

SIT744 Lecture 2

Math Review

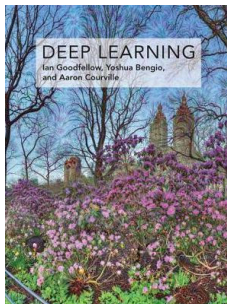
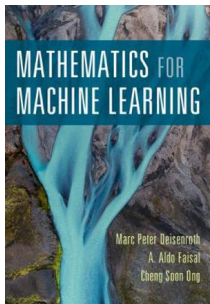
Nayyar Zaidi

Learning objectives

- Be familiar with the math behind deep learning theory

Plan

- Linear Algebra
- Derivatives and gradients



This lecture is based on chapters in these books (click to view):

- Mathematics for Machine Learning
 - Sections 5.1, 5.2, and 5.6
- Deep Learning Book
 - Chapter 2

Why linear algebra?

Linear algebra studies vector spaces and matrices

Tensors for efficient computing

- **Data** are represented as vectors, matrices, and tensors
- **Model parameters** are represented as vectors, matrices, and tensors

Scalars, Vectors, Matrices and Tensors

Scalar

A single number

Vector

An array of numbers

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

Matrices

A 2-D array

$$\begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix}$$

Tensors

An array with more than two axes

Data as vectors

Learning activity

Think about a machine learning problem. How is the data represented in the problem?

Data as vectors

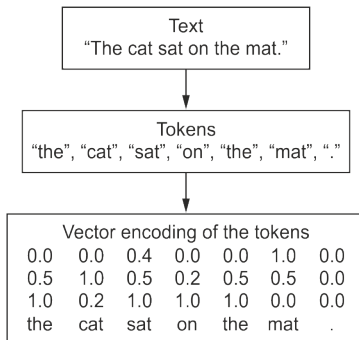


Figure 1: From text to tokens to vectors

Words represented as vectors



Figure 2: MNIST sample digits

Images are matrices

Dot product

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

Euclidean norm

$$\|\mathbf{x}\|_2 := \sqrt{\sum_{i=1}^n x_i^2} = \sqrt{\mathbf{x}^\top \mathbf{x}}$$

Other norms

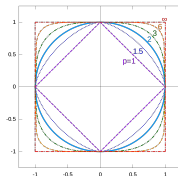


Figure 3: unit circle for different p-norms (src: wikipedia)

- Manhattan norm (L1 norm)

$$\|\mathbf{x}\|_1 := \sum_{i=1}^n |x_i|$$

- Lp-norm (if $p \geq 1$)

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

Lengths and distance

$\|\mathbf{x}\|_2$ is one way to measure the length of \mathbf{x} .

Distance between vectors

$$d(\mathbf{x}, \mathbf{y}) := \|\mathbf{x} - \mathbf{y}\| = \sqrt{\langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle}$$

Angle between vectors

$$\cos \omega = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

It is a way to measure similarity of two vectors.

Transpose of a matrix

The transpose of A , denoted A^T , is defined by

$$(A^T)_{i,j} = A_{j,i}$$

example

$$A = \begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,3} \\ A_{3,1} & A_{3,2} \end{bmatrix} \Rightarrow A^T = \begin{bmatrix} A_{1,1} & A_{2,1} & A_{3,1} \\ A_{1,2} & A_{2,2} & A_{3,2} \end{bmatrix}$$

Sum of two matrices

If A and B have the same shape, $\mathbf{C} = \mathbf{A} + \mathbf{B}$ is defined by

$$C_{i,j} = A_{i,j} + B_{i,j}$$

Scalar multiplication and sum

$\mathbf{D} = a \cdot \mathbf{B} + c$ is defined by

$$D_{i,j} = a \cdot B_{i,j} + c$$

A confusing notation common in deep learning

$\mathbf{C} = \mathbf{A} + \mathbf{b}$ is defined by

$$C_{i,j} = A_{i,j} + b_j$$

Matrix product

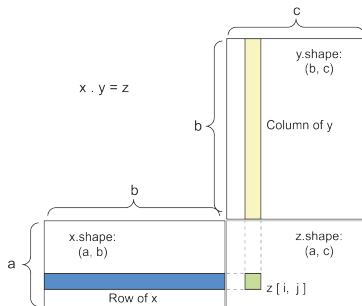


Figure 4: Matrix product

$\mathbf{C} = \mathbf{AB}$ is defined by

$$C_{i,j} = \sum_k A_{i,k} B_{k,j}$$

Matrix product is not commutative

$$\mathbf{AB} \neq \mathbf{BA}$$

However,

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

It is easier to think of matrix product as the composition of two functions.

