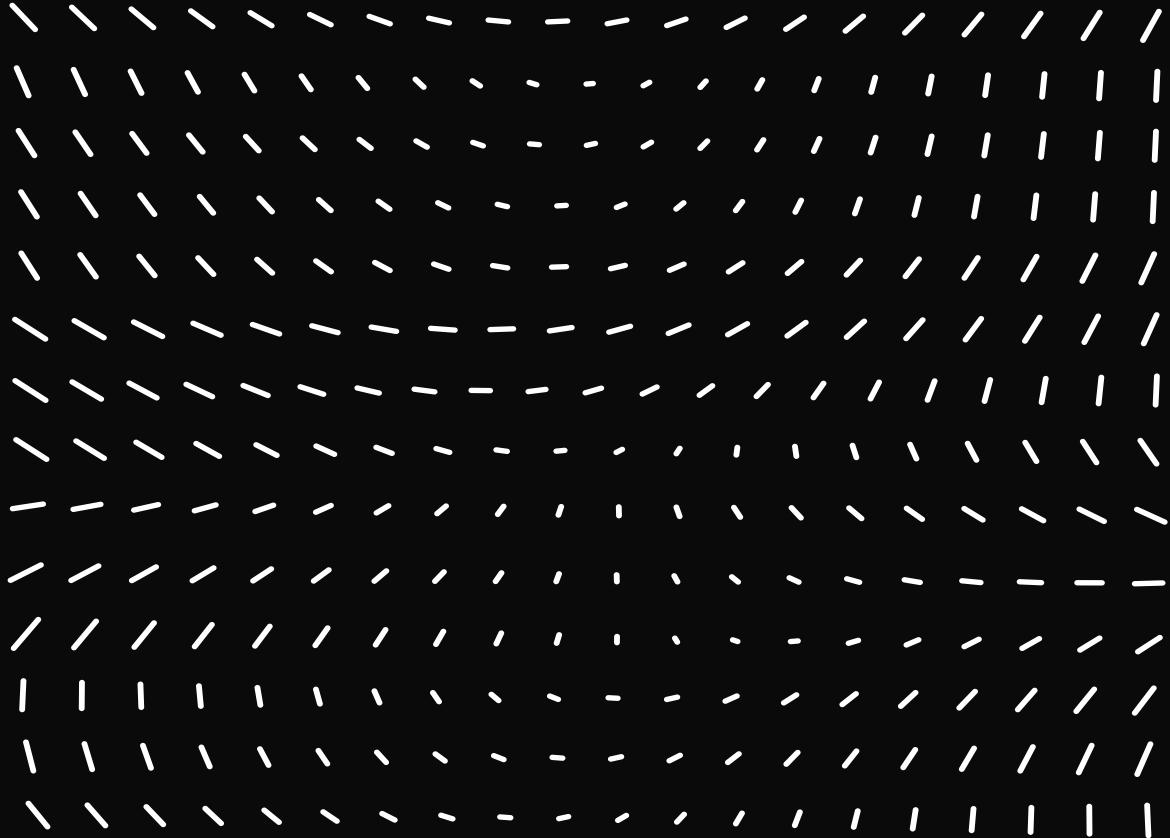


Linear Algebra



Vector Spaces

Span is a set of vectors that can be expressed as a linear combination.

$$\text{Span}(\{x_1, \dots, x_n\}) \triangleq \left\{ v : v = \sum_{i=1}^n d_i x_i, d_i \in \mathbb{R} \right\}$$

If vectors $\{x_1, \dots, x_n\}$ are set of linearly independent vectors, where each $x_i \in \mathbb{R}^n$, then span of the vectors will also be \mathbb{R}^n i.e., $\text{span}(\{x_1, \dots, x_n\}) = \mathbb{R}^n$.

Basis

A basis B is a set of linearly independent vectors that spans the whole space

$$\text{Span}(B) = \mathbb{R}^n$$

Linear Maps

$$f: V \rightarrow W \text{ such that } f(v+w) = f(v) + f(w)$$

and $f(av) = a f(v) \quad \forall v, w \in V$

Here 'a' is some constant. For ex: $f(2x)$

We can compute $y = f(x) \in \mathbb{R}^m$ for any $x \in \mathbb{R}^n$ as:

$$y = \left(\sum_{j=1}^n a_{1j} x_j, \dots, \sum_{j=1}^n a_{mj} x_j \right)$$

which is basically computing

$$y = Ax$$

If the function is invertible, then we can write

$$x = A^{-1}y$$

Range and Null space

In equation, $y = Ax$, the range of A is the span of columns in A . Formally, this is written as

$$\text{range}(A) \triangleq \{v \in \mathbb{R}^m : v \in Ax, x \in \mathbb{R}^n\}$$

Think of **range** as the set of vectors that can be reached or generated by A when multiplied with ' x '!

Null space is the set of vectors that get mapped to the null vector when multiplied by A

$$\text{nullspace}(A) \triangleq \{x \in \mathbb{R}^n : Ax = 0\}$$

Linear Projection

Given $y \in \mathbb{R}^m$ and a span of vectors $\{x_1, \dots, x_n\} \in \mathbb{R}^m$, the projection is finding the closest vector in the span w.r.t. to ' y ' in euclidean distance norm $\|v - y\|_2$ where $v \in \text{span}(\{x_1, \dots, x_n\})$

$$\text{Proj}(y; \{x_1, \dots, x_n\}) = \underset{v \in \text{span}(\{x_1, \dots, x_n\})}{\arg \min} \|y - v\|_2$$

To compute projection with a matrix,

$$\text{Proj}(y; A) = \underset{v \in \text{range}(A)}{\arg \min} \|v - y\|_2 = A(A^T A)^{-1} A^T y$$

This is an important formula for orthogonal projections in Least Squares linear regression. (11.2.2.2 in book)

7.1.3 Norms of a vector and matrix

"ways of measuring the size of a vector and a matrix"

Vector Norms

A norm is a measure of the length of the vector.

Denoted as $\|x\|$

Commonly used norms are:

$$p\text{-norm } \|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p} \text{ for } p \geq 1$$

if $p=1$,

$$1\text{-norm } \|x\|_1 = \sum_{i=1}^n |x_i|$$

if $p=2$,

$$2\text{-norm } \|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2} \quad \text{a.k.a Euclidean norm}$$

$$\text{Note: } \|x\|_2^2 = x^T x$$

$$\text{Max-norm } \|x\|_\infty = \max_i |x_i|$$

$$0\text{-norm } \|x\|_0 = \sum_{i=1}^n \mathbb{I}(|x_i| > 0) \quad \text{a.k.a pseudo norm}$$

It counts the number of non-zero elements in x .

Text book suggests, if we define $0^0=0$, we can write the pseudo norm as $\|x\|_0 = \sum_{i=1}^n x_i^0$. But how?

Note: we are only defining $0^0=0$ not for the rest of the numbers. So, if $x = [1, 0, 0, 1, 0]$

$$\begin{aligned} \|x\|_0 &= x_1^0 + x_2^0 + x_3^0 + x_4^0 + x_5^0 \\ &= 1^0 + 0^0 + 0^0 + 1^0 + 0^0 \\ &= 2 \checkmark \end{aligned}$$

Matrix Norms

Imagine a linear function such as $f(x) = Ax$ where $A \in \mathbb{R}^{m \times n}$

Here the function 'f' lengthens the matrix 'A' to any unit norm.

In matrix norm, we seek to find the maximum amount by which f could lengthen

the matrix. We refer this as induced norm of A which is given as

$$\|A\|_p = \max_{x \neq 0} \frac{\|Ax\|_p}{\|x\|_p} = \max_{\|x\|=1} \|Ax\|_p$$

For $p=2$, the matrix norm would be

$$\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)} = \max_i \sigma_i$$

σ_i is the i^{th} singular value.

What is λ_{\max} ?

It is the maximum eigenvalue

Nuclear Norm

aka Trace Norm defined as

$$\|A\|_* = \text{tr}(\sqrt{A^T A}) = \sum_i \sigma_i$$

where $\sqrt{A^T A}$ is the matrix square root

Since the singular values are always non-negative, we have

$$\|A\|_* = \sum_i |\sigma_i| = \|\sigma\|_1$$

There are different ways of calculating norms. One among them is **Schatten-p-norm**

Schatten p-norm

$$\|A\|_p = \left(\sum_i \sigma_i^p(A) \right)^{1/p}$$

Frobenius Norm

$$\|A\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^n a_{ij}^2} = \sqrt{\text{tr}(A^T A)} = \|\text{vec}(A)\|_2$$

In the case of Frobenius norm, we can think the matrix as a vector which is given as $\|A\| = \|\text{vec}(A)\|$. Above formula is the case of vector in 2-norm.

