

Chương 2

Đại số tuyến tính

2.1	Đại lượng vô hướng, vectơ, ma trận, và tenxơ	4
2.2	Nhân các vectơ và các ma trận	7
2.3	Ma trận đơn vị và ma trận nghịch đảo . . .	8
2.4	Độc lập tuyến tính và không gian bao tuyến tính	10
2.5	Chuẩn	12
2.6	Các loại ma trận và vectơ đặc biệt	13
2.7	Phân tích giá trị riêng	15
2.8	Phân tích giá trị kỳ dị	17
2.9	Ma trận giả nghịch đảo Moore–Penrose . .	18
2.10	Vết của ma trận	19
2.11	Định thức	20
2.12	Ví dụ: phân tích thành phần chính	20
2.13	Không gian Euclid	24

Đại số tuyến tính là một nhánh của toán học được sử dụng rộng rãi trong khoa học và kỹ thuật. Tuy nhiên, vì đại số tuyến tính là một dạng toán học liên tục chứ không phải toán rời rạc, nên nhiều nhà khoa học máy tính có ít kinh nghiệm với nó. Hiểu biết tốt về đại số tuyến tính là rất cần thiết để hiểu và làm việc với nhiều thuật toán học máy, đặc biệt là các thuật toán học sâu. Do đó, trước khi giới thiệu về học máy, ta cần tập trung vào các kiến thức nền tảng quan trọng về đại số tuyến tính.

Nếu bạn đã quen thuộc với đại số tuyến tính, hãy thoải mái bỏ qua chương này. Nếu bạn đã có kinh nghiệm với các khái niệm này nhưng cần một tài liệu tham khảo chi tiết để ôn lại các công thức chính, bạn có thể đọc cuốn *The Matrix Cookbook* (Petersen và Pedersen, 2006, [1]). Nếu bạn chưa từng tiếp xúc với đại số tuyến tính, chương này sẽ dạy bạn đủ để đọc cuốn sách này, nhưng bạn nên tham khảo thêm các tài liệu khác chỉ tập trung vào việc giảng dạy đại số tuyến tính, chẳng hạn như *Linear Algebra* (Shilov, 1977, [2]). Chương này sẽ bỏ qua nhiều chủ đề quan trọng của đại số tuyến tính không cần thiết để hiểu học máy.

2.1 Đại lượng vô hướng, vectơ, ma trận, và tenxơ

Một số đối tượng nghiên cứu của đại số tuyến tính:

- **Đại lượng vô hướng:** Đại lượng vô hướng đơn giản chỉ là một số đơn lẻ, khác với hầu hết các đối tượng khác được nghiên cứu trong đại số tuyến tính, thường là các mảng chứa nhiều số. Ta hay viết các số vô hướng bằng chữ cái thường, in nghiêng. Khi giới thiệu, ta sẽ chỉ rõ loại số của chúng. Ví dụ, ta nói “Giả sử $s \in \mathbb{R}$ là độ dốc của đường thẳng”, khi định nghĩa một số vô hướng có giá trị thực, hoặc “Giả sử $n \in \mathbb{N}$ là số lượng dữ liệu”, khi định nghĩa một số vô hướng tự nhiên.
- **Vectơ:** Vectơ là một mảng gồm các số. Các số này được sắp xếp theo thứ tự. Ta có thể xác định từng số riêng lẻ bằng chỉ số của nó trong thứ tự đó. Ta hay đặt tên cho các vectơ bằng chữ cái thường, in đậm, chẳng hạn như \mathbf{x} . Các phần tử của vectơ được xác định bằng cách viết tên của nó dưới dạng chữ in nghiêng, kèm theo chỉ số dưới. Phần tử đầu tiên của \mathbf{x} là x_1 , phần tử thứ hai là x_2 , và cứ tiếp tục như vậy. Ta cũng cần nói rõ loại số nào được lưu trữ trong vectơ. Nếu mỗi phần tử nằm trong \mathbb{R} , và vectơ có n phần tử, thì vectơ nằm trong tập hợp được tạo bởi tích Descartes của \mathbb{R} với chính nó n lần, ký hiệu là \mathbb{R}^n . Khi cần xác định rõ ràng các phần tử của vectơ, ta viết chúng dưới dạng một cột được bao quanh bởi dấu ngoặc vuông:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}. \quad (2.1)$$

