

23MAT204

Adithya Nair

October 14, 2024

# Contents

<b>I</b>	<b>Mathematics</b>	<b>3</b>
<b>1</b>	<b>Revision Linear Algebra</b>	<b>4</b>
1.1	Ways To Calculate Whether A Matrix Is Positive Definite . . . . .	4
1.2	Large Eigenvalue Computation . . . . .	4
1.3	Spectral Decomposition . . . . .	4
1.4	Singular Value Decomposition . . . . .	5
1.5	The Geometry of SVD . . . . .	5
1.6	Polar Decomposition . . . . .	5
1.7	Principal Component Analysis . . . . .	5
1.8	Norms Of Vectors And Matrices . . . . .	5
1.9	Best Low Rank Matrix . . . . .	7
<b>2</b>	<b>Multi variable Optimization</b>	<b>8</b>
2.1	Contour Curves . . . . .	8
2.2	Multi variable Calculus . . . . .	8
2.3	Newton's Method . . . . .	8
2.3	How solution to a linear system can be found as a solution of an optimization problem. .	9
2.4	Numerical Algorithm Of Gauss-Jacobi Method . . . . .	9
2.5	Gauss-Siedel Iteration Method . . . . .	10
2.6	Unidirectional/Line search . . . . .	10
2.7	Directions Of Change . . . . .	11
2.8	Directions Of Descent/Ascent In Numerical Algorithms. . . . .	11
2.9	Steepest Descent . . . . .	11
2.10	Numerical Method . . . . .	11
2.11	Conjugate Gradient Descent . . . . .	12
<b>3</b>	<b>Probability</b>	<b>13</b>
3.1	Probability mass function . . . . .	13
3.2	Cumulative Distribution Function . . . . .	13
3.3	Probability density function - For continuous inputs . . . . .	13
3.4	Joint Probability Distribution . . . . .	15
3.5	Double Integrals . . . . .	15
3.6	Joint Probability Density Function . . . . .	15
3.7	Independent Probabilities . . . . .	15
3.8	Binomial Distribution . . . . .	15
3.9	Poisson Distribution . . . . .	15
3.10	Continuous Uniform Distribution . . . . .	16
3.11	Exponential Distribution . . . . .	16
3.12	Gaussian/Normal Distribution . . . . .	16

3.13 Covariance And Correlation Coefficient . . . . . 16

3.14 Multiple Joint Distributions . . . . . 16

Part I

Mathematics

# Chapter 1

## Revision Linear Algebra

Every matrix satisfies its own characteristic matrix.

What this means is

$$\lambda^3 - 2\lambda^2 + \lambda - 4 = 0 \quad (1.1)$$

$$A^3 - 2A^2 + A - 4I = 0 \quad (1.2)$$

$$A^2 - 2A + I - 4A^{-1} = 0 \quad (1.3)$$

$$A^{-1} = \frac{1}{4}[A^2 - 2A + I] \quad (1.4)$$

$$X = A^{-1}b = \frac{1}{4}[A^2 - 2A + I]b \quad (1.5)$$

There are some implications behind this. Here's a problem, generate a random 3x3 matrix, find the ranks of  $(A - \lambda_1 I)$  and then  $(A - \lambda_1 I)(A - \lambda_2 I)$

What you will see is that the rank reduces upon multiplying roots of the characteristic equation

### 1.1 Ways To Calculate Whether A Matrix Is Positive Definite

1. The minors are positive then negative alternating.
2. The eigenvalues are positive.

### 1.2 Large Eigenvalue Computation

Large eigenvalues are computed by using a matrix multiplication method. This method involves:

### 1.3 Spectral Decomposition

$$S = Q\Lambda Q^T = \lambda_1 q_1 q_1^T + \lambda_2 q_2 q_2^T + \cdots + \lambda_n q_n q_n^T$$

This allows us to represent a matrix by the sum of n rank 1 matrices. This is used in square symmetric matrices.

This is the basis behind Principal Component Analysis.

## 1.4 Singular Value Decomposition

For rectangular matrices.

$$A = U\Sigma V^T = \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T + \dots + \sigma_r u_r v_r^T$$

This is mainly used for getting useful properties about the matrix without needing to perform significant operations on this matrix, such as the orthogonality.

