

MSE 546: Advanced Machine Learning

Sirisha Rambhatla

University of Waterloo

Lecture

Math Background: Linear Algebra and Optimization

Outline

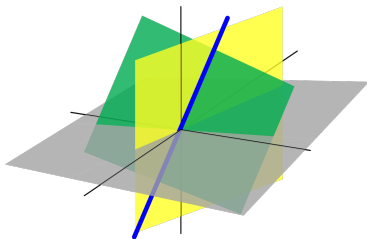
- 1 Math Background: Linear Algebra and Matrix Calculus
- 2 Math Background: Optimization for ML

Outline

- 1 Math Background: Linear Algebra and Matrix Calculus
- 2 Math Background: Optimization for ML

Ease of notation and much more!

The properties of data samples (vectors) in higher dimensions lead to key insights, e.g. dimensionality reduction. It also makes our lives easier in analysis.



For $w, x \in \mathbb{R}^n$ inner-product is represented as

$$\sum_{i=1}^n w_i x_i$$

Scalar Representation

$$w^\top x$$

Vector Representation

Linear Algebra

Linear algebra is the study of matrices and vectors, and other generalizations such as tensors.

- vector
- matrices
- tensors

Definition and Notation: Vectors

Vector: a vector $\mathbf{x} \in \mathbb{R}^n$ is a list of n numbers, usually written as a column vector:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

- a vector of all ones is denoted by $\mathbf{1}$
- a vector of all zeros is denoted by $\mathbf{0}$
- a unit vector \mathbf{e}_i is a vector of all 0's, except 1 at the i -th entry aka **one-hot** vector

Definition and Notation: Matrices

Matrix: A matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ with m rows and n columns is a 2d array of numbers, arranged as follows:

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

- if $m = n$ the matrix is said to be **square**
- \mathbf{A}_{ij}/A_{ij} or $\mathbf{A}_{i,j}/A_{i,j}$ to denotes the entry in i -th row and j -th column
- $\mathbf{A}_{i,:}$ to denote the i -th row and $\mathbf{A}_{:,j}$ to denote the j -th column.
- convention: all vectors as column vectors by default
- we use bold upper case letters \mathbf{A} to denote matrices, bold lower case letters \mathbf{a} to denote vectors, and non-bold letters a to denote scalars.

Definition and Notation

Matrices continued ...

We can view a matrix as a set of columns stacked along the horizontal axis:

$$\mathbf{A} = \left[\begin{array}{c|c|c|c} | & | & & | \\ \mathbf{A}_{:,1} & \mathbf{A}_{:,2} & \cdots & \mathbf{A}_{:,n} \\ | & | & & | \end{array} \right] \quad \mathbf{A} = [\mathbf{A}_{:,1}, \mathbf{A}_{:,2}, \dots, \mathbf{A}_{:,n}]$$

also view a matrix as a set of rows stacked along the vertical axis

$$\mathbf{A} = \left[\begin{array}{c|c|c} \text{---} & \mathbf{A}_{1,:}^T & \text{---} \\ \text{---} & \mathbf{A}_{2,:}^T & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{A}_{m,:}^T & \text{---} \end{array} \right] \quad \mathbf{A} = [\mathbf{A}_{1,:}; \mathbf{A}_{2,:}; \dots; \mathbf{A}_{m,:}]$$

Definition and Notation

Transpose of a Matrix

The transpose of a matrix results from “flipping” the rows and columns. Given a matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$, its transpose, written $\mathbf{A}^\top \in \mathbb{R}^{m \times n}$, is defined as $(\mathbf{A}^\top)_{ij} = (\mathbf{A})_{ji}$

$$(\mathbf{A}^\top)^\top = \mathbf{A}$$

$$(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top$$

$$(\mathbf{A} + \mathbf{B})^\top = \mathbf{A}^\top + \mathbf{B}^\top$$

- If for a square matrix $\mathbf{A} = \mathbf{A}^\top$, it is called **symmetric**

Vector Spaces

Vector Spaces A vector space is a collection of such vectors, which can be added together, and scaled by scalars (1-dimensional numbers), in order to create new points. These operations are defined to operate element-wise.