Ta có thể xem vectơ như là xác định các điểm trong không gian, với mỗi phần tử cung cấp tọa độ theo một trục khác nhau. Đôi khi, ta cần truy cập một tập hợp các phần tử của vectơ. Trong trường hợp này, ta định nghĩa một tập chứa các chỉ số và viết tập hợp đó dưới dạng chỉ số dưới. Ví dụ, để truy cập x_1 , x_3 và x_6 , ta định nghĩa tập hợp $S = \{1, 3, 6\}$ và viết là \mathbf{x}_S . Ta sử dụng dấu “—” để chỉ phần bù của một tập hợp. Ví dụ, \mathbf{x}_{-1} là vectơ chứa tất cả các phần tử của \mathbf{x} ngoại trừ x_1 , và \mathbf{x}_{-S} là vectơ chứa tất cả các phần tử của \mathbf{x} ngoại trừ x_1 , x_3 và x_6 .

- **Ma trận:** Ma trận là một mảng 2 chiều gồm các số, vì vậy mỗi phần tử được xác định bởi hai chỉ số thay vì chỉ một. Ta thường đặt tên cho ma trận bằng các chữ cái viết hoa và in đậm, chẳng hạn như \mathbf{A} . Nếu ma trận \mathbf{A} các phần tử giá trị thực có cỡ $m \times n$, tức là có m hàng và n cột, thì ta nói rằng $\mathbf{A} \in \mathbb{R}^{m \times n}$. Ta thường xác định các phần tử của ma trận bằng cách sử dụng tên của nó ở dạng in nghiêng nhưng không in đậm, và các chỉ số được liệt kê có dấu phẩy ngăn cách, hoặc có thể bỏ dấu phẩy nếu không gây nhầm lẫn. Ví dụ, $A_{1,1}$, hay A_{11} , là phần tử ở góc trên bên trái của \mathbf{A} , được đánh dấu hàng 1 cột 1, và $A_{m,n}$, hay A_{mn} , là phần tử ở góc dưới bên phải, ở hàng m cột n . Ta có thể xác định tất cả các số ở hàng i bằng cách viết dấu “:” thay cho chỉ số cột. Ví dụ, $\mathbf{A}_{i,:}$ biểu thị hàng ngang thứ i của ma trận \mathbf{A} . Tương tự, $\mathbf{A}_{:,j}$ là cột thứ j của \mathbf{A} . Khi cần xác định rõ ràng các phần tử của ma trận, ta viết chúng dưới dạng mảng và đặt trong dấu ngoặc vuông:

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \\ A_{31} & A_{32} \end{bmatrix} \quad (2.2)$$

Đôi khi ta có thể cần đánh chỉ số các biểu thức nhận giá trị ma trận không chỉ là một chữ cái. Trong trường hợp này, ta sử dụng chỉ số dưới bên phải biểu thức, nhưng không chuyển bất kỳ thứ gì thành chữ thường. Ví dụ, $f(\mathbf{A})_{ij}$ cho phần tử (i, j) của ma trận được tính bằng cách áp dụng hàm số f lên \mathbf{A} .

- **Tenxơ:** Trong một số trường hợp, ta cần một mảng với nhiều hơn hai trục. Trong trường hợp tổng quát, một mảng số được sắp xếp trên một lưới đều với số lượng trục thay đổi được gọi là tenxơ. Ta ký hiệu một tenxơ có tên là “A” bằng kiểu chữ này: \mathcal{A} . Ta xác định phần tử của \mathcal{A} tại tọa độ (i, j, k) bằng cách viết \mathcal{A}_{ijk} .

Một phép toán quan trọng trên ma trận là **chuyển vị**. Chuyển vị của một ma trận là hình ảnh phản chiếu của ma trận qua một đường chéo, gọi là **đường chéo chính**, chạy từ trên xuống dưới và sang phải, bắt đầu từ góc trên bên trái. [Hình 2.1](#) mô tả trực quan phép toán này.

