

# Linear Algebra Tutorial

## CS5011 - Machine Learning

Abhinav Garlapati    Varun Gangal

Department of Computer Science  
IIT Madras

January 23, 2016

# What is Linear Algebra

## Linear Algebra

Linear algebra is the branch of mathematics concerning vector spaces and linear mappings between such spaces. It includes the study of lines, planes, and subspaces, but is also concerned with properties common to all vector spaces.

Why do we study Linear Algebra?

- Provides a way to compactly represent & operate on sets of linear equations.
- In machine learning, we represent data as matrices and hence it is natural to use notions and formalisms developed in Linear Algebra.

- Consider the following system of equations:

$$\begin{aligned}4x_1 - 5x_2 &= -13 \\ -2x_1 + 3x_2 &= 9\end{aligned}$$

- In matrix notation, the system is more compactly represented as:

$$Ax = b$$

$$A = \begin{bmatrix} 4 & -5 \\ -2 & 3 \end{bmatrix}$$

$$b = \begin{bmatrix} -13 \\ 9 \end{bmatrix}$$

## Definition

A set  $V$  with two operations  $+$  and  $\cdot$  is said to be a **vector space** if it is closed under both these operations and satisfies the following eight axioms.

- 1 Commutative Law

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

- 2 Associative Law

$$(x + y) + z = x + (y + z), \quad \forall x, y, z \in V$$

- 3 Additive identity

$$\exists 0 \in V \quad \text{s.t.} \quad x + 0 = x, \quad \forall x \in V$$

- 4 Additive inverse

$$\forall x \in V, \quad \exists \tilde{x} \quad \text{s.t.} \quad x + \tilde{x} = 0$$

# Vector Space (Contd..)

## 5 Distributive Law

$$\alpha \cdot (x + y) = \alpha \cdot x + \alpha \cdot y, \quad \forall \alpha \in \mathbb{R}, x, y \in V$$

## 6 Distributive Law

$$(\alpha + \beta) \cdot x = \alpha \cdot x + \beta \cdot x, \quad \forall \alpha, \beta \in \mathbb{R}, x \in V$$

## 7 Associative Law

$$(\alpha\beta) \cdot x = \alpha \cdot (\beta \cdot x), \quad \forall \alpha, \beta \in \mathbb{R}, x \in V$$

## 8 Unitary Law

$$1 \cdot x = x, \quad \forall x \in V$$

## Definition

Let  $W$  be a subset of a vector space  $V$ . Then  $W$  is called a **subspace** of  $V$  if  $W$  is a vector space.

- Do we have to verify all 8 conditions to check whether a given subset of a vector space is a subspace?
- **Theorem:** Let  $W$  be a subset of a vector space  $V$ . Then  $W$  is a subspace of  $V$  if and only if  $W$  is non-empty and  $x + \alpha y \in W, \quad \forall x, y \in W, \alpha \in \mathbb{R}$

## Definition

Norm is any function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  satisfying:

- 1  $\forall x \in \mathbb{R}^n, \quad f(x) \geq 0$  (non-negativity)
- 2  $f(x) = 0$  iff  $x = 0$  (definiteness)
- 3  $\forall x \in \mathbb{R}^n, \quad f(tx) = |t|f(x)$  (homogeneity)
- 4  $\forall x, y \in \mathbb{R}^n, \quad f(x + y) \leq f(x) + f(y)$  (triangle inequality)

- Example -  $l_p$  norm

$$\|x\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}$$

- Matrices can have norms too - e.g., Frobenius norm

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n A_{ij}^2} = \sqrt{\text{tr}(A^T A)} \quad (1)$$

# Range Of A Matrix

- The **span** of a set of vectors  $X = \{x_1, x_2, \dots, x_n\}$  is the set of all vectors that can be expressed as a linear combination of the vectors in  $X$ .