Another way of solving Frobenius norm is through singular values

i.e.,

$$\|A\|_F = \sqrt{\sigma_1^2 + \sigma_2^2 + \sigma_3^2 + \dots + \sigma_n^2}$$

If A is a huge matrix, then computing norm could be expensive.

A stochastic approximation to the Frobenius norm could be created by

Hutchinson Trace Estimation technique i.e.,

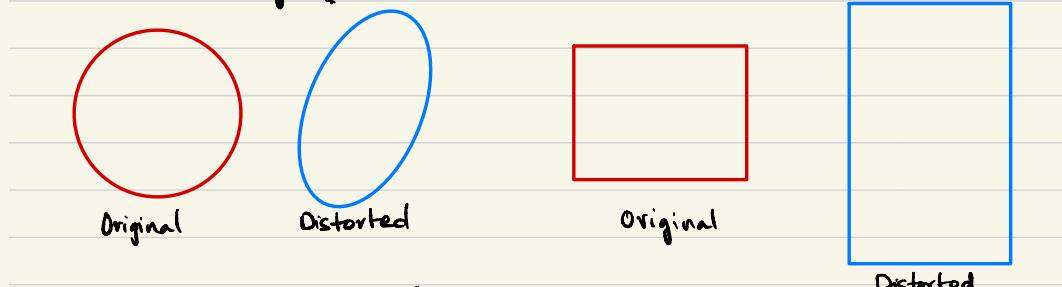
$$\|A\|_F^2 = \text{tr}(A^T A) = \mathbb{E}[v^T A^T A v] = \mathbb{E}[\|Av\|_2^2]$$

where $v \sim \mathcal{N}(0, I)$

Singular values are non-negative values of a matrix obtained after decomposition. Any matrix could be decomposed into the form $U\Sigma V^*$ where U and V are unitary values. Σ is a diagonal matrix consisting of Singular values. These values in the Σ diagonal matrix are listed in decreasing order. The maximum i.e., σ_1 is the norm of the matrix.

What is the difference between these norms?

To recap, norms are used to identify how much a function 'f' can lengthen a matrix. When we are provided with a matrix for a shape we could identify by how much the shape could distort using norm.



There are three types of norms:

- Induced norm: which measures what is the maximum of $\frac{\|Ax\|}{\|x\|}$ for any $x \neq 0$ (or, equivalently, the maximum of $\|Ax\|_1$ for $\|x\|=1$).
- Element-wise norm: which is like unwrapping the matrix A into a long vector, then calculating its vector norm
- Schatten-norm, which measures the vector norm of singular values of A .

How each norm is different?

Frobenius norm = Element-wise 2-norm = Schatten 2-norm

Induced 2-norm = Schatten ∞ -norm. aka Spectral norm.

7.1.4 Properties of a matrix

Trace of a square matrix

Let $A \in \mathbb{R}^{n \times n}$ be a square matrix. The trace of a matrix is the sum of the diagonal elements in the matrix. It is denoted as: $\text{tr}(A)$, and defined as:

$$\text{tr}(A) \triangleq \sum_{i=1}^n A_{ii}$$

Trace has the following properties:

$$\text{tr}(A) = \text{tr}(A^T)$$

$$\text{tr}(A+B) = \text{tr}(A) + \text{tr}(B)$$

$$\text{tr}(c) = c \text{tr}(A)$$

$$\text{tr}(AB) = \text{tr}(BA)$$

$$\text{tr}(A) = \sum_{i=1}^n \lambda_i \text{ where } \lambda_i \text{ are the eigenvalues of } A.$$

$$\text{tr}(ABC) = \text{tr}(BCA) = \text{tr}(CAB) - \text{Cyclic Permutation Property} *$$

Using cyclic permutation property, we can have the following trace trick with scalars (x, x^T) as follows:

$$x^T Ax = \text{tr}(x^T Ax) = \text{tr}(xx^T A)$$

Hutchinson Trace Estimator

For large matrices, it may be expensive to compute trace. But we can compute efficiently through matrix-vector products Ax , where ' v ' is a vector sampled randomly from uniform distribution. Here we could use Monte carlo approximation to $\text{tr}(A)$ using the following identity:

$$\text{tr}(A) = \text{tr}(AIE(vv^T)) = E[\text{tr}(Avv^T)] = E[\text{tr}(v^T Av)]$$

Definite Matrices

A symmetric matrix $M \in \mathbb{R}^{n \times n}$ could be classified into one of the 4 types of definite matrices

Positive-definite matrix: M is positive-definite if $x^T M x > 0$ $\forall x \in \mathbb{R}^n \setminus \{0\}$

$x \neq 0$ but all other numbers

Positive - semidefinite matrix: M is positive-semidefinite if $x^T M x \geq 0$ $\forall x \in \mathbb{R}^n$

Negative - definite matrix: M is negative-definite if $x^T M x < 0$ $\forall x \in \mathbb{R}^n \setminus \{0\}$

Negative - semidefinite matrix: M is negative-semidefinite if $x^T M x \leq 0 \forall x \in \mathbb{R}^n$

* You can also determine definiteness through eigenvalues.

Positive-definite if all eigenvalues are positive, Positive-semidefinite if all

eigenvalues are non-negative, Negative-definite if all eigenvalues are negative,

Negative-semidefinite if all eigenvalues are non-positive, Indefinite if the matrix has mix of positive and negative eigen values.

Determinant of a square matrix

When mathematicians were studying matrices, they were looking for a simple heuristic which capture the characteristics of a matrix. They were looking for a **determining** function for a matrix which now we refer as **determinant**.

A determinant of a 2×2 matrix given below is computed as

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \quad |A| = ad - bc$$

A determinant of a matrix $|A| \neq 0$ suggest that for a given system of linear equations, there exists a unique solution.

If a determinant $|A|=0$, then the matrix doesn't have an inverse.

A determinant can enable us to measure change in unit volume of a shape through linear transformation. The details of this heuristic are obscure at this point of time. But it just signifies the important of determinant.

Properties

$$|A| = |A^T|$$

$$|cA| = c|A|$$

$$|AB| = |A||B|$$

$$|A| = 0 \text{ iff } A \text{ is singular}$$

$$|A^{-1}| = 1/|A| \text{ if } A \text{ is not singular}$$

$$|A| = \prod_{i=1}^n \lambda_i \quad \lambda_i \rightarrow \text{eigenvalues of } A$$

What does it mean a matrix is Singular?

- A square matrix which is not invertible is called Singular Matrix
- Simple check if $|A|=0$ then A is singular

For a positive definite matrix, we can transform

$$A = L L^T \rightarrow \text{Cholesky Decomposition}$$

Lower triangular Upper triangular
Cholesky decomp Cholesky decomp

Because $|A| = |L|^2$ property

$$\left\{ \begin{array}{l} |A| = |L||L^T| = |L|^2 \\ \log|A| = 2 \log|L| = 2 \log \prod_i \lambda_i \\ = 2 \operatorname{trace}(\log(\operatorname{diag}(L))) \end{array} \right.$$

Rank of a matrix

- Ranks tells us whether we may have a chance at solving the system of linear equations.
 - Column Rank is the dimension of the space spanned by its columns.
 - Row Rank is the dimension of the space spanned by its rows
- What is the meaning of spanned by its row or columns?
- Consider this example matrix

$$\begin{bmatrix} 1 & 2 & 3 \\ 3 & 6 & 9 \end{bmatrix}$$

Here the row rank is 1
column rank is 1

How?

Row rank: The 2nd row is duplicate of 1st
So only 1 independent row

Column rank: The 3rd column is 3 times the
1st column. The 2nd column is
2 times the 1st column.

So only 1 independent column.

Recall "Span" definition:

Span is the set of linearly independent vectors
Hence, the rank of column is the length of the column span
and similarly the rank of rows is the length of row span.

- For any given matrix, the rank of column = rank of row. (As seen in above example)
- Hence, we generally denote simply rank(A) for a matrix.

Properties: (Let $A \in \mathbb{R}^{m \times n}, B \in \mathbb{R}^{n \times p}$)

- $\text{rank}(A) \leq \min(m, n)$. If $\text{rank}(A) = \min(m, n)$, then A is referred as "full rank"
- $\text{rank}(A) = \text{rank}(A^T) = \text{rank}(A^T A) = \text{rank}(A A^T)$ or "Rank deficient"
In my defenit?
- $\text{rank}(AB) \leq \min(\text{rank}(A), \text{rank}(B))$
- $\text{rank}(A+B) \leq \text{rank}(A) + \text{rank}(B)$

One can show that square matrix is invertible iff it is full rank
A matrix is said to be invertible if \exists Invertible matrix suggest that there
 $A \cdot A^{-1} = I$ exist a unique solution.

Condition Numbers

Condition numbers measures how numerically stable any computations involving A (our matrix) will be.

What does it mean numerically stable?