Singular values are the eigenvalues of a matrix which are **non-zero**

So a rectangular matrix  $m \times n$  would have  $AA^T$  which is  $m \times m$  and  $A^T A$  which is  $n \times n$ . The number of singular values is whether  $m$  or  $n$  is lower.

The trace of  $A^T A$  equal to the sum of all  $a_{ij}^2$ . The trace is the sum of all eigenvalues of  $A^T A$ , and For  $A_{m \times n}$  that's the sum of the eigenvalues squared.

This is known as the **Frobenius Norm**

This method can also be used to generate orthonormal bases for the 4 Fundamental Subspaces

## 1.5 The Geometry of SVD

Since we have an orthogonal matrix, diagonal and orthogonal matrix... It's a rotation step followed by a stretching step, finally ended by another rotation step.

## 1.6 Polar Decomposition

This decomposition is a special case of the Singular Value Decomposition. Where you decompose the elements into the polar form with an orthogonal matrix and a positive semi-definite matrix.

For 2D,

$$\vec{x} = r(\cos(\theta) + \sin(\theta))$$

$$\begin{aligned} A &= U\Sigma V^T \\ &= U(V^T V)\Sigma V^T \\ &= (UV^T)(V\Sigma V^T) \\ &= Q \times S \end{aligned}$$

Here S is the scaling matrix, and Q is the orthogonal matrix.

Although my assumptions might be wrong, this might be a way to generate rotation matrices for higher-dimension spaces. Might be useful in exploring the semantic spaces of LLMs.

## 1.7 Principal Component Analysis

Why do we use normalization while working with data? We use normalization to ensure that the units are eliminated while working with that given data. **This normalization is done by subtracting by the average and divide by the standard deviation.**

Matrix where the row and columns sums are equal.

## 1.8 Norms Of Vectors And Matrices

Norms are a way to measure the size of a vector/matrix. The norm of a vector which calculates the magnitude is known as the  $L^2$  norm or the Euclidean norm  $\|v\|_2$ . This number gives us the length of the vector.

In general, the  $\|u\|_p$  or the p-norm of a vector is  $[|u_1|^p + |u_2|^p + \dots + |u_n|^p]^{\frac{1}{p}}$ .

Now, the  $L_1$  norm would be the sum of the components of the vector. This is the sum of the projections to each corresponding axis.

The  $L_3$  norm would be  $[|u_1|^3 + |u_2|^3 + \dots + |u_n|^3]^{\frac{1}{3}}$

The  $L_\infty$  norm would be  $[|u_1|^\infty + |u_2|^\infty + \dots + |u_n|^\infty]^{\frac{1}{\infty}}$ . This returns the maximum vector component of the vector.

$$x + y \& = 1$$

Let's take the equation,  $\|v\|_1 = 1$   $x - y \& = 1$

$$x + y \& = 1$$

$$-x - y = 1$$

[./figures/l1eq1](#)

The S-norm of a vector  $\vec{x}$  is  $\vec{x}^T S \vec{x}$ . When S is a symmetric positive definite matrix, this S-norm is known as the energy of vector  $\vec{v}$

There are three types of Matrix norms:

1. Spectral Norm

2. Frobenius Norm =  $\sqrt{\sigma_1^2 + \sigma_2^2 + \dots + \sigma_r^2}$

3. Nuclear Norm

1. Spectral Norm We know that the vector norm for a vector  $\vec{x}$  is nothing but  $\vec{x}^T \vec{x}$ . We take this property.

$$\begin{aligned} \text{Max} \|A\|_2^2 &= \text{Max} \frac{\|Ax\|_2^2}{\|x\|_2^2} \\ &= \text{Max} \frac{x^T A^T A x}{x^T x} \\ &= \text{Max} \{\lambda_i(S)\} = \lambda_1 = \sigma_1^2 \end{aligned}$$

2. Frobenius Norm The Frobenius norm for a matrix M,  $\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$  is the equation  $\sqrt{a_{11}^2 + a_{12}^2 + a_{21}^2 + a_{22}^2}$