Vector Addition and Scaling For a vector $\mathbf{x} \in \mathbb{R}^n$

- **Addition:** $\mathbf{x} + \mathbf{y} = (x_1 + y_1, \dots, x_n + y_n)$
- **Scaling:** $c\mathbf{x} = (cx_1, \dots, cx_n)$, where $c \in \mathbb{R}$

Linear independence, spans and basis sets

For a set of vectors $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$:

Linear Independence: If none of these vectors can be represented as a linear combination of the remaining vectors.

$$\mathbf{x}_j \text{ cannot be represented as } \sum_{i=1, i \neq j}^{n-1} \alpha_i \mathbf{x}_i$$

Span: Span of a set of vectors $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ is defined as the set of all vectors that can be expressed as a linear combination of $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$

$$\text{span}(\{\mathbf{x}_1, \dots, \mathbf{x}_n\}) \triangleq \left\{ \mathbf{v} : \mathbf{v} = \sum_{i=1}^n \alpha_i \mathbf{x}_i, \alpha_i \in \mathbb{R} \right\}$$

Basis \mathcal{B} is a set of linearly independent vectors that span the whole space, i.e., $\text{span}(\mathcal{B}) = \mathbb{R}^n$, e.g. standard basis uses the coordinate vectors $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n$.

Range and Null Space of a matrix

Range Space aka Column Space Viewing $\mathbf{A} \in \mathbb{R}^{m \times n}$ as a set of n vectors in \mathbb{R}^m . The range of this matrix is the span of the columns of \mathbf{A} , i.e.

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

Null Space The set of all vectors $\mathbf{x} \in \mathbb{R}^n$ that get mapped to the null vector when multiplied by \mathbf{A} , i.e.,

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

The span of the rows of \mathbf{A} is the complement to the nullspace of \mathbf{A}

Vector and Matrix Norms

Vector Norms A norm of a vector $\|x\|$ is, informally, a measure of the “length” of the vector. More formally, a **norm** is any function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ that satisfies four properties:

- For all $x \in \mathbb{R}^n$, $f(x) \geq 0$ (**non-negativity**).
- $f(x) = 0$ if and only if $x = 0$ (**definiteness**)
- For all $x \in \mathbb{R}^n$, $t \in \mathbb{R}$, $f(tx) = |t|f(x)$ (**homogeneity**)
- For all $x, y \in \mathbb{R}^n$, $f(x + y) \leq f(x) + f(y)$ (**triangle inequality**)

Common Examples

- **p-norm** $\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$
- **2-norm aka Euclidean norm** $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$ also $\|x\|_2^2 = x^\top x$
- **1-norm** $\|x\|_1 = \sum_{i=1}^n |x_i|$
- **Max-norm** $\|x\|_\infty = \max_i |x_i|$

Vector and Matrix Norms

Matrix Norms

- **Frobenius Norm** If we think of a matrix as a vector, we can define the matrix norm in terms of a vector norm, $\|\mathbf{A}\| = \|\text{vec}(\mathbf{A})\|$. If the vector norm is the 2-norm, the corresponding matrix norm is the Frobenius norm:

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

Here $\text{tr}(\cdot)$ denotes the trace of the matrix.

Some Key Properties of Matrices

Rank of a Matrix The **column rank** of a matrix \mathbf{A} is the dimension of the space spanned by its columns, and the **row rank** is the dimension of the space spanned by its rows.

- For $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\text{rank}(\mathbf{A}) \leq \min(m, n)$. If $\text{rank}(\mathbf{A}) = \min(m, n)$, then \mathbf{A} is said to be **full rank**, otherwise it is called **rank deficient**.
- For $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{A}^T) = \text{rank}(\mathbf{A}^T \mathbf{A}) = \text{rank}(\mathbf{A} \mathbf{A}^T)$.
- For $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{n \times p}$, $\text{rank}(\mathbf{AB}) \leq \min(\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B}))$.
- For $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$, $\text{rank}(\mathbf{A} + \mathbf{B}) \leq \text{rank}(\mathbf{A}) + \text{rank}(\mathbf{B})$.

A square matrix is invertible iff* it is full rank.