- Let's consider a matrix 'A'. If we multiply it with a vector 'x', then we have a new matrix B

$$Ax = B$$

- Let's change the x by a small amount. Normally, what we expect is the B also changes by a small amount
- However, instead if the B changes by a large value then A is deemed to be numerically unstable
- A numerically stable matrix have outputs proportional to the change in perturbations.

Conditional number is denoted as $\kappa(A)$

$$\kappa(A) \triangleq \|A\| \cdot \|A^{-1}\|$$

where $\|A\|$ is the norm of matrix.

If $\kappa(A) \leq 1$, then it is well-conditioned

If $\kappa(A)$ is large, then it is ill-conditioned

A larger $\kappa(A)$ indicates it is nearly singular, whereas matrix with $\kappa(A)$ close to 1 is far from singular.

For example:

$$\text{let } A = 0.1 I_{100 \times 100}$$

the determinant

$$|A| = 10^{-100} \quad (\text{which means close to being singular})$$

but $\kappa(A) = 1$ (because scaling by 10 would lead to non-singular matrix)

If norm is euclidean-norm (aka L_2 -norm), then $\kappa(A)$ could also be computed as the ratio of large & small singular values

$$\kappa(A) = \frac{\sigma_{\max}}{\sigma_{\min}} = \sqrt{\frac{\lambda_{\max}}{\lambda_{\min}}} \quad \lambda \text{ eigenvalues}$$

If we recall, there are different types of norms. Condition number changes based on the type of norm

In real world, singular matrices are rare - meaning have low probability. But singular matrices acts as a boundary for positive & negative determinant matrices

A well-conditioned matrix A of an elliptical shape preserves the shape. Whereas an ill-conditioned matrix skews the unit sphere into a long thin cigar or needle shape.

Note: the quadratic objective function for the above ellipsoid shape is

$$f(x) = x^T A x$$

If $\kappa(A) = 1$, the shape will stay circular.

If $\kappa(A) \approx 10^3$, the shape could be a line or skewed ellipse

1.1.5 Special Type of Matrices

Diagonal Matrix

A matrix where all non-diagonal elements are 0.

Denoted as $D = (d_1, \dots, d_n)$

$$D = \begin{bmatrix} d_1 & & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_n \end{bmatrix} \quad \text{Ex: } \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \rightarrow \text{Identity matrix}$$

Identity matrix:

$I \in \mathbb{R}^{n \times n}$ — square matrix

$I = \text{diag}(1, 1, 1, \dots)$

$$AI = IA = A$$

Block Diagonal

Contains matrices on its main diagonal

$$\text{Eg: } \begin{pmatrix} A & 0 \\ 0 & B \end{pmatrix}$$

Band-Diagonal matrix

Contains entries only along the diagonal and k sides (width) of diagonal. For a tridiagonal 6×6 matrix ($k=3$)

$$\begin{bmatrix} A_{11} & A_{12} & 0 & 0 & 0 & 0 \\ A_{21} & A_{22} & A_{23} & 0 & 0 & 0 \\ 0 & A_{32} & A_{33} & A_{34} & 0 & 0 \\ 0 & 0 & A_{43} & A_{44} & A_{45} & 0 \\ 0 & 0 & 0 & A_{54} & A_{55} & A_{56} \\ 0 & 0 & 0 & 0 & A_{65} & A_{66} \end{bmatrix} \quad \text{Labeled } k=3$$

Triangular Matrices

An upper triangular matrix has non-zero entries above diagonal
 { vice versa for lower triangular matrix.

In triangular matrices, the diagonal elements are the eigenvalues of A. Therefore, the $|A| = \prod_i A_{ii}$

Positive Definite Matrices

Given $A \in \mathbb{R}^{n \times n}$ and $x \in \mathbb{R}^n$ (vector), the scalar obtained through the product $x^T A x$ is called **Quadratic Form**

$$x^T A x = \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j$$

Property:

$$x^T A x = (x^T A x)^T = x^T A^T x = x^T \left(\frac{1}{2} A + \frac{1}{2} A^T\right) x$$

Assumption

Matrices appearing often in quadratic form are symmetric

Positive Definite : $x^T A x > 0 \quad \forall x \in \mathbb{R}^n$

Negative Definite : $x^T A x < 0 \quad \forall x \in \mathbb{R}^n$

Indefinite : $x^T A x$ neither +ve nor -ve

for varying x values i.e., $x_1^T A x_1 < 0$ and $x_2^T A x_2 > 0$

Positive Semi-definite : $x^T A x \geq 0 \quad \forall x \in \mathbb{R}^n$

Negative Semi-definite : $x^T A x \leq 0 \quad \forall x \in \mathbb{R}^n$

- If A is positive definite, then $-A$ is negative definite {Vice versa}

- Similarly A is semi-positive definite, then $-A$ is negative semi-definite {Vice versa}

- If all eigen values are +ve, then A is +ve definite.

- If all elements are +ve, it does not mean it is +ve definite.

A +ve definite matrix can have -ve entries.

- Positive definite matrix diagonal values are larger than sum of the elements of row or column element. +ve definite are

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}| \quad \text{diagonally dominant}$$

- In 2d, any symmetric matrix is true definite iff $a>0, d>0$, and $ad>b^2$
for the matrix $\begin{pmatrix} a & b \\ b & d \end{pmatrix}$

Gram Matrix

Any matrix $A \in \mathbb{R}^{m \times n}$, the Gram matrix

$G = A^T A$ is always positive semidefinite
if $m \geq n$ and A is full rank ($\text{rank}(A) = \min(m, n) \rightarrow$ completely independent vectors in \mathbb{R}^n)
then $G = A^T A$ is true definite

Orthogonal Matrices

Two vectors $x, y \in \mathbb{R}^n$ are orthogonal if

$$x^T y = 0$$

A vector is normalized if

$$\|x\|_2 = 1 \quad (\text{normalized } x)$$

A set of vectors that is pairwise orthogonal and normalized is called **orthonormal**

A square matrix U is **orthogonal** if all its columns are orthonormal

Unitary: If the entries of U are complex then we call U as **Unitary**

U is orthogonal iff

$$U^T U = U U^T = I$$

If U is not square and $U \in \mathbb{R}^{n \times m}$, where $n < m$ and columns are orthonormal then

$$U^T U = I, \text{ but } U U^T \neq I$$

We call orthogonal only for square matrix

Euclidean Norm:

$$\|Ux\|_2 = \|x\|_2 \quad x \in \mathbb{R}^n \setminus \{0\}$$

Angle between two vectors are always preserved after transformed by an orthogonal matrix.

$$\cos(\alpha(x, y)) = \frac{x^T y}{\|x\| \|y\|} \quad \text{so, } \cos(\alpha(Ux, Uy)) = \frac{(Ux)^T (Uy)}{\|Ux\| \|Uy\|} = \frac{x^T y}{\|x\| \|y\|} = \cos(\alpha(x, y))$$

Transformations by orthogonal matrix are generalizations

Rotations if $|U| = 1$

Reflections if $|U| = -1$

Tip: Any square matrix can be transformed into orthogonal matrix using a technique called Gram Schmidt Orthogonalization.

7.2 Matrix Multiplication

Let $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{n \times p}$ and the product of the two matrices is given as: $C = AB \in \mathbb{R}^{m \times p}$

$$\text{where, } C_{ij} = \sum_{k=1}^n A_{ik} B_{kj}$$

- Time complexity: $O(mnp)$

- Matrix multiplication can be parallelized-

- GPUs and TPUs perform faster.

Properties:

$$\text{Associative } (AB)C = A(BC)$$

$$\text{Distributive } A(B+C) = AB+AC$$

$$\text{Not commutative } AB \neq BA$$

7.2.1 Vector-Vector products

Two vectors $x, y \in \mathbb{R}^n$, the quantity is called
inner product, dot product or scalar product

$$\langle x, y \rangle \triangleq x^T y = \sum_{i=1}^n x_i y_i$$

Dot products are commutative $x^T y = y^T x$

Given vectors $x \in \mathbb{R}^m, y \in \mathbb{R}^n$, xy^T is called outer product

It is a matrix whose entries are given by $(xy^T)_{ij} = x_i y_j$

$$xy^T \in \mathbb{R}^{m \times n} = \begin{bmatrix} x_1 y_1 & x_1 y_2 & \cdots & x_1 y_n \\ x_2 y_1 & x_2 y_2 & \cdots & x_2 y_n \\ \vdots & \vdots & \ddots & \vdots \\ x_m y_1 & x_m y_2 & \cdots & x_m y_n \end{bmatrix}$$

Matrix - Vector products

Given a matrix $A \in \mathbb{R}^{m \times n}$, $x \in \mathbb{R}^n$, the product between the two will be: $y = Ax \in \mathbb{R}^m$

$$m \begin{array}{|c|} \hline \boxed{\textcolor{green}{y}} \\ \hline \end{array} = m \begin{array}{|c|} \hline \textcolor{red}{\text{|||||}} \\ \hline A \\ \hline n \\ \hline \end{array} \begin{array}{|c|} \hline \textcolor{blue}{\text{|||}} \\ \hline x \\ \hline n \\ \hline \end{array}$$

y can be viewed as the linear combination of the columns of A where the coefficients of this combination are x :

$$y = Ax = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} x_1 + \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} x_2 + \dots + \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} x_n$$

These columns can be viewed as a set of **basis vectors** defining a linear subspace.