In other words, set of all vectors  $v$  such that  $v = \sum_{i=1}^{|X|} \alpha_i x_i, \alpha_i \in R$

- The **range** or **columnspace** of a matrix  $A$ , denoted by  $R(A)$  is the span of its columns. In other words, it contains all linear combinations of the columns of  $A$ . For instance, the columnspace of

$A = \begin{bmatrix} 1 & 0 \\ 5 & 4 \\ 2 & 4 \end{bmatrix}$  is the plane spanned by the vectors  $\begin{bmatrix} 1 \\ 5 \\ 2 \end{bmatrix}$  and  $\begin{bmatrix} 0 \\ 4 \\ 4 \end{bmatrix}$



# Nullspace Of A Matrix

## Definition

The nullspace  $N(A)$  of a matrix  $A \in \mathbb{R}^{m \times n}$  is the set of all vectors that equal 0 when multiplied by  $A$ . The dimensionality of the nullspace is also referred to as the **nullity** of  $A$ .

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

- Note that vectors in  $N(A)$  are of dimension  $n$ , while those in  $R(A)$  are of size  $m$ , so vectors in  $R(A^T)$  and  $N(A)$  are both of dimension  $n$ .

# Example

Consider the matrix

$$A = \begin{bmatrix} 1 & 0 \\ 5 & 4 \\ 2 & 4 \end{bmatrix}$$

The nullspace of  $A$  is made up of vectors  $x$  of the form  $\begin{bmatrix} u \\ v \end{bmatrix}$ , such that

$$\begin{bmatrix} 1 & 0 \\ 5 & 4 \\ 2 & 4 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

The nullspace here only contains the vector  $(0,0)$ .

## Another Example

Now, consider the matrix

$$B = \begin{bmatrix} 1 & 0 & 1 \\ 5 & 4 & 9 \\ 2 & 4 & 6 \end{bmatrix}$$

Here, the third column is a linear combination of the first two columns.  
Here, the nullspace is the line of all points  $x = c, y = c, z = -c$ .

## Definition

A set of vectors  $\{x_1, x_2, \dots, x_n\} \in \mathbb{R}^n$  is said to be **(linearly) independent** if no vector can be represented as a linear combination of the remaining vectors.

- i.e., if  $x_n = \sum_{i=1}^{n-1} \alpha_i x_i$  for some scalar values  $\alpha_1, \dots, \alpha_{n-1} \in \mathbb{R}$ , then we say that the vectors  $\{x_1, x_2, \dots, x_n\}$  are linearly dependent; otherwise, the vectors are linearly independent
- The **column rank** of a matrix  $A \in \mathbb{R}^{m \times n}$  is the size of the largest subset of columns of  $A$  that constitute a linearly independent set
- Similarly, **row rank** of a matrix is the largest number of rows of  $A$  that constitute a linearly independent set

# Properties Of Ranks

- For any matrix  $A \in \mathbb{R}^{m \times n}$ , it turns out that the column rank of  $A$  is equal to the row rank of  $A$ , collectively as the rank of  $A$ , denoted as **rank(A)**
- Some basic properties of the rank:
  - 1 For  $A \in \mathbb{R}^{m \times n}$ ,  $\text{rank}(A) \leq \min(m, n)$ .  
If  $\text{rank}(A) = \min(m, n)$ ,  $A$  is said to be **full rank**
  - 2 For  $A \in \mathbb{R}^{m \times n}$ ,  $\text{rank}(A) = \text{rank}(A^T)$
  - 3 For  $A \in \mathbb{R}^{m \times n}$ ,  $B \in \mathbb{R}^{n \times p}$ ,  $\text{rank}(AB) \leq \min(\text{rank}(A), \text{rank}(B))$
  - 4 For  $A, B \in \mathbb{R}^{m \times n}$ ,  $\text{rank}(A + B) \leq \text{rank}(A) + \text{rank}(B)$

# Orthogonal Matrices

- A square matrix  $U \in R^{n \times n}$  is **orthogonal** iff
  - All columns are mutually orthogonal  $v_i^T v_j = 0, \forall i \neq j$
  - All columns are normalized  $v_i^T v_i = 1, \forall i$
- If  $U$  is orthogonal,  $UU^T = U^T U = I$ . This also implies that the inverse of  $U$  happens to be its transpose.
- Another salient property of orthogonal matrices is that **they do not change** the Euclidean norm of a vector when they operate on it, i.e  $\|Ux\|_2 = \|x\|_2$ .  
Multiplication by an orthogonal matrix can be thought of as a pure rotation, i.e., it does not change the magnitude of the vector, but changes the direction.