* “iff” is not a typo, it means “if and only if” referring to necessary and sufficient conditions

Special Matrices

Diagonal Matrix A diagonal matrix is a matrix where all non-diagonal elements are 0. This is typically denoted by $\mathbf{D} = \text{diag}(d_1, d_2, \dots, d_n)$, with

$$\mathbf{D} = \begin{pmatrix} d_1 & & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_n \end{pmatrix}$$

Identity matrix is a square matrix with $\mathbf{I} = \text{diag}(1, 1, \dots, 1) \in \mathbb{R}^{n \times n}$

Positive Definite Matrices A symmetric matrix (square) is positive definite iff all of its eigenvalues are positive. These are denoted by $\mathbf{A} \succ 0$. We also have **positive semi-definite** matrices where $\mathbf{A} \succcurlyeq 0$. Similarly we have negative definite ($\mathbf{A} \prec 0$) and negative semi-definite ($\mathbf{A} \preccurlyeq 0$) matrices. An example is a Gram matrix $\mathbf{G} = \mathbf{A}^\top \mathbf{A}$ which is always positive semi-definite.

Special Matrices

Orthogonal Matrices

- Two vectors $x, y \in \mathbb{R}^n$ are **orthogonal** if $x^\top y = 0$.
- A vector $x \in \mathbb{R}^n$ is **normalized** if $\|x\|_2 = 1$
- A set of vectors that is pairwise orthogonal and normalized is called **orthonormal**.
- A square matrix $U \in \mathbb{R}^{n \times n}$ is **orthogonal** (aka **orthonormal**) if all its columns are **orthonormal**.
- Note the difference in meaning of the term **orthogonal** for vectors and matrices.
- If the entries of a matrix U are complex valued, we use the term **unitary** instead of orthogonal. And for these matrices

$$U^\top U = I = UU^\top$$

and

$$U^\top = U^{-1}$$

Vector-Vector Product

For a matrix A and x , their product

In terms of inner product

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

In terms of outer product:

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

Matrix-Vector Product

For a matrix \mathbf{A} and \mathbf{x} , their product

In terms of inner product

$$\mathbf{y} = \mathbf{Ax} = \begin{bmatrix} - & \mathbf{a}_1^\top & - \\ - & \mathbf{a}_2^\top & - \\ & \vdots & \\ - & \mathbf{a}_m^\top & - \end{bmatrix} \mathbf{x} = \begin{bmatrix} \mathbf{a}_1^\top \mathbf{x} \\ \mathbf{a}_2^\top \mathbf{x} \\ \vdots \\ \mathbf{a}_m^\top \mathbf{x} \end{bmatrix}$$

In terms of outer product:

$$\mathbf{y} = \mathbf{Ax} = \begin{bmatrix} | & | & \cdots & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_n \\ | & | & & | \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} | \\ \mathbf{a}_1 \\ | \end{bmatrix} x_1 + \begin{bmatrix} | \\ \mathbf{a}_2 \\ | \end{bmatrix} x_2 + \cdots + \begin{bmatrix} | \\ \mathbf{a}_n \\ | \end{bmatrix} x_n$$

Matrix-Matrix Product

For two matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{n \times p}$, their product $\mathbf{C} = \mathbf{AB} \in \mathbb{R}^{m \times p}$ give by

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

In terms of inner product

$$\mathbf{C} = \mathbf{AB} = \begin{bmatrix} - & \mathbf{a}_1^\top & - \\ - & \mathbf{a}_2^\top & - \\ & \vdots & \\ - & \mathbf{a}_m^\top & - \end{bmatrix} \begin{bmatrix} | & | & & | \\ \mathbf{b}_1 & \mathbf{b}_2 & \cdots & \mathbf{b}_p \\ | & | & & | \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1^\top \mathbf{b}_1 & \mathbf{a}_1^\top \mathbf{b}_2 & \cdots & \mathbf{a}_1^\top \mathbf{b}_p \\ \mathbf{a}_2^\top \mathbf{b}_1 & \mathbf{a}_2^\top \mathbf{b}_2 & \cdots & \mathbf{a}_2^\top \mathbf{b}_p \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{a}_m^\top \mathbf{b}_1 & \mathbf{a}_m^\top \mathbf{b}_2 & \cdots & \mathbf{a}_m^\top \mathbf{b}_p \end{bmatrix}$$

In terms of outer product:

$$\mathbf{C} = \mathbf{AB} = \begin{bmatrix} | & | & \cdots & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_n \\ | & | & & | \end{bmatrix} \begin{bmatrix} - & \mathbf{b}_1^\top & - \\ - & \mathbf{b}_2^\top & - \\ & \vdots & \\ - & \mathbf{b}_n^\top & - \end{bmatrix} = \sum_{i=1}^n \mathbf{a}_i \mathbf{b}_i^\top$$

Matrix-Matrix Product

Properties of matrix multiplication.