Basis vectors are a set of linearly independent vectors spanning whole space.

Matrix - Matrix Products

$$C = AB \quad \text{where } A \in \mathbb{R}^{m \times n}, B \in \mathbb{R}^{n \times p}$$

vector-vector inner product:

$$\text{Let } A = \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix} \quad B = \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix}$$

$$C = AB = \begin{bmatrix} a_1^T \\ a_2^T \end{bmatrix} \begin{bmatrix} b_1 & b_2 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} a_1^T b_1 & a_1^T b_2 \\ a_2^T b_1 & a_2^T b_2 \end{bmatrix}$$

$$a_1^T b_1 = \begin{bmatrix} 1 & 3 \end{bmatrix} \begin{bmatrix} 5 \\ 7 \end{bmatrix} = 5 + 21 = 26 \quad \begin{bmatrix} 26 & 30 \\ 38 & 44 \end{bmatrix}$$

$$a_2^T b_1 = \begin{bmatrix} 2 & 4 \end{bmatrix} \begin{bmatrix} 5 \\ 7 \end{bmatrix} = 10 + 28 = 38$$

$$a_1^T b_2 = \begin{bmatrix} 1 & 3 \end{bmatrix} \begin{bmatrix} 6 \\ 8 \end{bmatrix} = 6 + 24 = 30$$

$$a_2^T b_2 = \begin{bmatrix} 2 & 4 \end{bmatrix} \begin{bmatrix} 6 \\ 8 \end{bmatrix} = 12 + 32 = 44$$

Vector - vector outer product:

$$C = AB = \begin{bmatrix} 1 & 1 & \dots \\ a_1 & a_2 & \dots \\ 1 & 1 & \dots \end{bmatrix} \begin{bmatrix} -b_1^T \\ -b_2^T \\ \vdots \end{bmatrix} = \sum_{i=1}^n a_i b_i^T$$

Matrix - vector product:

$$C = AB = A \begin{bmatrix} 1 & 1 & \dots \\ b_1 & b_2 & \dots \\ 1 & 1 & \dots \end{bmatrix} = \begin{bmatrix} 1 & 1 & \dots \\ A b_1 & A b_2 & \dots \\ 1 & 1 & \dots \end{bmatrix}$$

$C = Ab_i$

Similarly,

$$C = AB = \begin{bmatrix} -a_1^T \\ -a_2^T \\ \vdots \end{bmatrix} B = \begin{bmatrix} -a_1^T B \\ -a_2^T B \\ \vdots \end{bmatrix}$$

$$C_i^T = a_i^T B$$

$$A^2 = A \cdot A$$

$$A^{02} = [A_{ij}^2] \rightarrow \text{elementwise square}$$

If A is a diagonal matrix, then

$$A^2 = A^{02}$$

Application: Manipulating data matrices

Given $X \in N \times D$ matrix

Sum across rows:

$$I_N^T X = (\sum_n x_{n1}, \dots, \sum_n x_{nD})$$

Sum across columns:

$$X I_D = \begin{pmatrix} \sum_d x_{1d} \\ \vdots \\ \sum_d x_{Nd} \end{pmatrix}$$

Sum all entries:

$$I_N^T X I_D = \sum_{ij} x_{ij}$$

Scaling Rows and Columns

Scaling Rows using $\text{diag}(s)$ matrix where 's' is a scaling factor given by $\text{diag}(s)x$

Scaling Columns using $\text{diag}(s)$ as $x \text{diag}(s)$

Remember, scaling is generally used for standardizing a matrix (or a.k.a normalization). To standardize a matrix:

$$\text{standardize}(x) = (x - \bar{x}\mu^T) \text{diag}(\sigma)^{-1} \rightarrow \text{why inverse?}$$

where $\bar{x}\mu = \bar{x}\bar{x}$ (mean), σ is standard deviation of vectors.

Sum of squares Used in 4.59, 7.236, 11.6

The sum of squares matrix for $x \in \mathbb{R}^{D \times D}$

$$S_0 \triangleq x^T x = \sum_{n=1}^N x_n x_n^T$$

Scatter matrix

$$S_{\bar{x}} \triangleq \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T = \left(\sum_n x_n x_n^T \right) - N \bar{x} \bar{x}^T$$

$$\text{Mean } \bar{x} = \frac{1}{N} x^T \mathbf{1}_N$$

Here we are computing mean centered data (\tilde{x}) from x

We can compute the centered data matrix using

$$\tilde{x} = x - \bar{x}\bar{x}^T = x - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T x = C_N x$$

where

$$C_N \triangleq \mathbf{1}_N \mathbf{1}_N^T - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T$$

C_N is the centering matrix

Scatter matrix can now be computed as:

$$S_{\bar{x}} = \tilde{x}^T \tilde{x} = x^T C_N^T C_N x = x^T C_N x$$

What is the use of Scatter matrix?

- $S_{\bar{x}}$ is the unnormalized covariance matrix
- Covariance matrix

$$C = \frac{1}{N-1} S_{\bar{x}}$$

Example:

$$c_0 = [0]$$

$$c_1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

$$c_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} - \frac{1}{3} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} \frac{2}{3} & -\frac{1}{3} & -\frac{1}{3} \\ -\frac{1}{3} & \frac{2}{3} & -\frac{1}{3} \\ -\frac{1}{3} & -\frac{1}{3} & \frac{2}{3} \end{bmatrix}$$

Gram Matrix

A matrix of inner products for x is called Gram Matrix (K)

$$K \triangleq x x^T$$

Inner products of mean-centered data vectors

$$\tilde{K} = \tilde{x} \tilde{x}^T$$

- ' K ' is the feature similarity matrix. If we're only given K and we would like to model the data, then first we need to **normalize** the data.
- In cases, where x is not provided and only ' K ' provided, we can normalize the data without \tilde{x} through centering matrix ' C '
- **Double centering** is a trick which involves multiplying centering matrix on both sides to achieve zero mean in the data. We double center $K \rightarrow \tilde{K}$ as follows

$$\tilde{K} = \tilde{x} \tilde{x}^T = C_N K C_N$$

Distance Matrix

Let $x \in \mathbb{R}^{N \times P}$ $y \in \mathbb{R}^{N \times P}$, the distance between the elements is computed as

$$D_{ij} = (x_i - y_j)^T (x_i - y_j) = \|x_i\|^2 - 2x_i^T y_j + \|y_j\|^2$$

In matrix, $\|x_i\|^2$ can be computed by $\text{diag}(x x^T) = \hat{x}$

$$\text{Similarly, } \|y_j\|^2 = \text{diag}(y y^T) = \hat{y}$$

Finally, the **distance matrix** can be computed as follows:

$$D = \hat{x} I_{Ny} - 2x y^T + I_N \hat{y}^T$$

How I_{Nx} and I_{Ny} is interpreted?

If $x = y$, then

$$D = \hat{x} I_N - 2x x^T + I_N \hat{x}^T$$

7.2.5 Kronecker Products

If $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{p \times n}$, the Kronecker product is given as
 $A \otimes B \in \mathbb{R}^{mp \times n}$

$$A \otimes B = \begin{bmatrix} a_{11}B & \cdots & a_{1n}B \\ \vdots & \ddots & \vdots \\ a_{m1}B & \cdots & a_{mn}B \end{bmatrix}$$

Check textbook for example

Properties:

$$(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$$

$$(A \otimes B) \text{vec}(C) = \text{vec}(BCA^T)$$

7.2.6 Einstein Summation

a.k.a einsum is a short-hand notation for specifying operations on matrix, vector, and tensors.

For example :

$$C_{ij} = \sum_k A_{ik} B_{kj}$$

Can be written as

$$C_{ij} = A_{ik} B_{kj}$$

We just dropped \sum_k

In programming, for example in numpy

`C = np.einsum('ik,kj->ij', A,B)`

Similar operations can be performed on tensor. Consider a tensor which has a batch (B) of sentences (M) where each sentence has word embeddings of size (D). To illustrate

Batch	Sentence 1	Sentence 2	Sentence 3
	word ₁ [0 0 1 0] word ₂ [0 1 0 0] word ₃ [1 0 0 0]		
		word ₁ [0 1 0 0] word ₂ [0 0 0 1] word ₃ [0 0 1 0]	

Let $W \in \mathbb{R}^{D \times P}$ is a word embedding matrix. We can multiply as :

$$F_{BMP} = S_{BMD} W_{DP}$$

7.3 Matrix Inversion

For matrix $A \in \mathbb{R}^{n \times n}$, the inverse is A^{-1}

$$A^{-1}A = AA^{-1} = I$$

A^{-1} exists iff $\det(A) \neq 0$. If $\det(A)=0$, then A is **singular matrix**.