3. Nuclear Norm The nuclear norm is the sum of the singular values of a matrix A.

For an identity matrix,

- The Spectral Norm is 1
- The Frobenius Norm is  $\sqrt{n}$
- The Nuclear Norm is  $n$

For an orthogonal matrix,

- The Spectral Norm is 1
- The Frobenius Norm is  $\sqrt{n}$
- The Nuclear Norm is  $n$

## 1.9 Best Low Rank Matrix

We say that a matrix is the best approximation of another matrix, based on the Frobenius Norm.. For a singular value decomposition  $A = U\Sigma V^T$ , if we assume that the singular values are arranged in descending order... We can select the singular value range where the values are significant contributors to the final matrix.

We can then reduce the size of  $U, \Sigma$  and  $V^T$  into a smaller matrix  $B$ , based on the number of singular values chosen which would give the best approximation.

Let  $A = U\Sigma V^T$  where  $\Sigma : \sigma_1 \geq \sigma_2 \geq \dots \sigma_n$ , then  $B = U_{m \times m} \Sigma V^T_{n \times n}$  is a best rank- $k$  approx. to  $A$ . Where,  $S$  is a diagonal matrix of  $n \times n$  where  $s_i = \sigma_i (i = 1 \dots k)$  else  $s_i = 0$ , by best  $B$  is a solution to  $\min_B \|A - B\|_F$  where  $\text{rank}(B) = k$



## Chapter 2

# Multi variable Optimization

We minimize  $f(x_1, x_2, \dots, x_3)$ .

### 2.1 Contour Curves

### 2.2 Multi variable Calculus

The points where  $\nabla f = \vec{0}$  are called the stationary points. The Hessian matrix should be positive definite for a minima and negative definite for a maxima. A point where there is no change, is known as a saddle point.

$$f(x, y) = x^3 + y^3 + 2x^2 + 4y^2 + 6$$

$$\frac{\partial f}{\partial x} = 3x^2 + 4x$$

$$\frac{\partial^2 f}{\partial x} = 6x + 4$$

$$\frac{\partial f}{\partial y} = 3y^2 + 8y$$

$$\frac{\partial f}{\partial y \partial x} = 0$$

$$\frac{\partial^2 f}{\partial y} = 6y + 8$$

$$x = 0, \frac{-4}{3}$$

$$y = 0, \frac{-8}{3}$$

### 2.3 Newton's Method

Newton's Method is a numerical method to find a minimum of a function.

Iterative formula:

$$x_{n+1} = x_n - \frac{f'(x_n)}{f''(x_n)}$$

Find the minimum of  $f(x) = x^2 + \frac{54}{x}$  using Newton's method.

Newton's method can be practically done for all functions by evaluating the first and second derivatives numerically and then apply the formula

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \text{ to find root of } f(x) = 0$$

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \text{ to find minimum of } f(x) \text{ (root of } f'(x) = 0 \text{) to find minimum of } f(x) \text{ (root of } f'(x) = 0 \text{) to find minimum}$$

Find an objective function whose solution is specified as  $\vec{x} : A\vec{x} = b$

Take for example,  $\frac{1}{2}x^T Ax - b^T x + c$ ,  $x \in R^n$ ,  $A = A^T$

## 2.4 Numerical Algorithm Of Gauss-Jacobi Method

Input A = [a<sub>ij</sub>], X0 = x<sup>(0)</sup>, tolerance TOL, maximum number of iterations,

$$Ax = b$$

We separate A, A = D - L - U

$$(D - L - U)x = b$$

$$(Dx = b + (L + U)x$$

$$(x = D^{-1}b + D^{-1}(L + U)x$$

$$x^{(k)} = Tx^{(k-1)} + c$$

The code to implement the method:

```
A = [5,-2,3;-3,9,1;2,-1,-7];
```

```
b = [-1;2;3];
```

```
n = 100; % No. of iterations
```

```
D = diag(diag(A));
```

```
L = -tril(A,-1);
```

```
U = -triu(A,1);
```

```
X = zeros(size(A)(1),n);
```

```
T = inv(D)*(L+U);
```

```
c = inv(D)*b;
```

```
for i = 2:n
```

```
    X(:,i) = T*X(:,i-1) + c;
```

```
    if(X(:,i) == X(:,i-1))
```

```
        i-1
```

```
        break;
```

```
    end
```

```
end
```

## 2.5 Gauss-Siedel Iteration Method

The iterative method is somewhat similar to <Gauss-Jacobi method>.