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \\ A_{31} & A_{32} \end{bmatrix} \Rightarrow \mathbf{A}^T = \begin{bmatrix} A_{11} & A_{21} & A_{31} \\ A_{12} & A_{22} & A_{32} \end{bmatrix}$$

Hình 2.1: Chuyển vị của ma trận có thể được xem như là hình ảnh phản chiếu qua đường chéo chính.

Ký hiệu chuyển vị của ma trận \mathbf{A} là \mathbf{A}^T , và nó được định nghĩa như sau

$$(\mathbf{A}^T)_{ij} = A_{ji} \quad (2.3)$$

Vectơ có thể được xem như ma trận chỉ chứa một cột. Do đó, chuyển vị của một vectơ là một ma trận chỉ có một hàng. Đôi khi, ta định nghĩa một vectơ bằng cách viết các phần tử của nó dưới dạng một ma trận hàng, sau đó sử dụng phép toán chuyển vị để biến nó thành một vectơ cột dạng chuẩn, ví dụ: $\mathbf{x} = [x_1, x_2, x_3]^T$, hoặc $\mathbf{x} = (x_1, x_2, x_3)^T$.

Một đại lượng vô hướng có thể được xem như một ma trận chỉ có một phần tử. Từ đó, ta có thể thấy rằng một đại lượng vô hướng chính là chuyển vị của chính nó: $a = a^T$.

Ta có thể cộng các ma trận với nhau, miễn là chúng có cùng cỡ, bằng cách cộng các phần tử tương ứng của chúng: $\mathbf{C} = \mathbf{A} + \mathbf{B}$, trong đó $C_{ij} = A_{ij} + B_{ij}$.

Ta cũng có thể cộng một đại lượng vô hướng với ma trận hoặc nhân ma trận với đại lượng vô hướng, bằng cách thực hiện phép toán đó trên từng phần tử của ma trận: $\mathbf{D} = a \cdot \mathbf{B} + c$, trong đó $D_{ij} = a \cdot B_{ij} + c$.

Trong ngữ cảnh học sâu, ta cũng sử dụng một số ký hiệu ít tính quy ước hơn. Ta cho phép cộng ma trận với vectơ, tạo ra một ma trận khác: $\mathbf{C} = \mathbf{A} + \mathbf{b}$, trong đó $C_{ij} = A_{ij} + b_j$. Nói cách khác, vectơ \mathbf{b} được cộng vào từng hàng của ma trận. Cách viết tắt này loại bỏ nhu cầu phải định nghĩa một ma trận với \mathbf{b} được sao chép vào mỗi hàng trước khi thực hiện phép cộng. Việc sao chép ngầm \mathbf{b} vào nhiều vị trí này được gọi là **phát sóng** (broadcasting).

2.2 Nhân các vectơ và các ma trận

Một trong những phép toán quan trọng nhất liên quan đến ma trận là phép nhân hai ma trận. Tích ma trận của các ma trận \mathbf{A} và \mathbf{B} là một ma trận thứ ba \mathbf{C} . Để phép nhân này được xác định, ma trận \mathbf{A} phải có số cột bằng với số hàng của ma trận \mathbf{B} . Nếu \mathbf{A} có cỡ $m \times n$ và \mathbf{B} có cỡ $n \times p$, thì \mathbf{C} sẽ có cỡ $m \times p$. Ta viết tích ma trận bằng cách đặt hai hoặc nhiều ma trận cạnh nhau, ví dụ

$$\mathbf{C} = \mathbf{AB}. \quad (2.4)$$

Phép toán nhân được định nghĩa bởi

$$C_{ij} = \sum_{k=1}^n A_{ik} B_{kj}. \quad (2.5)$$

Lưu ý rằng tích chuẩn của hai ma trận không chỉ đơn thuần là một ma trận chứa tích của các phần tử riêng lẻ. Phép toán như vậy cũng có và được gọi là **tích từng phần tử** hay **tích Hadamard**, và được ký hiệu là $\mathbf{A} \odot \mathbf{B}$.