Properties

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

$$(AB)^{-1} = B^{-1}A^{-1}$$

$$(A^{-1})^T = (A^T)^{-1} \triangleq A^{-T}$$

For $A \in \mathbb{R}^{2 \times 2}$, the A^{-1} will be

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}, \quad A^{-1} = \frac{1}{|A|} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$

For block diagonal

$$\begin{pmatrix} A & 0 \\ 0 & B \end{pmatrix}^{-1} = \begin{pmatrix} A^{-1} & 0 \\ 0 & B^{-1} \end{pmatrix}$$

7.4 Eigenvalue Decomposition (EVD)

Let $A \in \mathbb{R}^{n \times n}$, then the eigenvalue of A is denoted as $\lambda \in \mathbb{R}$ and eigenvector of A will be $u \in \mathbb{R}^n$, only if

$$Au = \lambda u, \quad u \neq 0$$

Intuitively what this means is for a given vector ' u ' if we multiply with a matrix we will get a new vector of the same size of u but only scaled by a factor of λ .

If A is a rotation matrix, u is the axis of rotation, then multiplying the two matrices will yield a new scale for axis of rotation.

The scalar component $c \in \mathbb{R}$ for eigenvector u :

$$A(cu) = cAu = c\lambda u = \lambda cu$$

We usually assume eigenvector is normalized to have length 1.

(λ, u) is an eigenvalue - eigenvector pair of A only if

$$(\lambda I - A)u = 0 \rightsquigarrow \text{How? } \lambda I u - A u = 0 \\ (\lambda \neq 0) \quad \quad \quad A u = \lambda I u$$

$(\lambda I - A)u = 0$ will have non-zero solution to u iff $(\lambda I - A)$ has a non-empty null space, which will be the case if $(\lambda I - A)$ is singular.

i.e., $\det(\lambda I - A) = 0$ Remember: we're saying singular only for $\lambda I - A$. However, A still needs to

↓
- This called characteristic equation. be non-singular i.e. $|A| \neq 0$ of A .

- The n solutions of this eqⁿ are n eigenvalues λ_i and u_i are corresponding eigenvectors
- It is standard to sort the eigenvalues by their magnitude.

} Somewhat closely
reminds of why we consider only the first K principal components of a matrix to represent the whole matrix

Properties:

$$\text{tr}(A) = \sum_{i=1}^n \lambda_i$$

$$\det(A) = \prod_{i=1}^n \lambda_i$$

$$\text{rank}(A) = \text{nz}(\lambda)$$

eigenvalue of A^{-1} = $1/\lambda_i$ and eigenvector $A^{-1}u_i = 1/\lambda_i u_i$

In triangle matrix or diagonal matrix, eigenvalues are the diagonal values.

Diagonalization

We can write all the eigenvector equations as

$$AU = U\Lambda \rightsquigarrow \text{How? Note: previously we have realized there could be}$$

where $U \in \mathbb{R}^{n \times n}$

$$U = \begin{bmatrix} | & | & | & | \\ u_1 & u_2 & u_3 & \dots & u_n \\ | & | & | & | \end{bmatrix} \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$$

many equations from the matrix.

If the eigenvectors of A are linearly independent, then matrix U will be invertible (which is basically going other side as inverse form)

$$A = U\Lambda U^{-1} \rightsquigarrow \text{Any matrix in this form is called Diagonalizable.}$$

7.4.3 Eigenvalues and Eigenvectors of Symmetric Matrices

If A is a symmetric matrix, then the eigenvectors can be shown as orthonormal i.e.

$$u_i^T u_j = 0 \text{ iff } i \neq j$$

$$u_i^T u_i = u_i \cdot u_i = 1 \text{ iff } i = j \text{ (diagonal elements)}$$

Therefore, $U^T U = U^T U = I \Rightarrow$ Hence, U is orthogonal matrix

Note previously we represented $A = U \Lambda U^T$

Now we are saying A can also be represented as

$$\begin{aligned} A &= U \Lambda U^T \\ &= \begin{pmatrix} 1 & 1 & \dots & 1 \\ u_1 & u_2 & \dots & u_n \\ 1 & 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{pmatrix} \begin{pmatrix} -u_1^T \\ -u_2^T \\ \vdots \\ -u_n^T \end{pmatrix} \\ &= \sum_{i=1}^n \lambda_i u_i u_i^T \end{aligned}$$

Here the assumption is $U^T = U^{-1}$.

With that assumption, the A^{-1} will be

$$A^{-1} = U \Lambda^{-1} U^T = \sum_{i=1}^n \frac{1}{\lambda_i} u_i u_i^T$$

Checking For Positive Definiteness

Recall positive definiteness is calculated by

$$x^T A x > 0$$

Similarly we can substitute A with eigenvalue & eigenvector

$$x^T A x = x^T U \Lambda U^T x \quad \text{let } y = U^T x \quad \therefore y^T = x^T U$$

$$= y^T \Lambda y$$

$$= \sum_{i=1}^n \lambda_i y_i^2$$

$\because y^2$ will always be positive, the sign of λ_i determines the definiteness of matrix. i.e., if all (λ_i) $\lambda_i > 0$ +ve definite $\lambda_i \geq 0$ the semidefinite.
 $\lambda_i < 0$ -ve definite $\lambda_i \leq 0$ -ve semidefinite
if λ 's contains the 0's then indefinite matrix.

7.4.4 Geometry of Quadratic Forms

Quadratic form is a function written for a matrix A as

$$f(x) = x^T A x$$

We know $A = U \Lambda U^T$, so

$$f(x) = x^T A x = x^T U \Lambda U^T x = y^T \Lambda y = \sum_{i=1}^n \lambda_i y_i^2$$

The level sets of $f(x)$ define hyper-ellipsoids

For example: $\lambda_1 y_1^2 + \lambda_2 y_2^2 = r$ (eqn for a 2d ellipse)

Eigenvectors determine orientation of the ellipse and eigenvalues determine how elongated it is.

7.4.5 Standardizing and whitening data

- Standardizing is a preprocessing technique applied to the data so that each column has zero mean and unit variance.
- Standardizing data will have unit variance but it does not remove the correlation between the columns.
- To remove correlation, we approach a process called "Whitening".
- Let's say, if we are training on images, the raw input is redundant since adjacent pixel values are highly correlated.
- The goal of whitening is to make the input less redundant i.e., features are less correlated and have same variance.

Let the covariance matrix of data be

$$\Sigma = \frac{1}{N} X^T X$$

Now Σ can be diagonalized as

$$\Sigma = E D E^T$$

which is basically decomposing any matrix into eigenvectors & eigenvalues

Parallelly, we decompose X using SVD

$$X = U S V^T$$

When we compare the decomposed forms,

$$E = V \quad \text{and} \quad D = S^2 \quad \text{How?}$$

We define PCA whitening matrix as:

$$W_{PCA} = D^{-1/2} E^T \quad \text{→ pulled out of thin air}$$

Now we construct a transformed vector

$$y = W_{PCA} Z \quad \text{which will be a de-correlated version of } X$$

To check the covariance of the whitening:

$$\begin{aligned} \text{cov}[y] &= W_{PCA} Z W_{PCA}^T = W_{PCA} \Sigma W_{PCA}^T \\ &= (D^{-1/2} E^T) (E D E^T) (E D^{-1/2}) = I \\ \text{We know } W_{PCA} &= D^{-1/2} E^T \quad \downarrow \quad \downarrow \quad \downarrow \quad \text{This is } W_{PCA}^{-1} \\ \Sigma &= E D E^T \end{aligned}$$

Rotating the whitening matrix will still maintain the whitening property.

If R is rotation, then $W_{ZCA} = R W_{PCA}$ is the new rotated matrix. The whitening property $W^T W = \Sigma^{-1}$ is preserved.

$$\text{We know } W_{PCA} = D^{-1/2} E^T$$

$$\text{Let's assume we rotated } W_{ZCA} = R W_{PCA} = R D^{-1/2} E^T$$

$$\begin{aligned} \text{Let's also assume } R = E \quad W_{ZCA} &= E D^{-1/2} E^T = \Sigma^{-1/2} = V S^{-1} V^T \\ &\quad \downarrow \quad \downarrow \\ &\quad \Sigma = E D E^T \end{aligned}$$

This is called Mahalanobis whitening or ZCA (zero phase component analysis)

- When applied to images the ZCA transformed vectors will still look like images.

7.4.6 Power Method

Power method is a simple iterative method to compute the largest eigenvalue and the corresponding vector for a matrix.