The only difference is that the new values are computed by already existing new values.

```
A = [5,-2,3;-3,9,1;2,-1,-7];
b = [-1;2;3];

n = 100; % No. of iterations
D = diag(diag(A));
L = -tril(A,-1);
U = -triu(A,1);

X = zeros(size(A)(1),n);
T = inv(D-L)*U;
c = inv(D-L)*b;

for i = 2:n
    X(:,i) = T*X(:,i-1) + c;
    if(X(:,i) == X(:,i-1))
        display("The iteration converges at:")
        i-1
        break;
    end
end
end
```

$$(D - L)x^k = Ux^{k-1} + b$$
$$x^k = (D - L)^{-1}\vec{b} + (D - L)^{-1}(U\vec{x})$$

If the matrix is 'diagonally dominant', where the magnitude of the diagonal elements should be greater than the sum of the magnitude of the other elements(absolute value) in the same row.

## 2.6 Unidirectional/Line search

For a given function,

Find the minima of  $f(x) = (x - 1)^2 + (y - 2)^2$  starting from  $(0, 0)$  along the direction of x-axis.

How to perform uni-directional search:

- Write the parametric representation of the line search,

$$\vec{s}(t) = \vec{a} + t\vec{b}$$

- Write the function in terms of  $t$ ,  $f(\vec{s}(t))$ . this is a single variable function.
- Find the minimum of  $f(\vec{s}(t))$ ,  $t^*$
- Obtain the solution for unidirectional search by substituting  $t^*$  in  $\vec{s}(t)$

Find the minimum of the function  $f(x) = (x - 2)^2 - y$ , starting from the point  $(-1, 0)$  along  $(1, 1)^T$

## 2.7 Directions Of Change

A direction given by the vector  $\{x^{(k)}\}$  is a descent direction only if the function value decreases along that direction from the point,  $x^{(k)}$ , when

$$\nabla f(x^{(k)}) \cdot d < 0$$

A direction given by the vector  $\{x^{(k)}\}$  is an ascent direction only if the function value increases along that direction from the point,  $x^{(k)}$ , when

$$\nabla f(x^{(k)}) \cdot d > 0$$

Question: Check whether the given function have a descent direction  $x = (2, -1)$  along the given directions  $d_1$  and  $d_2$

$$f = 2x_1^2 + x_2^2 - 2x_1x_2 + 2x_1^3 + x_2^4$$

## 2.8 Directions Of Descent/Ascent In Numerical Algorithms.

In any numerical algorithm to find the optimum of an unconstrained problem, the iterative formula used is:

$$x^{k+1} = x^k + \alpha^k d^k$$

Where  $\alpha^k$  is the step length in step k and  $d^k$  is the direction of descent in step k, if it's a minimization problem or a direction of ascent if it's a maximization problem.

Newton's Method:  $x_{k+1} = x_k - (H(x_k)^{-1} \nabla f(x_k))$  Gradient Descent -  $x_{k+1} = x_k + \alpha_k d_k$  where  $d_k = -\nabla f(x_k)$  and  $\alpha$  = using line search

A direction along which the function increases rapidly from a point is given by the gradient of the function at that point

## 2.9 Steepest Descent

A direction along which the function increases rapidly from a point is given by the gradient of the function at that point

A direction along which the function decreases rapidly from a point is given by the negative of the gradient of the function at that point

The iterative formula of method of steepest descent is:

$$x^{k+1} = x^k + \alpha^k d^k$$

The descent direction  $d^k$  is along the steepest descent direction,  $d^k = -\nabla f(x^k)$ . The step length  $\alpha^k$  is obtained by performing a unidirectional search from  $x^k$  along the direction  $d^k$ , by minimizing,

This method produces successive directions that are perpendicular to each other.

Near the minima, during line search the convergence is very slow.