Matrix multiplication is

- Associative: $(\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC})$
- Distributive: $\mathbf{A}(\mathbf{B} + \mathbf{C}) = \mathbf{AB} + \mathbf{AC}$

BUT NOT

- Commutative: $\mathbf{AB} \neq \mathbf{BA}$

Inverse of a Matrix

The inverse of a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is denoted \mathbf{A}^{-1} , and is the unique matrix such that

$$\mathbf{A}^{-1}\mathbf{A} = \mathbf{I} = \mathbf{A}\mathbf{A}^{-1}$$

Note that \mathbf{A}^{-1} exists if and only if $\det(\mathbf{A}) \neq 0$. If $\det(\mathbf{A}) = 0$, it is called a **singular matrix**. The following are properties of the inverse; all assume that $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ are non-singular:

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

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

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

Eigenvalue decomposition (EVD)

Review of some standard material on the eigenvalue decomposition or EVD of square (real-valued) matrices.

Given a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, we say that $\lambda \in \mathbb{R}$ is an eigenvalue of \mathbf{A} and $\mathbf{u} \in \mathbb{R}^n$ is the corresponding eigenvector if

$$\mathbf{A}\mathbf{u} = \lambda\mathbf{u}, \quad \mathbf{u} \neq \mathbf{0}$$

The **characteristic equation** is given by

$$\det(\lambda\mathbf{I} - \mathbf{A}) = 0$$

Eigenvalue decomposition (EVD)

Diagonalization We can write all the eigenvector equations simultaneously as

$$\mathbf{A}\mathbf{U} = \mathbf{U}\mathbf{\Lambda}$$

where the columns of $\mathbf{U} \in \mathbb{R}^{n \times n}$ are the eigenvectors of \mathbf{A} and $\mathbf{\Lambda}$ is a diagonal matrix whose entries are the eigenvalues of \mathbf{A} , i.e.,

$$\mathbf{U} \in \mathbb{R}^{n \times n} = \begin{bmatrix} | & | & \cdots & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_n \\ | & | & & | \end{bmatrix}, \quad \mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$$

If the eigenvectors of \mathbf{A} are linearly independent, then the matrix \mathbf{U} will be invertible, so we can write

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^{-1}$$

A matrix that can be written in this form is called **diagonalizable**.

Singular value decomposition (SVD)

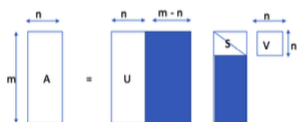
SVD generalizes EVD to rectangular matrices. Any (real) matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$ can be decomposed as

$$\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T = \sigma_1 \begin{pmatrix} | \\ \mathbf{u}_1 \\ | \end{pmatrix} \begin{pmatrix} - & \mathbf{v}_1^T & - \end{pmatrix} + \cdots + \sigma_r \begin{pmatrix} | \\ \mathbf{u}_r \\ | \end{pmatrix} \begin{pmatrix} - & \mathbf{v}_r^T & - \end{pmatrix}$$

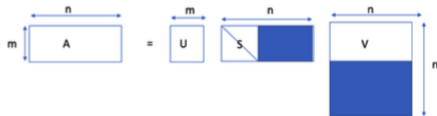
where

- $\mathbf{U} \in \mathbb{R}^{m \times m}$ is an orthonormal matrix ($\mathbf{U}^T \mathbf{U} = \mathbf{U} \mathbf{U}^T = \mathbf{I}_m$)
- $\mathbf{V} \in \mathbb{R}^{n \times n}$ is an orthonormal matrix ($\mathbf{V}^T \mathbf{V} = \mathbf{V} \mathbf{V}^T = \mathbf{I}_n$)
- $\mathbf{S} \in \mathbb{R}^{m \times n}$ contains $r = \min(m, n)$ singular values $\sigma_i \geq 0$ on the main diagonal, 0s filling the rest of the matrix

Singular value decomposition (SVD)



(a)



(b)

Figure 7.8: SVD decomposition of a matrix, $A = USV^T$. The shaded parts of each matrix are not computed in the economy-sized version. (a) Tall skinny matrix. (b) Short wide matrix.

Gradients

Consider a function that maps a vector to a scalar, $f : \mathbb{R}^n \rightarrow \mathbb{R}$. The gradient of a function at a point \mathbf{x} is the vector of its partial derivatives:

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

To emphasize the point at which the gradient is evaluated, we can write

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

- Operator “ ∇ ” (*nabla*) maps a function f to another function $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^n$.
- Since \mathbf{g} is a vector valued function, it is known as a **vector field**.