# Quadratic Form of Matrices

- Given a square matrix  $A \in \mathbb{R}^{n \times n}$  and a vector  $x \in \mathbb{R}^n$ , the scalar value  $x^T A x$  is called a **quadratic form**
- A symmetric matrix  $A \in \mathbb{S}^n$  is positive definite (PD) if for all non-zero vectors  $x \in \mathbb{R}^n$ ,  $x^T A x > 0$
- Similarly, positive semidefinite if  $x^T A x \geq 0$ , negative definite if  $x^T A x < 0$  and negative semidefinite if  $x^T A x \leq 0$
- One important property of positive definite and negative definite matrices is that they are always full rank, and hence, invertible.
- **Gram matrix:** Given any matrix  $A \in \mathbb{R}^{m \times n}$ , matrix  $G = A^T A$  is always positive semidefinite.  
Further if  $m \geq n$ , then  $G$  is positive definite.

# Eigenvalues & Eigenvectors

- Given a square matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\lambda$  is said to be an eigenvalue of  $\mathbf{A}$  and vector  $\vec{x}$  the corresponding eigenvector if

$$A\vec{x} = \lambda\vec{x}$$

- Geometrical interpretation**

We can think of the eigenvectors of a matrix  $A$  as those vectors which upon being operated by  $A$  are only scaled but not rotated.

- Example**

$$A = \begin{bmatrix} 6 & 5 \\ 1 & 2 \end{bmatrix}, \vec{x} = \begin{bmatrix} 5 \\ 1 \end{bmatrix}$$

$$A\vec{x} = \begin{bmatrix} 35 \\ 7 \end{bmatrix} = 7\vec{x}$$



# Characteristic Equation

- Trivially, the  $\vec{0}$  vector would always be an eigenvector of any matrix. Hence, we only refer only to non-zero vectors as eigenvectors.
- Given a matrix  $A$ , how do we find all eigenvalue-eigenvector pairs?

$$A\vec{x} = \lambda\vec{x}$$

$$A\vec{x} - \lambda I\vec{x} = 0$$

$$(A - \lambda I)\vec{x} = 0$$

The above will hold iff

$$|(A - \lambda I)| = 0$$

This equation is also referred to as the characteristic equation of  $A$ . Solving the equation gives us all the eigenvalues  $\lambda$  of  $A$ . Note that these eigenvalues can be **complex**.

- ① The trace  $\text{tr}(A)$  of a matrix  $A$  also equals the sum of its  $n$  eigenvalues.

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

- ② The determinant  $|A|$  is equal to the product of the eigenvalues.

$$|A| = \prod_{i=1}^n \lambda_i$$

- ③ The rank of a matrix is equal to the number of non zero eigenvalues of  $A$ .
- ④ If  $A$  is invertible, then the eigenvalues of  $A^{-1}$  are of form  $\frac{1}{\lambda_i}$ , where  $\lambda_i$  are the eigenvalues of  $A$ .

- **Theorem**

*If a real matrix  $A^{n \times n}$  has all eigenvalues distinct, then all its eigenvectors are linearly independent*

- **Proof**

We will do a proof by means of contradiction. Suppose a matrix  $A$  has  $n$  distinct eigenvalues, but a set of  $k$  eigenvectors is linearly dependent, and  $k$  is chosen so that it is the smallest such set that is linearly dependent.

$$\begin{aligned}\sum_{i=1}^{i=k} a_i \vec{v}_i &= \vec{0} \\ (A - \lambda_k I) \sum_{i=1}^{i=k} a_i \vec{v}_i &= \vec{0} \\ \sum_{i=1}^{i=k} (A - \lambda_k I) a_i \vec{v}_i &= \vec{0} \\ \sum_{i=1}^{i=k} a_i (\lambda_i - \lambda_k) \vec{v}_i &= \vec{0}\end{aligned}$$

Since the eigenvalues are distinct,  $\lambda_i \neq \lambda_k \forall i \neq k$ . Thus the set of  $(k - 1)$  eigenvectors is also linearly dependent, violating our assumption of it being the smallest such set. This is a result of our incorrect starting assumption. Hence proved by contradiction.