Tích vô hướng giữa hai vectơ \mathbf{x} và \mathbf{y} có cùng số chiều là tích ma trận $\mathbf{x}^T \mathbf{y}$. Ta có thể xem tích ma trận $\mathbf{C} = \mathbf{AB}$ bằng cách tính C_{ij} là tích vô hướng giữa hàng thứ i của \mathbf{A} và cột thứ j của \mathbf{B} .

Các phép toán nhân ma trận có nhiều tính chất hữu ích giúp cho việc phân tích toán học các ma trận trở nên thuận tiện hơn. Ví dụ, phép nhân ma trận có tính chất kết hợp:

$$(\mathbf{AB}) \mathbf{C} = \mathbf{A} (\mathbf{BC}). \quad (2.6)$$

Nó cũng có tính chất phân phối:

$$\begin{aligned} \mathbf{A} (\mathbf{B} + \mathbf{C}) &= \mathbf{AB} + \mathbf{AC} \\ (\mathbf{A} + \mathbf{B}) \mathbf{C} &= \mathbf{AC} + \mathbf{BC}. \end{aligned} \quad (2.7)$$

Tuy nhiên, phép nhân ma trận không có tính chất giao hoán (điều kiện $\mathbf{AB} = \mathbf{BA}$ không phải lúc nào cũng đúng), khác với phép nhân vô hướng. Tuy nhiên, tích vô hướng giữa hai vectơ thì có tính giao hoán:

$$\mathbf{x}^T \mathbf{y} = \mathbf{y}^T \mathbf{x} \quad (2.8)$$

Chuyển vị của tích ma trận có một tính chất đơn giản:

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T. \quad (2.9)$$

nhân vectơ đó với ma trận này. Ta nói ma trận đơn vị bảo toàn các vectơ n chiều và ký hiệu là I_n . Như vậy, $I_n \in \mathbb{R}^{n \times n}$, và

$$\forall \mathbf{x} \in \mathbb{R}^n, \quad I_n \mathbf{x} = \mathbf{x}. \quad (2.13)$$

Cấu trúc của ma trận đơn vị rất đơn giản: tất cả các phần tử trên đường chéo chính đều bằng 1, trong khi tất cả các phần tử khác đều bằng 0, như hình dưới đây:

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Hình 2.2: Ví dụ về **ma trận đơn vị**: Đây là I_3

Với ma trận \mathbf{A} cỡ $m \times n$, nghịch đảo của \mathbf{A} được ký hiệu là \mathbf{A}^{-1} , và nó được định nghĩa là ma trận sao cho

$$\mathbf{A}^{-1} \mathbf{A} = I_n \quad (2.14)$$

Ta có thể giải phương trình (2.10) theo các bước sau:

$$\begin{aligned} \mathbf{Ax} &= \mathbf{b} \\ \mathbf{A}^{-1} (\mathbf{Ax}) &= \mathbf{A}^{-1} \mathbf{b} \\ (\mathbf{A}^{-1} \mathbf{A}) \mathbf{x} &= \mathbf{A}^{-1} \mathbf{b} \\ I_n \mathbf{x} &= \mathbf{A}^{-1} \mathbf{b} \\ \mathbf{x} &= \mathbf{A}^{-1} \mathbf{b} \end{aligned} \quad (2.15)$$

Tất nhiên, quá trình này phụ thuộc vào khả năng tìm được \mathbf{A}^{-1} . Ta sẽ thảo luận về các điều kiện để \mathbf{A}^{-1} tồn tại trong phần sau.

Khi \mathbf{A}^{-1} tồn tại, có một số thuật toán để tìm nó dưới dạng đóng. Về lý thuyết, cùng một ma trận nghịch đảo có thể được sử dụng để giải phương trình nhiều lần cho các giá trị khác nhau của \mathbf{b} . Tuy nhiên, \mathbf{A}^{-1} chủ yếu hữu ích như một công cụ lý thuyết và không nên thực sự được sử dụng trong hầu hết các phần mềm. Vì \mathbf{A}^{-1} chỉ có thể được biểu diễn với độ chính xác giới hạn trên máy tính, các thuật toán sử dụng cả giá trị của \mathbf{b} thường có thể đưa ra ước tính chính xác hơn về \mathbf{x} , chẳng hạn **thuật toán Gauss–Seidel**.