Directional Derivative

The directional derivative measures how much the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ changes along a direction \boldsymbol{v} in space. It is defined as

$$D_{\boldsymbol{v}} f(\boldsymbol{x}) = \nabla f(\boldsymbol{x}) \cdot \boldsymbol{v}$$

“ \cdot ” denotes the dot product or the inner product.

Jacobian

Consider a function that maps a vector to another vector, $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$. The Jacobian matrix of this function is a $m \times n$ matrix of partial derivatives:

$$\mathbf{J}_f(\mathbf{x}) = \frac{\partial \mathbf{f}}{\partial \mathbf{x}^\top} \triangleq \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{pmatrix} = \begin{pmatrix} \nabla f_1(\mathbf{x})^\top \\ \vdots \\ \nabla f_m(\mathbf{x})^\top \end{pmatrix}$$

Note that *here* we lay out the results in the same orientation as the output of f ; this is sometimes called **numerator layout** or the Jacobian formulation.

Hessian

Consider a twice differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, the Hessian matrix is a symmetric $n \times n$ matrix of second partial derivatives:

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

We see that the Hessian is the Jacobian of the gradient.

Functions that map scalars to scalars

Consider a differentiable function $f : \mathbb{R} \rightarrow \mathbb{R}$

$$\begin{aligned}\frac{d}{dx} [f(x) + g(x)] &= \frac{df(x)}{dx} + \frac{dg(x)}{dx} \\ \frac{d}{dx} [f(x)g(x)] &= f(x) \frac{dg(x)}{dx} + g(x) \frac{df(x)}{dx} \\ \frac{d}{dx} f(u(x)) &= \frac{du}{dx} \frac{df(u)}{du}\end{aligned}$$

The last relationship here is known as the **chain rule of calculus**. This is important for backpropagation algorithm to train neural networks.

Product Rule for Vectors

$$\frac{\partial u(\mathbf{x})^\top v(\mathbf{x})}{\partial \mathbf{x}} = \left[\frac{\partial v(\mathbf{x})}{\partial \mathbf{x}} \right] u(\mathbf{x}) + \left[\frac{\partial u(\mathbf{x})}{\partial \mathbf{x}} \right] v(\mathbf{x})$$

Useful Relations

Consider a differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$

$$\frac{\partial(\mathbf{a}^\top \mathbf{x})}{\partial \mathbf{x}} = \mathbf{a}$$

$$\frac{\partial(\mathbf{b}^\top \mathbf{A} \mathbf{x})}{\partial \mathbf{x}} = \mathbf{A}^\top \mathbf{b}$$

$$\frac{\partial(\mathbf{x}^\top \mathbf{A} \mathbf{x})}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}$$

Useful Relations

Consider a differentiable function $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$

$$\frac{\partial f}{\partial \mathbf{X}} = \begin{pmatrix} \frac{\partial f}{\partial x_{11}} & \cdots & \frac{\partial f}{\partial x_{1n}} \\ \vdots & & \\ \frac{\partial f}{\partial x_{m1}} & \cdots & \frac{\partial f}{\partial x_{mn}} \end{pmatrix}$$

Outline

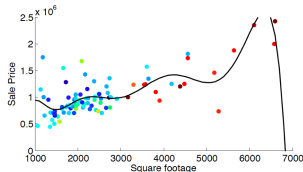
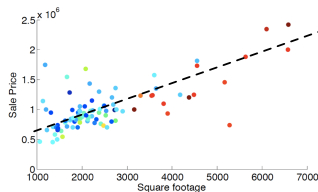
- 1 Math Background: Linear Algebra and Matrix Calculus
- 2 Math Background: Optimization for ML

How do we fit a model to a given dataset?

- The core problem in machine learning is parameter estimation.
- This requires solving an **optimization problem**, where we try to find the values for a set of variables $\theta \in \Theta$, that **minimize** a scalar-valued **loss function or cost function** $\mathcal{L} : \theta \rightarrow \mathbb{R}$

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} \mathcal{L}(\theta)$$

- The parameter space $\Theta \subset \mathbb{R}^D$, where D is the number of variables. Thus we are interested in **continuous optimization**.