# Diagonalization

Given a matrix  $A$ , we consider the matrix  $S$  with each column being an eigenvector of  $A$

$$S = \begin{bmatrix} \vdots & \vdots & \dots & \vdots \\ \vec{v}_1 & \vec{v}_2 & \dots & \vec{v}_n \\ \vdots & \vdots & \dots & \vdots \end{bmatrix}$$

$$AS = \begin{bmatrix} \vdots & \vdots & \dots & \vdots \\ \lambda_1 \vec{v}_1 & \lambda_2 \vec{v}_2 & \dots & \lambda_n \vec{v}_n \\ \vdots & \vdots & \dots & \vdots \end{bmatrix}$$

$$AS = \begin{bmatrix} \vdots & \vdots & \dots & \vdots \\ \vec{v}_1 & \vec{v}_2 & \dots & \vec{v}_n \\ \vdots & \vdots & \dots & \vdots \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \dots \\ \vdots & \ddots & \dots \\ 0 & \dots & \lambda_n \end{bmatrix}$$

# Diagonalization

$$AS = S\Lambda$$

$$A = S\Lambda S^{-1}$$

- $S^{-1}AS$  is diagonal
- Note that the above result is dependent on  $S$  being invertible. In the case where the eigenvalues are distinct, this will be true since the eigenvectors will be linearly independent

# Properties of Diagonalization

- 1 A square matrix  $A$  is said to be **diagonalizable** if  $\exists S$  such that  $A = S\Lambda S^{-1}$ .
- 2 Diagonalization can be used to simplify computation of the higher powers of a matrix  $A$ , if the diagonalized form is available

$$A^n = (S\Lambda S^{-1})(S\Lambda S^{-1}) \dots (S\Lambda S^{-1})$$

$$A^n = S\Lambda^n S^{-1}$$

$\Lambda^n$  is simple to compute since it is a diagonal matrix.

# Eigenvalues & Eigenvectors of Symmetric Matrices

- Two important properties for a symmetric matrix  $A$ :
  - 1 All the eigenvalues of  $A$  are real
  - 2 The eigenvectors of  $A$  are orthonormal, i.e., matrix  $S$  is orthogonal.  
Thus,  $A = S\Lambda S^T$ .
- Definiteness of a symmetric matrix depends entirely on the sign of its eigenvalues. Suppose  $A = S\Lambda S^T$ , then

$$x^T A x = x^T S \Lambda S^T x = y^T \Lambda y = \sum_{i=1}^n \lambda_i y_i^2$$

- Since  $y_i^2 \geq 0$ , sign of expression depends entirely on the  $\lambda_i$ 's. For example, if all  $\lambda_i > 0$ , then matrix  $A$  is positive definite.



# Eigenvalues of a PSD Matrix

Consider a positive semi definite matrix  $A$ . Then,  $\forall \vec{x}$  which are eigenvectors of  $A$ .

$$\vec{x}^T A \vec{x} \geq 0$$

$$\lambda \vec{x}^T \vec{x} \geq 0$$

$$\lambda \|\vec{x}\|^2 \geq 0$$

Hence, all eigenvalues of a PSD matrix are non-negative.

# Singular Value Decomposition

- 1 We saw that diagonalization is applicable only to square matrices. We need some analogue for rectangular matrices too, since we often encounter them, e.g the Document-Term matrix. For a rectangular matrix, we consider left singular and right singular vectors as two bases instead of a single base of eigenvectors for square matrices.
- 2 The Singular Value Decomposition is given by  $A = U\Sigma V^T$  where  $U \in R^{m \times m}$ ,  $\Sigma \in R^{m \times n}$  and  $V \in R^{n \times n}$ .