2.4 Độc lập tuyến tính và không gian bao tuyến tính

Để \mathbf{A}^{-1} tồn tại, phương trình (2.10) phải có đúng một nghiệm cho mọi giá trị của \mathbf{b} . Tuy nhiên, cũng có thể hệ phương trình vô nghiệm hoặc có vô số nghiệm cho một số giá trị của \mathbf{b} . Không thể có trường hợp phương trình có nhiều hơn một nghiệm nhưng ít hơn vô số nghiệm cho một giá trị cụ thể của \mathbf{b} ; vì nếu cả \mathbf{x} và \mathbf{y} đều là nghiệm, thì

$$\mathbf{z} = \alpha \mathbf{x} + (1 - \alpha) \mathbf{y} \quad (2.16)$$

cũng là một nghiệm với số thực α bất kỳ.

Để phân tích phương trình có bao nhiêu nghiệm, ta có thể xem các cột của \mathbf{A} chỉ các hướng khác nhau mà ta có thể di chuyển từ gốc tọa độ (điểm được xác định bởi vectơ gồm toàn số 0), và xác định có bao nhiêu cách để đạt đến \mathbf{b} . Theo cách nhìn này, mỗi phần tử của \mathbf{x} xác định ta nên di chuyển bao xa theo từng hướng, với x_i xác định khoảng cách di chuyển theo hướng của cột thứ i :

$$\mathbf{Ax} = \sum_{i=1}^n x_i \mathbf{A}_{:,i} \quad (2.17)$$

Nhìn chung, phép toán kiểu này được gọi là **tổ hợp tuyến tính**. Tổng quát, một tổ hợp tuyến tính của một tập hợp các vectơ $\{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)}\}$ được xác định bằng cách nhân mỗi vectơ $\mathbf{v}^{(i)}$ với một hệ số vô hướng tương ứng và cộng các kết quả lại:

$$\sum_{i=1}^n c_i \mathbf{v}^{(i)} \quad (2.18)$$

Bao tuyến tính sinh bởi một tập các vectơ là tập hợp tất cả các điểm có thể đạt được bằng cách lấy tổ hợp tuyến tính của các vectơ ban đầu.

Việc xác định liệu phương trình $\mathbf{Ax} = \mathbf{b}$ có nghiệm hay không tương đương với việc kiểm tra xem \mathbf{b} có nằm trong tập sinh bởi các cột của \mathbf{A} hay không. Tập hợp sinh này được gọi là **không gian cột** hoặc **miền giá trị** của \mathbf{A} .

Để hệ phương trình $\mathbf{Ax} = \mathbf{b}$ có nghiệm cho mọi giá trị của $\mathbf{b} \in \mathbb{R}^m$, ta cần không gian cột của \mathbf{A} phải bao phủ toàn bộ \mathbb{R}^m . Nếu bất kỳ điểm nào trong \mathbb{R}^m bị loại khỏi không gian cột, điểm đó là một giá trị \mathbf{b} tiềm năng mà không có nghiệm. Việc yêu cầu rằng không gian cột của \mathbf{A} phải bao phủ toàn bộ \mathbb{R}^m ngay lập tức ngụ ý rằng \mathbf{A} phải có ít nhất m cột, tức là $n \geq m$. Nếu không, số chiều của không gian cột sẽ nhỏ hơn m . Ví dụ, xét một ma trận cỡ 3×2 . Mục tiêu \mathbf{b} là 3 chiều, nhưng \mathbf{x} chỉ có 2 chiều, vì vậy việc thay đổi giá trị của \mathbf{x} chỉ cho phép ta truy ra một mặt

phẳng 2 chiều trong \mathbb{R}^3 . Phương trình có nghiệm nếu và chỉ nếu \mathbf{b} nằm trên mặt phẳng đó.