Fitting a line (linear model) to Data. Fitting a polynomial (degree > 1)

Optimization

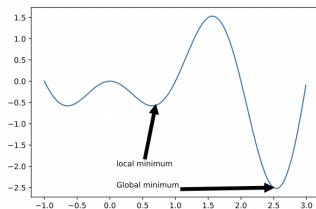
Definition: Optimization refers to choosing the best element from some set of available alternatives. A general form is as follows:

$$\begin{array}{ll} \text{minimize} & \mathcal{L}(\boldsymbol{\theta}) \\ \text{subject to} & f_i(\boldsymbol{\theta}) \leq 0, i = 1, \dots, m \\ & h_i(\boldsymbol{\theta}) = 0, i = 1, \dots, p. \end{array} \quad (1)$$

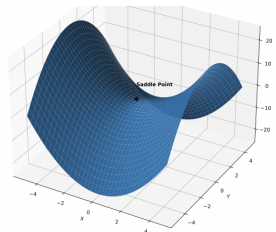
Difficulties:

- ❶ Local or global optimum?
- ❷ Difficulty to find a feasible point,
- ❸ Stopping criteria,
- ❹ Poor convergence rate,
- ❺ Numerical issues

Local and Global Optimization



(a)



(b)

Figure 8.1: (a) Illustration of local and global minimum in 1d. Generated by [extrema_fig_1d.ipynb](#). (b) Illustration of a saddle point in 2d. Generated by [saddle.ipynb](#).

Local and Global Optimization

- A point that satisfies the equation below is called a global optimum. Finding such a point is called global optimization

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} \mathcal{L}(\theta)$$

- This is in general intractable, and we just try to find a **local optimum**. For continuous problems, this is defined to be a point θ^* which has lower (or equal) cost than “nearby” points.
- A local minimum could be surrounded by other local minima with the same objective value; this is known as a **flat local minimum**.
- A point is said to be a **strict local minimum** if its cost is strictly lower than those of neighboring points.

Local and Global Optimization

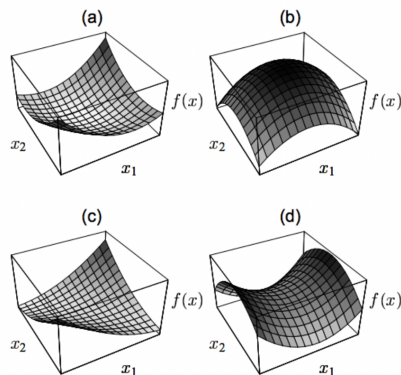


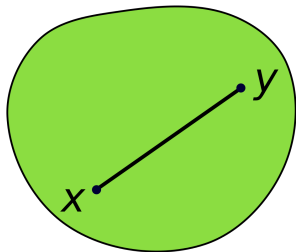
Figure 8.6: The quadratic form $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$ in 2d. (a) \mathbf{A} is positive definite, so f is convex. (b) \mathbf{A} is negative definite, so f is concave. (c) \mathbf{A} is positive semidefinite but singular, so f is convex, but not strictly. Notice the valley of constant height in the middle. (d) \mathbf{A} is indefinite, so f is neither convex nor concave. The stationary point in the middle of the surface is a saddle point. From Figure 5 of [She94].

Optimality Conditions for local vs global optima

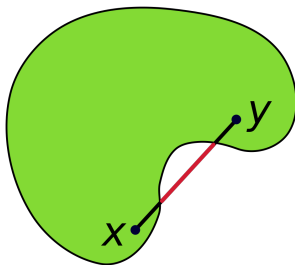
- For continuous, twice differentiable functions, we can precisely characterize the points which correspond to local minima.
- Let $\mathbf{g}(\boldsymbol{\theta}) = \nabla \mathcal{L}(\boldsymbol{\theta})$ be the gradient vector, and $\mathbf{H}(\boldsymbol{\theta}) = \nabla^2 \mathcal{L}(\boldsymbol{\theta})$ be the Hessian matrix.
- Consider a point $\boldsymbol{\theta}^* \in \mathbb{R}^D$, and let $\mathbf{g}^* = \mathbf{g}(\boldsymbol{\theta})|_{\boldsymbol{\theta}^*}$ be the gradient at that point, and $\mathbf{H}^* = \mathbf{H}(\boldsymbol{\theta})|_{\boldsymbol{\theta}^*}$ be the corresponding Hessian.
- One can show that the following conditions characterize every local minimum:
 - **Necessary condition:** If $\boldsymbol{\theta}^*$ is a local minimum, then we must have $\mathbf{g}^* = 0$ (i.e., $\boldsymbol{\theta}^*$ must be a stationary point), and \mathbf{H}^* must be positive semi-definite.
 - **Sufficient condition:** If $\mathbf{g}^* = 0$ and \mathbf{H}^* is positive definite, then $\boldsymbol{\theta}^*$ is a local optimum.