Let A be a matrix with orthonormal eigenvectors and eigenvalues. We know $|d_1| > |d_2| > \dots > |d_m| \geq 0$.

So $A = U\Lambda U^T$. Let v_0 be a random vector within the same range as A . We can assume that $Ax_0 = v_0$ for some suitable x . Hence we can write v_0 in the form of A as:

$$v_0 = U(\Lambda U^T x) = a_1 u_1 + \dots + a_m u_m \quad (\text{Here } a_i \text{ are scalars } \in \mathbb{R})$$

$\Lambda \in \mathbb{R}^{m \times m}$

Now here comes the trick. If we repeatedly multiply v with A at certain time ' t ', we would have $v_{t+1} \xrightarrow{\text{This is norm not modulus}} v_t$ i.e., $|v_{t+1} - v_t| \leq \epsilon$ where ϵ is often known as 'tolerance' is a very small value e.g.: $\epsilon = 10^{-5}$. Formally

$$v_t \propto Av_{t-1} \quad (\text{Note: we normalize } v \text{ at each time step after multiplying with } A)$$

$$v_0 \propto a_1 d_1^0 u_1 + a_2 d_2^0 u_2 + \dots + a_m d_m^0 u_m$$

$$v_1 \propto a_1 d_1^1 u_1 + a_2 d_2^1 u_2 + \dots + a_m d_m^1 u_m$$

$$\vdots \quad \vdots \quad \vdots \quad \vdots$$

$$v_t \propto a_1 d_1^t u_1 + a_2 d_2^t u_2 + \dots + a_m d_m^t u_m$$

$$= d_1^t (a_1 u_1 + a_2 (d_2/d_1)^t u_2 + \dots + a_m (d_m/d_1)^t u_m)$$

$$= d_1^t a_1 u_1 \quad \rightsquigarrow \text{because } d_1 \geq d_2 \geq d_3 \geq \dots$$

$$\text{and } d_1 \leq d_2 \leq \frac{d_3}{d_1} \leq \dots \leq 1$$

which suggests that a random vector v_0 will eventually converge to u_1 (the largest eigen vector of A) at some timestep ' t '. The only requirement being $v_0^T u_1 \neq 0$, which will be true for random v_0 with high probability.

What does it mean when a random vector v_0 to u_1 is zero?

Well, if you look in terms of cosine similarity = $\frac{a \cdot b}{|a| |b|}$

If the $v_0^T u_1 = 0$, then the vector doesn't belong in the matrix space.

Previously we stated that $Ax = \lambda x$ for some value of x . Let's say we do several iterations of v_t i.e., $v_t \leftarrow Av_{t-1}$. For some timestep t , we have v_t such that

$$v_t \approx \lambda \cdot x \text{ and remember } v_t \leftarrow A^t a_i u_i$$

where x is our approximation of a dominant eigenvector of A .

We now use Rayleigh quotient to find the eigenvalue. The Rayleigh Quotient is defined as

$$R(A, x) \triangleq \frac{x^T Ax}{x^T x}$$

Hence for u_i as x

$$R(A, u_i) = \frac{u_i^T A u_i}{u_i^T u_i} = \frac{u_i^T \cancel{\lambda u_i}}{u_i^T u_i} = \frac{\lambda_i u_i^T u_i}{u_i^T u_i} = \lambda_i$$

\rightarrow we know $Ax = \lambda x$
so $Au_i = \lambda u_i$

1.4.1 Deflation

Using power method, we obtain first eigenvector and eigenvalue. In order to obtain the following values we approach the deflation trick.

$$\text{For } (\lambda_1, u_1) \Rightarrow B = (I - u_1 u_1^T) A$$

$$= A - u_1 u_1^T A$$

$$= A - \lambda_1 u_1 u_1^T$$

Note: $u_1 u_1^T$ is an outer product which will yield a matrix of the same size(A).

Using the matrix ' B ', we apply power method to compute λ_2 and u_2 .

Similarly for $(\lambda_3, u_3), (\lambda_4, u_4), \dots$

$$\text{For } (\lambda_3, u_3) \Rightarrow C = B - \lambda_2 u_2 u_2^T$$

$$\text{For } (\lambda_4, u_4) \Rightarrow D = C - \lambda_3 u_3 u_3^T$$

& so on.

Deflation technique is used to implement PCA and Sparse PCA.

7.4.8 Eigenvectors optimize quadratic forms

This is where we dive into matrix calculus to solve an optimization problem.

Here is the following constrained optimization problem:

$$\max_{x \in \mathbb{R}^n} x^T A x \quad \text{subject to } \|x\|_2^2 = 1$$

why $\|x\|_2^2$?

One way to solve such optimization problem $\|x\|_2$ is the L2-norm. $\|x\|_2=1$ indicates is through **Lagrangian Optimization**.

To maximize/minimize a function $f(x)$ sides does not change the fact. We'll have which is subject to $g(x)$, the Lagrange $\|x\|_2^2 = 1 \Rightarrow x^T x = 1$ (we know $x^T x$ is matrix squared)

is defined as

What happens if choose $\|x\|_2^2 > 1$ (like 100)?

$$L(x, \lambda) \triangleq f(x) + \lambda g(x)$$

This is called as **Lagrange multiplier** but I guess it just complicates things.

To solve the quadratic form, the Lagrangian is given as:

$$L(x, \lambda) = x^T A x + \lambda (1 - x^T x)$$

$f(x) \quad g(x)$

Thereby, we can establish that x^* to be an optimal point to the problem. That is by finding the point where the gradient of Lagrangian is at zero. Formally given as

$$\nabla_x L(x, \lambda) = 2A^T x - 2\lambda x = 0$$

The above equation is familiar

Remember that,

Eigenvectors (A) =

$$2A^T x - 2\lambda x = 0$$

Eigenvectors (A^T)

$$2A^T x = 2\lambda x$$

$$A^T x = \lambda x$$

The only points which could maximize or minimize A are its eigenvectors.

Hence, we are essentially finding the eigenvectors that maximize or minimize through Lagrangian optimization. We don't need to perform Lagrangian optimization to find the vectors responsible for maximizing or minimizing the quadratic form. This is more of a theoretical confirmation.

But also what other forms we might encounter which are other than ellipsoid shapes? For such forms Lagrangian is one way to find the vectors responsible.

7.5 Singular Value Decomposition (SVD)

So far we have discussed Eigen Value Decomposition (EVD) for square matrices. Now we will look into generalizing this decomposition for rectangular matrices.

7.5.1 Basics

A matrix $A \in \mathbb{R}^{m \times n}$ can be decomposed as

$$A = USV^T$$

$$= \sigma_1(u_1)(-v_1^T) + \sigma_2(u_2)(-v_2^T) + \dots + \sigma_r(u_r)(-v_r^T)$$

- Here $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$ and $S \in \mathbb{R}^{m \times n}$
- The matrix S contains $r = \min(m, n)$ singular values which are $\sigma_i \geq 0$ ($\sigma_i \in S$). The rest all are 0s.
- Columns of U are left singular vectors and columns of V are right singular vectors
- This is called Singular Value Decomposition.

Why don't we call these eigenvalues and what's the difference between them?

If a matrix consists of $m > n$ dimensions, then only n singular values are computed and the rest $m-n$ are ignored i.e., zero-ed. This leads to a simple economy sized SVD known as thin SVD.

Time complexity for SVD: $O(\min(m^2, m \cdot n))$

7.5.2 Connection between SVD and EVD

If A is a real, symmetric, and positive definite matrix, then

Singular values are equal to the eigenvalues

Singular vectors are equal to the eigenvectors

$$\underbrace{A}_{\text{SVD}} = \underbrace{USV^T}_{\text{EVD}} = \underbrace{USU^{-1}}_{\text{EVD}}$$

For $A^T A$ and $A A^T$, we can make following deductions:

$$A^T A = V S^T U^T U S V^T$$

$$A A^T = U S V^T V S^T U^T$$

$$A^T A = V S^T S V^T = V D_n V^{-1}$$

$$A A^T = U S S^T U^T = U D_m U^{-1}$$

$$(A^T A)V = V S^T S = V D_n$$

$$(A A^T)U = U S S^T = U D_m$$

So eigenvectors of $A A^T$ is right singular vector and eigenvalues are D_m (diagonal $n \times n$ matrix)
 Similarly eigenvectors of $A^T A$ is left singular vector and eigenvalues are D_n (diagonal $m \times m$ matrix)

In summary,

$$U = \text{eigenvector } (AA^T)$$

$$V = \text{eigenvector } (A^TA)$$

$$\Omega_m = \text{eigenval } (AA^T)$$

$$\Omega_n = \text{eigenval } (A^TA)$$

For an economy sized SVD, we can define

$$D = S^2 = S^TS = SS^T$$

Note: EVD might not exist for some square matrices, however, SVD does exist for them.