Điều kiện $n \geq m$ chỉ là điều kiện cần để mọi điểm đều có nghiệm. Nó không phải là điều kiện đủ, vì có thể một số cột là dư thừa. Ví dụ, xét một ma trận cỡ 2×2 có hai cột giống nhau. Ma trận này có không gian cột giống với không gian cột của ma trận 2×1 chỉ chứa một cột giống nhau ở trên. Nói cách khác, không gian cột vẫn chỉ là một đường thẳng và không bao phủ toàn bộ \mathbb{R}^2 , mặc dù có hai cột.

Sự dư thừa kiểu này được gọi là **phụ thuộc tuyến tính**. Một tập hợp các vectơ được gọi là **độc lập tuyến tính** nếu không có vectơ nào trong tập là tổ hợp tuyến tính của các vectơ còn lại. Nếu ta thêm một vectơ vào một tập hợp mà vectơ đó là tổ hợp tuyến tính của các vectơ khác trong tập, thì vectơ mới này không thêm điểm nào vào bao tuyến tính của tập đó. Điều này có nghĩa là để không gian cột của ma trận bao phủ toàn bộ \mathbb{R}^m , ma trận phải chứa ít nhất một tập hợp gồm m cột độc lập tuyến tính. Điều kiện này vừa là điều kiện cần vừa là điều kiện đủ để phương trình (2.10) có nghiệm với mọi giá trị của \mathbf{b} . Lưu ý rằng yêu cầu là tập hợp phải có đúng m cột độc lập tuyến tính, chứ không phải ít nhất m . Không có tập hợp các vectơ m chiều nào có thể có hơn m cột độc lập tuyến tính, nhưng một ma trận có nhiều hơn m cột có thể chứa nhiều hơn một tập hợp như vậy.

Để ma trận có nghịch đảo, ta cần đảm bảo rằng phương trình (2.10) có *nhiều nhất* một nghiệm cho mỗi giá trị của \mathbf{b} . Để làm được điều đó, ta cần đảm bảo rằng ma trận có nhiều nhất m cột. Nếu không, sẽ có nhiều hơn một cách tham số hóa cho mỗi nghiệm.

Điều này có nghĩa là ma trận phải là **ma trận vuông**, tức là $m = n$ và tất cả các cột phải tuyến tính độc lập. Một ma trận vuông có các cột phụ thuộc tuyến tính được gọi là **ma trận suy biến**.

Nếu \mathbf{A} không phải là ma trận vuông hoặc là ma trận vuông nhưng suy biến, vẫn có thể giải được phương trình. Tuy nhiên, ta không thể sử dụng phương pháp nghịch đảo ma trận để tìm nghiệm.

Đến lúc này, ta đã thảo luận về nghịch đảo ma trận khi nhân từ bên trái. Cũng có thể định nghĩa nghịch đảo được nhân từ bên phải:

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{I}. \quad (2.19)$$

Đối với các ma trận vuông, nghịch đảo trái và nghịch đảo phải là bằng nhau.

2.5 Chuẩn

Đôi khi ta cần đo chiều dài của một vectơ. Trong học máy, ta thường đo chiều dài của các vectơ bằng một hàm gọi là **chuẩn**. Chuẩn L^p được cho bởi

$$\|\mathbf{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}} \quad (2.20)$$

với $p \in \mathbb{R}, p \geq 1$.

Các chuẩn, bao gồm chuẩn L^p , là những hàm gán tương ứng mỗi vectơ với một số không âm. Về mặt trực quan, chuẩn của một vectơ \mathbf{x} đo khoảng cách từ gốc tọa độ đến điểm \mathbf{x} . Chặt chẽ hơn, một chuẩn là bất kỳ hàm f nào thỏa mãn các tính chất sau:

- $f(\mathbf{x}) = 0 \Rightarrow \mathbf{x} = \mathbf{0}$
- $f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y})$ (**bất đẳng thức tam giác**)
- $\forall \alpha \in \mathbb{R}, \quad f(\alpha \mathbf{x}) = |\alpha| f(\mathbf{x})$

Chuẩn L^2 , với $p = 2$, được gọi là **chuẩn Euclid**. Nó đơn giản là khoảng cách Euclid từ gốc tọa độ đến điểm được xác định bởi \mathbf{x} . Chuẩn L^2 được sử dụng thường xuyên trong học máy đến mức thường được ký hiệu đơn giản là $\|\mathbf{x}\|$, với chỉ số 2 bị lược bỏ. Một đại lượng cũng phổ biến khi đánh giá chiều dài của một vectơ là chuẩn L^2 bình phương, có thể được tính đơn giản là $\mathbf{x}^T \mathbf{x}$.