Convex Sets and Functions

Convex Set A set S is convex iff for any $x, y \in S$, $\lambda x + (1 - \lambda)y \in S$ for any $\lambda \in [0, 1]$,



Convex Polygon



Non-Convex Polygon

Convex Sets and Functions

Convex Functions: Function $f(\cdot)$ is a convex function, if for a convex set \mathcal{S} , any two points $\mathbf{x}, \mathbf{y} \in \mathcal{S}$ and any $\lambda \in [0, 1]$ satisfy

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}).$$

Concave Functions A function f is said to be *concave* if $-f$ is convex.

Non-convex Functions A function which is neither convex nor concave.

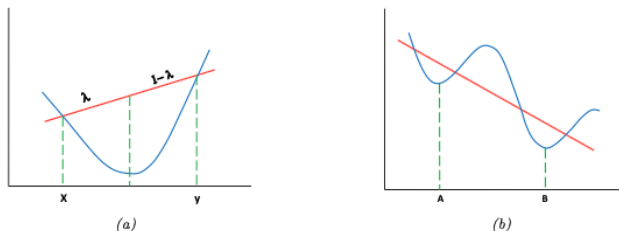
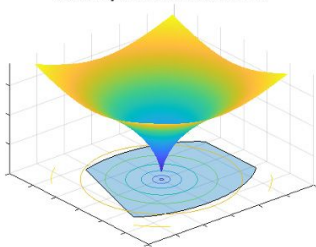


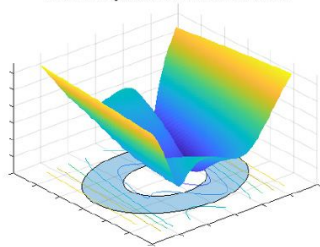
Figure 8.5: (a) Illustration of a convex function. We see that the chord joining $(x, f(x))$ to $(y, f(y))$ lies above the function. (b) A function that is neither convex nor concave. **A** is a local minimum, **B** is a global minimum.

Convex Sets and Functions

Convex Objective and Convex Constraints



Nonconvex Objective and Nonconvex Constraints



Convex and Non-convex Objectives (image source)

Convex Optimization

Convex Optimization is minimization (maximization) of a convex (concave) function over a convex set, e.g., Linear Programming (LP), Quadratic Programming (QP), and Semi-Definite Programming (SDP).

A number of popular convex optimization algorithms:

- Gradient descent
- Conjugate gradient
- Newton's method
- Quasi-Newton method
- Subgradient method

How do we know if a function f is convex?

Theorem 8.1.1. *Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice differentiable over its domain. Then f is convex iff $\mathbf{H} = \nabla^2 f(\mathbf{x})$ is positive semi definite (Section 7.1.5.3) for all $\mathbf{x} \in \text{dom}(f)$. Furthermore, f is strictly convex if \mathbf{H} is positive definite.*

For example, consider the quadratic form

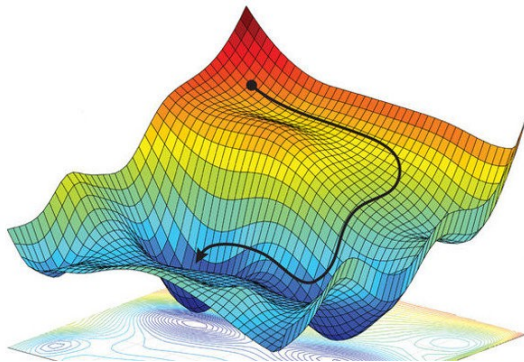
$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$$

This is

- **convex** if \mathbf{A} is *positive semi definite*, and
- **strictly convex** if \mathbf{A} is *positive definite*
- It is neither convex nor concave if \mathbf{A} has eigenvalues of mixed sign.

Deep learning involves non-convex optimization

BUT it is extremely hard to get any type of guarantees for non-convex optimizations. Therefore, convex optimization is still very popular and practical in a wide range of domains!