7.5.3 Pseudo Inverse

Remember for a matrix to be invertible (A^{-1}), the matrix has to be a square matrix. Finding inverse of a matrix suggests that there exists a unique solution for the matrix. Similarly, to find inverse for a non-square matrix, pseudo-inverse is computed for a matrix. It is denoted as A^+ .

Properties

$$AA^+A = A$$

$$A^+A A^+ = A^+$$

$$(AA^T)^T = AA^T$$

$$(A^+A^T)^T = A^+A$$

For a matrix $A \in \mathbb{R}^{m \times n}$,

if $m=n$ (i.e., square matrix) then,

$$A^+ = A^{-1}$$

if $m > n$,

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

if $m < n$

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

This A^+ is called Left inverse because when you multiply A on right $A^+A = I$

This is called right inverse when you multiply A on left $AA^+ = I$

Computing A^+ using SVD:

$$A = USV^T \Rightarrow A^+ = V S^T U^T$$

$$A^+ = A^{-1} = (USV^T)^{-1} = V S^{-1} U^T$$

} So computing SVD decomposition and inverting would yield A^+ . Here $S^{-1} = \text{diag}(\sigma_1^{-1}, \dots, \sigma_r^{-1}, 0, \dots, 0)$
rank
computed upto rank

Here A is full-rank matrix
i.e. $\text{rank}(A) = \min(m, n)$.

Full rank \approx linearly independent vectors

7.5.4 SVD and the range and null space of a matrix

Here we show that left and right vectors form an orthonormal basis for the range and null space.

Let's unpack this statement for a bit.

What does it mean vectors form an orthonormal basis?

- Two vectors are orthogonal if their dot products are equal to zero. Basis is the set of linearly independent vectors.

Orthonormal basis is set of linearly independent vectors which are orthogonal and their unit norm is 1. Remember, unit norm, usually for columns of a matrix, is the total length of vectors in Euclidean space. We can also say such matrix is normalized or standardized.

Now what does it mean orthonormal basis for range and null space?

- Recall that Range is the set of columns that can be written as a linear combination to generate a vector.
- Nullspace is the set of columns that map to zero to null space

We know

$$A = USV^T$$

$$Ax = (USV^T)x = \sum_{j: \sigma_j > 0} \sigma_j (v_j^T x) u_j = \sum_{j=1}^r \sigma_j (v_j^T x) u_j \quad \begin{array}{l} r \text{ is the rank of } A \\ \text{Dimension : } r \end{array}$$

Therefore, from this we can say Ax is the linear combination of left singular vectors u_1, \dots, u_r i.e., $\text{range}(A) = \text{Span}(\{u_j : \sigma_j > 0\})$

To find null space, let's consider vector $y \in \mathbb{R}^n$ that is a linear combination solely of right singular vectors for zero singular values

$$y = \sum_{j: \sigma_j = 0} c_j v_j = \sum_{j=r+1}^n c_j v_j$$

Dimension : $n-r$

Therefore $\text{nullspace}(A) = \text{Span}(\{v_j : \sigma_j = 0\})$

$$\text{Dimension}(\text{range}(A)) + \text{Dimension}(\text{nullspace}(A)) = r + n - r = n$$

i.e., rank + nullity = n

This is called rank-nullity theorem.

7.5.5 Truncated SVD

We know

$$A = USV^T$$

If we pick the top k singular values then,

$$A_k = U_k S_k V_k^T$$

The error between A and A_k

$$\|A - A_k\|_2 \approx 0$$

If the error is close to zero, then it is a good approximation.

Usually when the $k = r$ i.e., $\text{Rank}(A)$, then error will be close to zero

If $k < r$ and error is negligible, which means certain non-zero singular values are decaying towards zero. (common in natural data)

We refer $k < r$ decomposed matrix as Truncated SVD.

Total parameters needed to represent a $M \times N$ matrix using a rank k approximation is

K vectors of M + K vectors of N + K singular values

$$MK + NK + K = K(M+N+1) \text{ parameters}$$

For ex: 100×50 matrix has SVD. Let $k = 10$, then

$$10(100+50+1) = 1510 \text{ parameters. Which is } \frac{1510}{3000} = 33.7\% \text{ of original matrix size.}$$

7.6.1 LU Factorization

An matrix can be factorized into upper and lower triangular matrices

$$A = LU$$

In L the upper diagonal values will be zero

In U the lower diagonal values will be zero

However, before applying LU Factorization we may need to apply permute the entries. The reason being, let's assume the following matrix

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{bmatrix}$$

If $a_{11}=0$, then either l_{11} or $u_{11}=0$ which indicates either L or U is singular which means L or U doesn't have unique solution.

To avoid this, the matrix needs to be permuted so that the a_{11} is non-zero. This is repeated for subsequent steps. We denote this process as

$$PA = LU$$

where P is a permutation matrix where $P_{ij}=1$ if row j gets permuted to row i. This is called partial pivoting.

7.6.2 QR Decomposition

Another way to decompose a matrix $A \in \mathbb{R}^{m \times n}$ is by finding a series of orthonormal vectors a_1, a_2, \dots that span successive subspaces of A columns i.e.,

$$\begin{pmatrix} | & | & | \\ a_1 & a_2 & \dots & a_n \\ | & | & | & | \end{pmatrix} = \underbrace{\begin{pmatrix} | & | & | \\ a_1 & a_2 & \dots & a_n \\ | & | & | & | \end{pmatrix}}_{\text{orthonormal vectors}} \underbrace{\begin{pmatrix} r_{11} & r_{12} & \dots & r_{1n} \\ r_{21} & r_{22} & \dots & r_{2n} \\ \vdots & & & \vdots \\ r_{m1} & r_{m2} & \dots & r_{mn} \end{pmatrix}}_{\text{coefficients}}$$

Here we can write

$$a_1 = r_{11}q_1$$

$$a_2 = r_{21}q_1 + r_{22}q_2$$

:

$$a_n = r_{n1}q_1 + \dots + r_{nn}q_n$$

We note as

$$A = \hat{Q}\hat{R} \quad \text{where } \hat{Q} \in \mathbb{R}^{m \times n}, \hat{R} \in \mathbb{R}^{n \times n} \rightarrow \text{Reduced QR or economy sized QR}$$

$$A = QR \quad \text{where } Q \in \mathbb{R}^{m \times (n+m-n)} \text{ i.e., } Q \in \mathbb{R}^{m \times m}, R \in \mathbb{R}^{m \times n} \rightarrow \text{Full QR}$$

This is a square matrix which satisfies $Q^T Q = I$

7.6.3 Cholesky Decomposition

In this decomposition, any matrix can be factorized into
 $A = R^T R$ \sim Also can be written as $A = LL^T$ ($L = RT$)
where R is upper triangle matrix

This looks similar to LU factorization. i.e., $A = LU$. But remember
in LU Factorization, the $L \neq U$.

Cholesky Factorization is also known as Matrix Square Root
Likely?

Time Complexity: $O(V^3)$

Application: Sampling from MVN

Given a matrix $A \in \mathbb{R}^{m \times m}$ and $\Sigma = A^T A \in \mathbb{R}^{m \times m}$.
Here Σ is the covariance of A . We can then decompose this
covariance matrix Σ using cholesky decomposition. Let $L \in \mathbb{R}^{m \times m}$ be
the decomposed matrix.

So we would like to get a random multivariate normal
distribution $y \sim \mathcal{N}(\mu, \Sigma)$ where μ is the random mean $\in \mathbb{R}^m$.
We could totally achieve this by simply random sampling
 $x \sim \mathcal{N}(0, I)$ and then set

$$y = Lx + \mu$$

\downarrow

This yields random multivariate normal distribution
which has the column mean \sim random sample mean
and you can change the dimension of x and μ
while keeping the ' L ' constant to randomly sample
multivariate distributions from this technique.

Cholesky decomposition can capture the desired
correlations faithfully and can be used for generating
correlated data.

- x is random matrix
- L is obtained from covariance matrix decomposition of random matrix.
- The L make sure the features responsible for constructing Σ i.e. (A) are also responsible for generating new matrices with a random vector.
- Which is why covariance of Σ and y are same.
- Now covariance will remain same for any size of y wrt Σ

7.7 Solving systems of linear equations

Let's consider the following equations

$$3x_1 + 2x_2 - x_3 = 1$$

$$2x_1 - 2x_2 + 4x_3 = -2$$

$$-x_1 + 4x_2 - x_3 = 0$$

which can be represented as

$$Ax = b$$

$$A = \begin{pmatrix} 3 & 2 & -1 \\ 2 & -2 & 4 \\ -1 & 4 & -1 \end{pmatrix} \quad b = \begin{pmatrix} 1 \\ -2 \\ 0 \end{pmatrix}$$

$$\text{Solving for } x \text{ yield } x = [1, -2, 0]$$

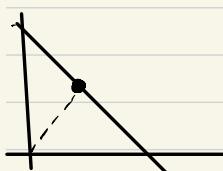
In general, a matrix $A \in \mathbb{R}^{m \times n}$

if $m=n$, then its full rank and have atleast 1 unique soln.