Chuẩn L^2 bình phương thuận tiện hơn trong việc xử lý về mặt toán học và tính toán so với chính chuẩn L^2 . Ví dụ, các đạo hàm của chuẩn L^2 bình phương theo từng phần tử của \mathbf{x} chỉ phụ thuộc vào phần tử tương ứng của \mathbf{x} , trong khi tất cả các đạo hàm của chuẩn L^2 phụ thuộc vào toàn bộ vectơ. Trong nhiều tình huống, chuẩn L^2 bình phương có thể không được ưa chuộng vì nó tăng rất chậm gần gốc tọa độ. Trong một số ứng dụng học máy, việc phân biệt giữa các phần tử có giá trị chính xác bằng không và các phần tử nhỏ nhưng không bằng không là rất quan trọng. Trong những trường hợp này, ta sử dụng một hàm tăng với cùng tốc độ ở mọi vị trí, nhưng vẫn giữ được tính đơn giản về mặt toán học: chuẩn L^1 . Chuẩn L^1 được xác định đơn giản như sau:

$$\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|. \quad (2.21)$$

Chuẩn L^1 thường được sử dụng trong học máy khi sự khác biệt giữa các phần tử bằng 0 và khác 0 là rất quan trọng. Mỗi khi một phần tử của \mathbf{x} dịch chuyển khỏi 0 một lượng ε , chuẩn L^1 sẽ tăng thêm ε .

Trong lĩnh vực học máy, đôi khi ta đo độ lớn của một vectơ bằng cách đếm số phần tử khác 0 của nó. Một số tác giả gọi hàm này là “chuẩn L^0 ”, nhưng đây là thuật ngữ không chính xác. Số phần tử khác 0 trong một vectơ không phải là một chuẩn, vì việc nhân vectơ với một hệ số α không làm thay đổi số phần tử khác 0. Chuẩn L^1 thường được sử dụng như một sự thay thế cho số phần tử khác 0.

Một chuẩn khác thường xuất hiện trong học máy là chuẩn L^∞ , còn được gọi là **chuẩn max**, xác định bởi giá trị tuyệt đối của phần tử có độ lớn lớn nhất trong vectơ:

$$\|\mathbf{x}\|_\infty = \max_i |x_i|. \quad (2.22)$$

Đôi khi ta cũng muốn đo độ lớn của một ma trận. Trong bối cảnh học sâu, cách phổ biến nhất để làm điều này là sử dụng **chuẩn Frobenius**, một chuẩn ít được biết đến:

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i,j=1}^n A_{ij}^2}, \quad (2.23)$$

tương tự như chuẩn L_2 của một vectơ.

Tích vô hướng của hai vectơ có thể được viết lại theo các chuẩn. Cụ thể:

$$\mathbf{x}^T \mathbf{y} = \|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \cos \theta \quad (2.24)$$

trong đó θ là góc giữa \mathbf{x} và \mathbf{y} .

2.6 Các loại ma trận và vectơ đặc biệt

Một số loại ma trận và vectơ đặc biệt rất hữu ích.

Ma trận đường chéo bao gồm hầu hết các phần tử bằng 0 và chỉ có các phần tử khác 0 trên đường chéo chính. Ma trận \mathbf{D} là ma trận đường chéo khi và chỉ khi $D_{i,j} = 0$ cho mọi $i \neq j$. Ta đã thấy một ví dụ về ma trận chéo: ma trận đơn vị, trong đó tất cả các phần tử trên đường chéo chính đều bằng 1. Ta viết $\text{diag}(\mathbf{v})$ để biểu diễn ma trận vuông chéo mà các phần tử trên đường chéo được xác định bởi các phần tử của vectơ \mathbf{v} . Ma trận đường chéo đáng chú ý vì việc nhân với ma trận đường chéo rất hiệu quả về tính toán. Để tính $\text{diag}(\mathbf{v}) \mathbf{x}$, ta chỉ cần nhân từng