Identity and Inverse Matrices

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

Here \mathbf{A}^{-1} is the **matrix inverse** of \mathbf{A} .

If $\mathbf{Ax} = \mathbf{b}$ and \mathbf{A}^{-1} exist, then

$$\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$$

The L^p norm of a vector \mathbf{x} is defined by

$$\|\mathbf{x}\|_p = \left(\sum_i |x_i|^p \right)^{\frac{1}{p}}$$

Common L^p norms

- $\|\mathbf{x}\|_2$
- $\|\mathbf{x}\|_1$
- $\|\mathbf{x}\|_\infty = \max_i |x_i|$

Eigendecomposition

Let \mathbf{A} be a square matrix. A nonzero vector \mathbf{v} is an **eigenvector** of \mathbf{A} if

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}.$$

Here λ is a scalar and is called the **eigenvalue** for \mathbf{v} .

The **eigendecomposition** of \mathbf{A} is then given by

$$\mathbf{A} = \mathbf{V} \text{diag}(\lambda) \mathbf{V}^{-1}.$$

Eigendecomposition for real symmetric matrix

$$\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$$

optimize quadratic expressions

$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x} \text{ subject to } \|\mathbf{x}\|_2 = 1$$

- $\max f$ is the maximum eigenvalue
- $\min f$ is the minimum eigenvalue

Machine learning and optimisation

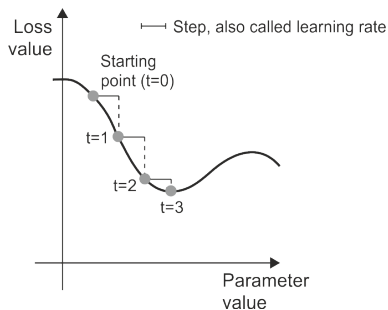


Figure 5: Learning as optimisation

Most machine learning algorithms are optimisation algorithms

- Find an \mathbf{x} to **minimise** $f(\mathbf{x})$, the loss function.

Other names for the loss function

- Cost function
- Error function

Global minimum

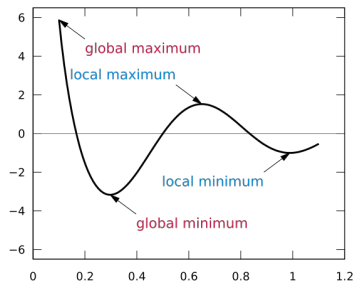


Figure 6: Global minimum (src: wikipedia)

Global minimum is realised at

$$\mathbf{x}^* = \arg \min f(\mathbf{x})$$

Derivative and gradient

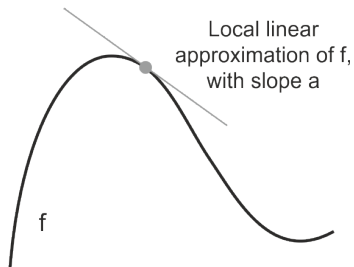
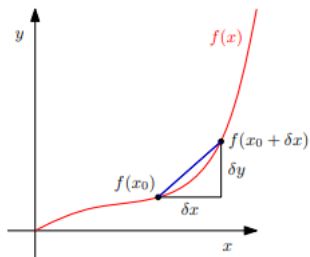


Figure 7: Derivative

Derivative $\frac{df(x)}{dx}$ (or $f'(x)$) measures the slope of f at x .

Derivative



$$\frac{df}{dx} := \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

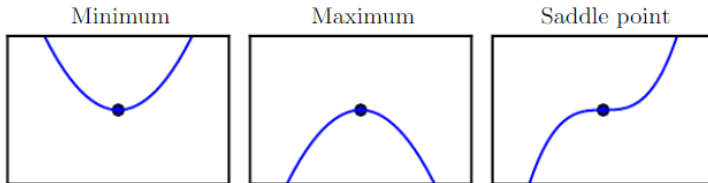
Gradient is for multi-variate functions

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(\mathbf{x}) \\ \vdots \\ \frac{\partial f}{\partial x_n}(\mathbf{x}) \end{bmatrix}$$

Partial derivative

$$\frac{\partial f(\mathbf{x})}{\partial x_1} = \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(\mathbf{x})}{h}$$

Stationary points



Consider a single variable. If x^* renders a global minimum, then the derivative $f'(x^*) = 0$.

Gradient-Based optimisation

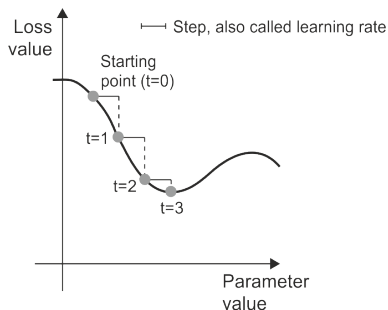


Figure 8: Method of steepest descent

Iteratively update a guess until converging to a stationary point

$$x' = x - \epsilon f'(x)$$

- ϵ is the learning rate.