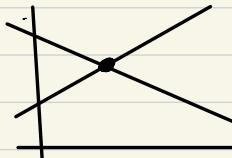
if $m < n$, then system is undetermined, so no unique soln.

if $m > n$, then system is over determined as there are more constraints than unknown. So no lines will intersect

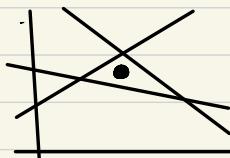
Let's assume $n=2$



$m=1$
(undetermined)



$m=2$
(unique solution)



$m=3$
(no unique solution
as over-determined)

7.7.1 Solving square systems

When $m=n$, we can solve the system of linear equations using LU decomposition i.e.,

$$Ax = b$$

$$LUx = b \quad (\because A = LU) \quad \text{therefore, solving } y \text{ should be sufficient enough.}$$

$$Ux = L^{-1}b$$

$$\text{Let } y \triangleq L^{-1}b$$

$$Ux = y \Rightarrow x = U^{-1}y$$

$y = L^{-1}b$ or $Ly = b$
evaluating $Ly_1 = b_1$ will help recursively
solve the rest.

7.7.3 Solving overconstrained systems (least square estimation)

When $m > n$, we have an overdetermined solution. So we'll typically not have an exact solution. But we may have an approximate solution.

For a given system $Ax = b$, we can approximate a solution by reducing the cost function i.e., least squares objective

$$f(x) = \frac{1}{2} \|Ax - b\|_2^2$$

→ What happens to the rest of the derivations if cross entropy is chosen?

The gradient for this cost function will be

$$g(x) = \frac{\partial}{\partial x} f(x) = A^T(Ax - b)$$

An optimum can be found by solving $g(x) = 0$

$$A^T(Ax - b) = 0$$

$$A^T A x = A^T b$$

$$\hat{x} = (A^T A)^{-1} A^T b$$

We know that,

$$(A^T A)^{-1} A^T = A^+ \text{ (left pseudo inverse)}$$

This is the ordinary least squares solution

In order to verify whether the solution is unique, we can differentiate $g(x)$, which is basically known as

"Hessian" is given as:

$$H(x) = \frac{\partial}{\partial x} g(x) = \frac{\partial^2}{\partial x^2} f(x) = A^T A$$

- If $H(x)$ is positive definite, then we'll have a unique solution.

- If A is full rank, then H will be positive definite, since for any random vector $v > 0$

$$v^T (A^T A) v = (Av)^T (Av) = \|Av\|^2 > 0$$

- Therefore in the case full rank, the least squares objective has unique global minimum.

- So we can show that the gradient plots for full rank matrices will be a \cup curve with no bumps or intermediate valleys like this

7.8 Matrix Calculus

Calculus deals with computing "rates of change" of functions as we vary their inputs.

$$f(x_1) = \hat{y}_1 \quad f(x_2) = \hat{y}_2$$

What changed to cause the new observation

7.8.1 Derivatives

Derivatives are all about change in the output with respect to change in the input.

You are making coffee. You added 1 tsp of coffee. Your coffee is level 2 strong on a scale of 1-5 (based on your own theory coffee strength).

$$\text{So here } x_1 = 1 \text{ (tsp)}$$

$$y_1 = 2 \text{ (level)}$$

Now you add one more tsp and the coffee is level 4 strong now.

$$\text{So here the new observations } x_2 = 2 \text{ (tsp)}$$

$$y_2 = 4 \text{ (level)}$$

By how much the coffee strength changed w.r.t. number of teaspoons.

$$\frac{\text{Change in coffee strength}}{\text{Change in number of tps}} = \frac{4-2}{2-1} = 2$$

This is also the slope

$$\text{slope} = \frac{y_2 - y_1}{x_2 - x_1} = \frac{\Delta y}{\Delta x}$$

Formally, we can write as

$$\frac{\Delta y}{\Delta x} = \frac{f(x + \Delta x) - f(x)}{\Delta x} \quad \text{aka "Step size" (h)}$$

We can write derivative of a function as

$$f'(x) \triangleq \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

Remember h and Δx are same
Just trying to stick to text book notations

We can consider,

$$f(x+h) \approx f(x) + f'(x) h$$

New output \approx old output + change in output * perturbation

Finite Difference Approximation

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = \lim_{h \rightarrow 0} \frac{f(x+h/2) - f(x-h/2)}{h} = \lim_{h \rightarrow 0} \frac{f(x) - f(x-h)}{h}$$

forward difference central difference backward difference

\downarrow

When we extra top of coffee to $\rightarrow t_1$ When we added 3 tops and look what happened at $t_1 \leftarrow t_2$ When we look back at where we mixed up to $\leftarrow t_1$

Smaller the step size, the better the estimate

If h is too small ≈ 0 , then numerical errors may occur.

f' notation is called Lagrange notation.

$\frac{dy}{dx}$ notation is called Leibnitz notation

7.8.2 Gradients

Before look into gradients, let's understand what are partial derivatives are.

If we have more than one variable we would like to observe rate of change w.r.t each variable

$$\text{eg: } x^2 + y^3 = z$$

Partial derivative w.r.t x would be

$$\frac{\partial}{\partial x} (x^2 + y^3) = 2x$$

Similarly,

$$\frac{\partial}{\partial y} (x^2 + y^3) = 3y^2$$

In a matrix $A \in \mathbb{R}^{m \times n}$, we have n different features. So when there is a function output, we would like to understand rate of change w.r.t. to each feature as

$$\frac{\partial f}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(x + h e_i) - f(x)}{h}$$

e_i is the i^{th} unit vector

The **gradient** of a function at a point x is the vector of its partial derivatives

$$g = \frac{\partial f}{\partial x} = \nabla f = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$$

For a specific point, the gradient can be written as

$$g(x^*) \triangleq \left. \frac{\partial f}{\partial x} \right|_{x^*}$$

7.8.3 Directional derivative

Change of gradient along a direction v is given as:

$$\underbrace{D_v f(x)}_{\text{Directional derivative}} = \nabla f(x) \cdot v$$

7.8.5 Jacobian

If there is a function $f: \mathbb{R}^n \rightarrow \mathbb{R}^d$, the first-order partial derivatives of this multivariate function can be written as

$$J_f(x) = \frac{\partial f}{\partial x^T} \triangleq \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \vdots & & \vdots \\ \frac{\partial f_d}{\partial x_1} & \frac{\partial f_d}{\partial x_2} & \dots & \frac{\partial f_d}{\partial x_n} \end{pmatrix} = \begin{pmatrix} (\nabla f_1(x))^T \\ \vdots \\ (\nabla f_d(x))^T \end{pmatrix}$$

7.8.5.1 Multiplying Jacobians and vectors

Jacobian vector product is multiplying vector on the right side

$$J_f(x)v = \begin{pmatrix} \nabla f_1(x) \\ \vdots \\ \nabla f_d(x)^T \end{pmatrix} v$$

where J_f is the Jacobian for $f: \mathbb{R}^n \rightarrow \mathbb{R}^d$

We can approximate $J_f(x)v$ to f with just two calls.

$$\hookrightarrow f(x) \in \mathbb{R}^{n \times d}$$

Vector Jacobian Product is basically left multiplying the vector i.e.,

$$u^T J_f(x) = u^T \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_d} \right)$$

Here $J \in \mathbb{R}^{n \times d}$ and $u \in \mathbb{R}^n$

$$\begin{array}{l} J_f \in \mathbb{R}^{n \times d} \\ v \in \mathbb{R}^d \end{array} \quad \text{Same dim}$$

$$J_f(x)v \in \mathbb{R}^{n \times d}$$

So what is v & how is it approximating to f ?

If $n \geq d$, VJP is efficient

If $n \leq d$, JVP is efficient.

Jacobian of a composition

If we have composition of functions, for instance

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

Jacobian of $h(x)$ would be:

$$J_h(x) = J_g(f(x)) J_f(x) \rightsquigarrow \text{Chain Rule}$$

7.8.6 Hessian

Any function which is twice differentiable, we call the second partial derivative matrix of the function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ as "Hessian"

$$H_f = \frac{\partial^2 f}{\partial x^2} = \nabla^2 f = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

Hessian is the Jacobian of gradient.

But what does Hessian reveal?

7.8.7 Gradients of commonly used functions

Refer textbook pg.265-266 for all the formulas related to mapping differential of a function into simple vector or matrix transformations. These are referred as identities.

When implementing the differentials of a function, you don't need to use explicit differential library to compute the output. One can rely on the identities.