## 2.10 Numerical Method

To implement this programmatically, we can write... For  $Ax = b$ , we can convert this to an equivalent optimization problem.

$$f(x) = \frac{1}{2} x^T A x + -b^T x + c$$

We have a linear system of equations, which we can be solved by  $Ax = b$

We can write  $g_0 = Ax_0 - b$

## 2.11 Conjugate Gradient Descent

The problem with normal gradient descent is that there is no way to predict when the function converges. There's a method that can reliably find the convergence in a predictable fashion. That's the conjugate gradient method.

First we must understand the conjugate direction and the terms associated with them

1. Krylov Subspace Let  $A \in R^{n \times n}, b \in R^n$ , the Krylov subspace  $K_j(A, b)$  is defined as  $K_j(A, b) = \text{span}\{\vec{b}, A\vec{b}, A^2\vec{b}, A^3\vec{b}, \dots, A^{j-1}\vec{b}\}$ . Thus  $K_j(A, b)$  is subspace of  $\mathbb{R}^n$
2. Krylov Matrix The Krylov Matrix is just the matrix consisting of the basis vectors as the column vectors. The Krylov subspace is the column space of such a matrix.
3. Motivation For Krylov Subspaces We know the Cayley Hamilton theorem, and its use in calculating the inverse from the characteristic equation Interestingly, the calculation for the inverse leads to a linear combination of the basis vectors of the Krylov subspace when solving linear systems.
4. Conjugate Directions For A real symmetric matrix  $n \times n$  with rank  $n$   
The directions  $d_0 d_1 \dots d_{n-1}$  are said to be A-conjugate if,

$$d_i A d_j = 0$$

for  $i \neq j$

It's a new definition for orthogonality(not orthonormal)

5. Numerical Method Input  $A, b, x^0$

Compute  $\vec{r}_0 = b - Ax^0$

```
A = [1 2 3; 2 3 4; 3 4 5];
```

```
b = [6 9 12]';
```

```
x = randi([-9,9], length(b), 1);
```

```
r = b - A*x;
```

```
d = r;
```

```
alpha = (r'*r)/(d'*(A*d));
```

```
beta = 0;
```

```
for i = 1:size(A,1)
```

```
    x = x + alpha*d;
```

```
    rnew = r - alpha*(A*d);
```

```
    beta = (rnew'*rnew)/(r'*r);
```

```
    d = rnew + beta*d;
```

```
    r = rnew;
```

```
    alpha = (r'*r)/(d'*(A*d));
```

```
    if sqrt(rnew) < 1e-10
```

```
        break;
```

```
    end
```

```
end
```

```
x
```

```
r
```

# Chapter 3

## Probability

### 3.1 Probability mass function

For discrete inputs The probability mass function for a discrete random variable  $x$  is defined as,  $f(x) = P(X = x)$

1. Properties

- $0 \leq f(x) \leq 1$
- $\sum f(x) = 1$
- $p(\text{at least } n \text{ outcomes}) = p(x \geq n) = f(n) + \dots + f(\text{end})$

2. Mean Of Random Variable

$$U_x = E(x) = \sum x f(x)$$

3. Variance Of Random Variable

$$\sigma_x^2 = E((x - \mu_x^2)) = \sum (x - \mu)^2 f(x) = \sum x^2 f(x) - \mu^2$$

$$\sigma_x^2 = E(X^2) - \mu^2$$

### 3.2 Cumulative Distribution Function

A cumulative distribution function is just adding the previous terms to the new terms.

### 3.3 Probability density function - For continuous inputs

For continuous inputs

The probability density function for a continuous random variable  $x$  is defined as,  $P(a \leq x \leq b) = \int_a^b f(x) dx$

1. Binomial Distribution The Bernoulli Process

- (a) Repeated trials,
- (b) Each trial results in an outcome that may be classified as a success
- (c) The trials are independent.