Gradient descent

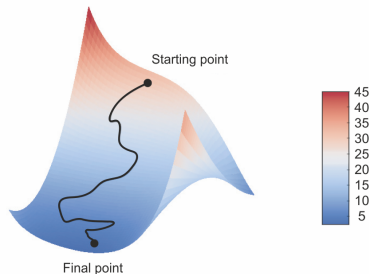


Figure 9: Steepest descent with two variables

$$\mathbf{x}' = \mathbf{x} - \epsilon \nabla_{\mathbf{x}} f(\mathbf{x})$$

Differentiation Rules

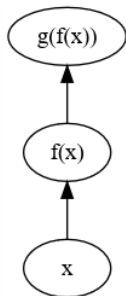


Figure 10: Chain rule

Chain rule

$$(g(f(x)))' = (g \circ f)'(x) = g'(f(x))f'(x)$$

Differentiation Rules

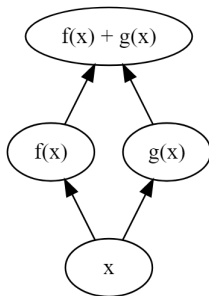


Figure 11: Sum rule

Sum rule

$$(f(x) + g(x))' = f'(x) + g'(x)$$

Differentiation Rules

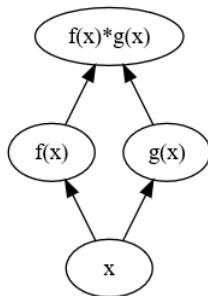


Figure 12: Product rule

Product rule

$$(f(x)g(x))' = f'(x)g(x) + f(x)g'(x)$$

Derivatives of common functions



$$\frac{d}{dx}123 = 0$$



$$\frac{d}{dx}x = 1$$



$$\frac{\partial}{\partial w}wx = x$$



$$\frac{\partial}{\partial b}(wx + b) = 1$$



$$\frac{d}{dx}e^x = e^x$$

Example

What is

$$\frac{d}{dx}x^2$$

- Hint: Product rule

Example

$$\frac{d}{dx}x^2 = \frac{d}{dx}(x \cdot x) = x \frac{d}{dx}x + x \frac{d}{dx}x = 2x$$

Example

How about

$$\frac{d}{dx} x^n$$

Example

$$\frac{d}{dx}x^n = n\frac{d}{dx}x^{n-1}$$

Special case

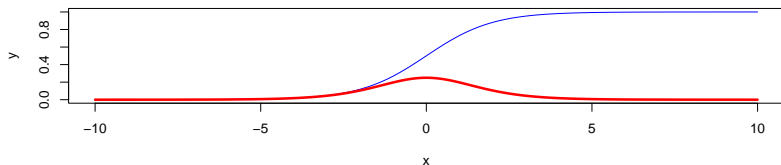
$$\frac{d}{dx}x^{-1} = -\frac{d}{dx}x^{-2}$$

Example

What is

$$\frac{d}{dx} \left(\frac{e^x}{e^x + 1} \right)$$

Example



$$\frac{d}{dx} \left(\frac{e^x}{e^x + 1} \right) = \frac{e^x}{(e^x + 1)^2}$$

Vanishing Gradients Problem

$\frac{d}{dx} \left(\frac{e^x}{e^x + 1} \right)$ is getting very small with a moderately sized x .

- $\frac{e^{10}}{(e^{10} + 1)^2} = 0.000045$

Chain rule implies that the gradient quickly vanishes with more than one layers with the sigmoid activation function.

Chain rule with computational graph

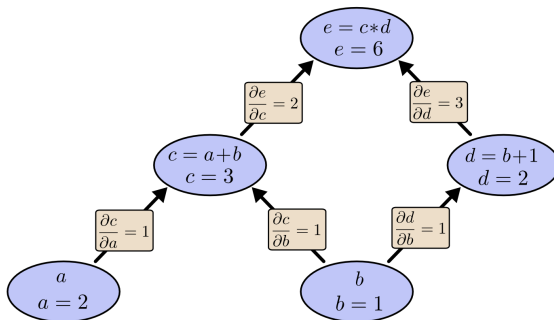


Figure 13: Backpropagation (src:colah.github.io)

Backpropagation is the default algorithm for computing gradients in deep learning.

Backpropagation is simply applying chain-rule on the computational graph.

Summary

- Linear algebra and multivariable differential calculus are fundamental math tools for deep learning
- Linear algebra is useful for representing data and transformations
- Differential calculus is useful for training deep learning models