \mathbf{u}_1 & \cdots & \mathbf{u}_k \ | \ \mathbf{u}_{k+1} & \cdots & \mathbf{u}_m \end{bmatrix} \begin{bmatrix} \sigma_1 & \cdots & 0 & & \\ \vdots & \ddots & \vdots & \mathbf{0}_{k \times (n-k)} \\ 0 & \cdots & \sigma_k & & \\ \hline & \mathbf{0}_{(m-k) \times k} & | \ \mathbf{0}_{(m-k) \times (n-k)} \end{bmatrix} \begin{bmatrix} \mathbf{v}_1^T \\ \vdots \\ \mathbf{v}_k^T \\ \hline \mathbf{v}_{k+1}^T \\ \vdots \\ \mathbf{v}_n^T \end{bmatrix}$$

Then we claim that

- $\{\mathbf{u}_1, \cdots, \mathbf{u}_k\}$ is an orthonormal basis for col(A).
- $\{\mathbf{u}_{k+1}, \cdots, \mathbf{u}_m\}$ is an orthonormal basis for $col(A)^{\perp} = null(A)$.
- + $\{\mathbf{v}_{k+1},\cdots,\mathbf{v}_n\}$ is an orthonormal basis for null(A).
- $\{\mathbf{v}_1,\cdots,\mathbf{v}_k\}$ is an orthonormal basis for $null(A)^{\perp}=col(A^T).$

Explanation: Recall that $\{\mathbf{v}_1, \dots, \mathbf{v}_k, \mathbf{v}_{k+1}, \dots, \mathbf{v}_n\}$ and \mathbb{R}^n is an orthonormal basis for \mathbb{R}^n , and $\{\mathbf{u}_1, \dots, \mathbf{u}_k, \mathbf{u}_{k+1}, \dots, \mathbf{u}_m\}$ is an orthonormal basis for \mathbb{R}^m such that

$$A^TA\mathbf{v}_i=\sigma_i^2\mathbf{v}_i, i=1,\cdots,k; \qquad A^TA\mathbf{v}_j=\mathbf{0}, \ j=k+1,\cdots,n$$

As $null(A^TA) = null(A)$, we also have $A\mathbf{v}_j = \mathbf{0}, \ j = k + 1, \cdots, n$.

Also recall that $\mathbf{u}_i = \frac{A\mathbf{v}_i}{\|A\mathbf{v}_i\|} = \frac{A\mathbf{v}_i}{\sigma_i}$, $i = 1, \dots, k$, and $\{\mathbf{u}_1, \dots, \mathbf{u}_k; \mathbf{u}_{k+1}, \dots, \mathbf{u}_m\}$ is an extension to an orthonormal basis for \mathbb{R}^m .

In particularly, we have $A\mathbf{v}_i = \sigma_i \mathbf{u}_i$, $i = 1, \dots, k$. As $\sigma_j = 0$ for $j = k + 1, \dots, n$, and $A\mathbf{v}_j = \mathbf{0}$, $j = k + 1, \dots, n$, we can also write $A\mathbf{v}_i = \sigma_i \mathbf{u}_i$, for all i.

Besides one can prove that $AA^T \mathbf{u}_i = \sigma_i^2 \mathbf{u}_i, \, \forall i.$ Indeed, we have

$$AA^T\mathbf{u}_i = AA^T\frac{A\mathbf{v}_i}{\sigma_i} = \frac{1}{\sigma_i}A(A^TA\mathbf{v}_i) = \frac{1}{\sigma_i}A(\sigma_i^2\mathbf{v}_i) = \sigma_i^2\frac{A\mathbf{v}_i}{\sigma_i} = \sigma_i^2\mathbf{u}_i, \ i = 1, \cdots, k$$

For $j = k + 1, \cdots, n$, $\sigma_j^2 = 0$, and by SVD algorithm, we know that

$$span\{\mathbf{u}_{k+1},\cdots,\mathbf{u}_m\}=col(A)^{\perp}=null(A^T)$$

That is $A^T \mathbf{u}_j = \mathbf{0}$, and consequently $AA^T \mathbf{u}_j = \mathbf{0}, j = k + 1, \dots, m$. Thus we see that $AA^T \mathbf{u}_i = \sigma_i^2 \mathbf{u}_i$ holds for all *i*.

In conclusion, we see that $\mathbf{v}'_i s$ are eigenvectors of $A^T A$ (with corresponding non zero eigenvalues given by $\sigma_1^2 \ge \cdots \ge \sigma_k^2 > 0$); while $\mathbf{u}'_i s$ are eigenvectors of AA^T (with the same set of eigenvalues as that of $A^T A$). Recall that AA^T and $A^T A$ have the same set of eigenvalues.

With the above preparations, SVD for A follows easily. Indeed, we have

$$\begin{split} AV &= A[\mathbf{v}_1, \cdots, \mathbf{v}_k, \mathbf{v}_{k+1}, \cdots, \mathbf{v}_n] = [A\mathbf{v}_1, \cdots, A\mathbf{v}_k, \mathbf{0}, \cdots, \mathbf{0}] = [\sigma_1 \mathbf{u}_1, \cdots, \sigma_k \mathbf{u}_k, \mathbf{0}, \cdots, \mathbf{0}] \\ &= \underbrace{[\mathbf{u}_1, \cdots, \mathbf{u}_k; \mathbf{u}_{k+1}, \cdots, \mathbf{u}_m]}_{U} \underbrace{\begin{bmatrix} \sigma_1 & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \sigma_k \\ \hline \mathbf{0}_{(m-k) \times k} & \mathbf{0}_{(m-k) \times (n-k)} \end{bmatrix}}_{\Sigma} \end{split}$$

Multiplying both sides by V^T from the right, and using the orthogonal property

 $VV^T = I_n$, we get finally that

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