Where  $n$  is the number of trials,  $x$  is the random variable distribution,  $q = 1 - p$   ${}^nC_x p^x q^{n-x}$

We have,

${}^nC_x$  Different ways  $x$  successes can occur among  $n$  trials

$p^x$  – probability of getting  $x$  successes

$(1 - p)^{n-x}$  – probability of getting  $n - x$  failures

## 2. Cumulative Distribution Function

$$\int_0^x f(y)dy$$

## 3. Mean Of Random Variable

$$U_x = E(x) = \sum x f(x)$$

## 4. Properties

- $f(x) \geq 0$
- $\int_x f(x) = 1$ , where  $\int_x$  is the range for which  $x$  is defined.
- For a continuous random variable, the probability at a point is 0, the probability that  $x$  takes a discrete value is zero.

## 5. Uniform Distribution A probability distribution which is constant through the interval, given by the function,

$$f(x) = \frac{1}{b - a}$$

$$Mean = \frac{a + b}{2}$$

$$Variance = \frac{(b - a)^2}{12}$$

$$F(x) = \int_{-\infty}^x f(x)dx = \int_{-\infty}^x \frac{1}{b - a}$$

$$f(X) = \begin{cases} 0, & -\infty < x < 0 \\ x, & 0 < x < 1 \\ 1, & x > 1 \end{cases}$$

## 6. Variance Of Random Variable

$$\sigma_x^2 = E((x - \mu_x^2)) = \sum (x - \mu)^2 f(x) = \sum x^2 f(x) - \mu^2$$

$$\sigma_x^2 = E(X^2) - \mu^2$$

### 3.4 Joint Probability Distribution

This is given for two variables, a function  $f(x, y)$  gives us a probability,  $P(X = x_1, X = x_2)$

1. Properties

- $0 \leq f(x_1, y_1) \leq 1$
- $\sum_x \sum_y f(x, y) = 1$

2. Marginal Distribution Function Marginal distribution functions are the probability distributions of the variables themselves.

- The distribution of  $x_1$  is given by the columns, and  $x_2$  is given by the rows. In other words you're fixing the values and adding the remaining.

3. Mean Of Marginal Values

$$Mean(X_1) = E(X_1) = \sum_x x_1 f_{1(x_1)}$$

### 3.5 Double Integrals

$$\int_x \int_y f(x, y) dy dx$$

1. Area Of A Region

$$R = \int_R 1 dx dy$$

### 3.6 Joint Probability Density Function

For a continuous scalar field,  $P(a \leq x \leq b, c \leq y \leq d)$

1. Properties

- $\int \int f(x, y) dx dy = 1$

### 3.7 Independent Probabilities

Two events  $A$  and  $B$  are said to be independent if  $P(A|B) = P(A)$  or  $P(B|A) = P(B)$

Two discrete random variables  $x_1$  and  $x_2$  are said to be independent if  $f(x_1, x_2) = f(x_1)f_2(x_2) \forall x_1 \& x_2$

### 3.8 Binomial Distribution

The Bernoulli process has these properties:

1. The experiment consists of repeated trials
2. Each trial results in an outcome that may be classified as a success or failure.

### 3.9 Poisson Distribution

$$\lim_{n \rightarrow \infty} P(X = x) = \frac{e^{-\lambda} \lambda^x}{x!}$$



### 3.10 Continuous Uniform Distribution

$$f(x) = \frac{1}{b-a}$$

### 3.11 Exponential Distribution

Let  $N$  be a random variable with a Poisson Distribution with mean  $\lambda x$  and  $X$  be the random variable that gives the time between each arrival or count. Then,

$$P(X > x) = P(N = 0) = \frac{e^{-\lambda x}(\lambda x)^0}{0!} = e^{-\lambda x}$$

PDF,

$$\lambda e^{-\lambda x}$$

### 3.12 Gaussian/Normal Distribution

- It is one of the most common probability distribution

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

### 3.13 Covariance And Correlation Coefficient

Correlation is a measure of how much two data points depend on each other.

$$\rho = \frac{Cov(X_1, X_2)}{\sigma_1 \sigma_2}$$

### 3.14 Multiple Joint Distributions

- Multiple discrete random variables

Multinomial distributions,

$$P(X_1 = x_1, X_2 = x_2, \dots, X_k = x_k) = \frac{n!}{x_1! x_2! x_3! \dots x_k!} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k}$$