# Singular Value Decomposition

- 1  $U$  is such that the  $m$  columns of  $U$  are the eigenvectors of  $AA^T$ , also known as the left singular vectors of  $A$ .
- 2  $V$  is such that the  $n$  columns of  $V$  are the eigenvectors of  $A^T A$ , also known as the right singular vectors of  $A$ .
- 3  $\Sigma$  is a rectangular diagonal matrix with each element being the square root of an eigenvalue of  $AA^T$  or  $A^T A$

Significance: SVD allows us to construct a lower rank approximation of a rectangular matrix. We choose only the top  $r$  singular values in  $\Sigma$ , and the corresponding columns in  $U$  and rows in  $V^T$

## 1 The Gradient

Consider a function  $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ . The gradient  $\nabla_A f(A)$  denotes the matrix of partial derivatives with respect to every element of the matrix  $A$ . Each element is given by  $(\nabla_A f(A))_{ij} = \frac{\partial f(A)}{\partial A_{ij}}$

## 2 The Hessian

Suppose a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  takes in vectors and returns real numbers. The Hessian, denoted as  $\nabla_x^2 f(x)$  or  $H$  is the  $n \times n$  matrix of partial derivatives.  $(\nabla_x^2 f(x))_{ij} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j}$ . Note that the Hessian is always symmetric.

- 3 Note that the Hessian is not the gradient of the gradient, since the gradient is a vector, and we cannot take the gradient of the vector. However, if we do take elementwise gradients of every element of the gradient, then we can construct the Hessian.

# Differentiating Linear and Quadratic Functions

If  $f(x) = b^T x$ , for some constant  $b \in \mathbb{R}^n$ . Let us find the gradient of  $f$ .

$$f(x) = \sum_{i=1}^{i=n} b_i x_i$$
$$\frac{\partial f(x)}{\partial x_k} = b_k$$

We can see that  $\frac{\partial b^T x}{\partial x} = b$ . We can intuitively see how this relates to differentiating  $f(x) = ax$  with respect to  $x$  when  $a$  and  $x$  are real scalars.

# Differentiating Linear and Quadratic Functions

Consider the function  $f(x) = x^T Ax$  where  $x \in \mathbb{R}^n$  and  $A \in \mathbb{R}^{n \times n}$  is a known symmetric matrix.

$$f(x) = \sum_{i=1}^{i=n} \sum_{j=1}^{j=n} A_{ij} x_i x_j$$

$$\frac{\partial f(x)}{\partial x_k} = \frac{\partial}{\partial x_k} \left[ \sum_{i \neq k} \sum_{j \neq k} A_{ij} x_i x_j + \sum_{i \neq k} A_{ik} x_i x_k + \sum_{j \neq k} A_{kj} x_k x_j + A_{kk} x_k^2 \right]$$

$$\frac{\partial f(x)}{\partial x_k} = \sum_{i \neq k} A_{ik} x_i + \sum_{j \neq k} A_{kj} x_j + 2A_{kk} x_k$$

$$\frac{\partial f(x)}{\partial x_k} = \sum_{i=1}^n A_{ik} x_i + \sum_{j=1}^n A_{kj} x_j$$

$$\frac{\partial f(x)}{\partial x_k} = 2 \sum_{i=1}^n A_{ki} x_i$$

# Differentiating Linear and Quadratic Functions

Thus  $\nabla_x(x^T Ax) = 2Ax$ . Now, let us find the Hessian  $H$ .

$$\frac{\partial}{\partial x_k} \frac{\partial f(x)}{\partial x_l} = \frac{\partial}{\partial x_k} \left( 2 \sum_{i=1}^{i=n} A_{li} x_i \right) = 2A_{kl}$$

Hence,  $\nabla_x^2(x^T Ax) = 2A$ .