A non-convex function (image source)

Gradient and Stochastic Gradient Descent

We describe two simple yet extremely popular methods

- **Gradient Descent (GD)**: simple and fundamental
- **Stochastic Gradient Descent (SGD)**: faster, effective for large-scale problems

Gradient is sometimes referred to as *first-order* information of a function. Therefore, these methods are called *first-order methods*.

Gradient Descent (GD)

Goal: minimize a function $F(\mathbf{w})$. For instance, the loss function with the choice of \mathbf{w}

Algorithm: keep moving in the *negative gradient direction*

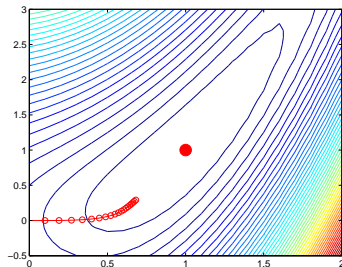
Start from some $\mathbf{w}^{(0)}$. For $t = 0, 1, 2, \dots$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla F(\mathbf{w}^{(t)})$$

where $\eta > 0$ is called step size or learning rate

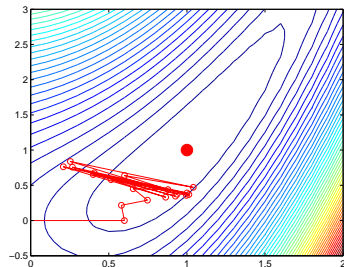
- in theory η should be set in terms of some parameters of F
- in practice we just try several small values

What does GD look like?



reasonable η decreases function value

Slow: each update makes one pass of the entire training set!



but large η is unstable

Stochastic Gradient Descent (SGD)

GD: keep moving in the negative gradient direction

SGD: keep moving in some *noisy* negative gradient direction

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \tilde{\nabla} F(\mathbf{w}^{(t)})$$

where $\tilde{\nabla} F(\mathbf{w}^{(t)})$ is a random variable (called **stochastic gradient**) s.t.

$$\mathbb{E} \left[\tilde{\nabla} F(\mathbf{w}^{(t)}) \right] = \nabla F(\mathbf{w}^{(t)}) \quad (\text{unbiasedness})$$

The gradient calculation uses only a sample, is therefore “noisy”, variants use a mini-batch of samples.

Key point: it could be *much faster to obtain a stochastic gradient!*

Convergence Guarantees

Many for both GD and SGD on convex objectives.

They tell you at most how many iterations you need to achieve

$$F(\mathbf{w}^{(t)}) - F(\mathbf{w}^*) \leq \epsilon$$

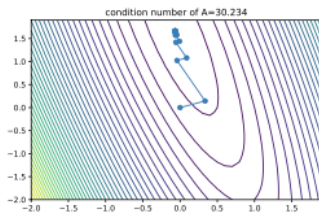
Even for *non-convex objectives* such as those encountered in deep learning, many recent works show effectiveness of GD/SGD.

Convergence also depends on the Condition Number κ

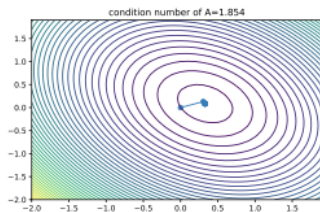
Condition number for a matrix \mathbf{A} is defined as

$$\kappa(\mathbf{A}) = \frac{\sigma_{\max}(\mathbf{A})}{\sigma_{\min}(\mathbf{A})}.$$

Condition number of the Hessian is a great indicator of how the landscape of the function (to be optimized) looks like. The closer to 1 the more *well conditioned* the matrix.



(a)



(b)

Figure 8.12: Illustration of the effect of condition number κ on the convergence speed of steepest descent with exact line searches. (a) Large κ . (b) Small κ . Generated by [lineSearchConditionNum.ipynb](#).

Reading

[PiML 1]

- 7.1.2.1, 7.1.2.1, 7.1.2.4, 7.1.3, 7.1.4.1, 7.1.4.2, 7.1.4.3, 7.1.5.1, 7.1.5.3, 7.1.5.4
- 7.2.1-3
- 7.3.1
- 7.4.1-2
- 7.5.1
- 7.8.2, 7.8.3, 7.8.5 (intro), 7.8.6, 7.8.7
- 8.1 (introduction), 8.1.1 (introduction), 8.1.1.1, 8.1.2 (material from slides only), 8.1.3 (8.1.3.1-8.1.3.3)

For the subsections, results covered in the slides are important, the rest will be covered on-need basis.

References I