phần tử x_i với v_i . Nói cách khác, $\text{diag}(\mathbf{v})\mathbf{x} = \mathbf{v} \odot \mathbf{x}$. Việc tìm nghịch đảo của ma trận vuông đường chéo cũng rất dễ dàng. Ma trận nghịch đảo chỉ tồn tại khi tất cả các phần tử trên đường chéo chính đều khác không, và trong trường hợp đó, $\text{diag}(\mathbf{v})^{-1} = \text{diag}\left(\left[\frac{1}{v_1}, \dots, \frac{1}{v_n}\right]^T\right)$. Trong nhiều trường hợp, ta có thể xây dựng một số thuật toán học máy rất tổng quát cho các ma trận bất kỳ, từ một thuật toán tương tự nhưng đơn giản hơn bằng cách giới hạn một số ma trận là ma trận đường chéo.

Không phải tất cả các ma trận đường chéo đều phải là ma trận vuông. Có thể tạo ra một ma trận đường chéo hình chữ nhật. Ma trận đường chéo không vuông không có ma trận nghịch đảo, nhưng vẫn có thể nhân chúng với vectơ và ma trận. Đối với một ma trận chéo không vuông \mathbf{D} , tích $\mathbf{D}\mathbf{x}$ sẽ bao gồm việc co giãn từng phần tử của \mathbf{x} , và, hoặc là nối thêm một số số không vào kết quả nếu \mathbf{D} nhiều hàng hơn cột, hoặc loại bỏ một số phần tử cuối cùng của vectơ nếu \mathbf{D} ít hàng hơn cột.

Ma trận đối xứng là ma trận bằng ma trận chuyển vị của chính nó:

$$\mathbf{A} = \mathbf{A}^T. \quad (2.25)$$

Ma trận đối xứng thường xuất hiện khi các phần tử được tạo ra bởi một hàm có hai tham số mà không phụ thuộc vào thứ tự của các tham số. Ví dụ, nếu \mathbf{A} là một ma trận các khoảng cách, với A_{ij} biểu diễn khoảng cách từ điểm i đến điểm j , thì $A_{ij} = A_{ji}$ vì các hàm khoảng cách là đối xứng.

Vectơ đơn vị là một vectơ có **chuẩn bằng 1**:

$$\|\mathbf{x}\|_2 = 1. \quad (2.26)$$

Vectơ \mathbf{x} và \mathbf{y} gọi là **trực giao (vuông góc)** với nhau nếu $\mathbf{x}^T \mathbf{y} = 0$. Nếu cả hai vectơ đều có chuẩn khác không, điều này có nghĩa là chúng tạo với nhau một góc 90 độ. Trong không gian \mathbb{R}^n , tối đa có thể có n vectơ trực giao lẫn nhau với chuẩn khác 0. Nếu các vectơ không chỉ trực giao mà còn có chuẩn đơn vị, ta gọi chúng là các vectơ **trực chuẩn**.

Ma trận trực giao là ma trận vuông có các hàng trực chuẩn lẫn nhau và các cột cũng trực chuẩn lẫn nhau.

$$\mathbf{A}^T \mathbf{A} = \mathbf{A} \mathbf{A}^T = \mathbf{I}. \quad (2.27)$$

Điều này dẫn đến

$$\mathbf{A}^{-1} = \mathbf{A}^T, \quad (2.28)$$

vì vậy các ma trận trực giao thu hút sự quan tâm do nghịch đảo của chúng rất dễ tính toán. Cần chú ý kỹ đến định nghĩa của ma trận trực giao. Khác với tên gọi, các hàng của chúng không chỉ trực giao mà còn trực chuẩn. Không có thuật ngữ đặc biệt nào cho ma trận mà các hàng hoặc cột của nó chỉ trực giao nhưng không trực chuẩn.

2.7 Phân tích giá trị riêng

Nhiều đối tượng toán học có thể được hiểu rõ hơn bằng cách phân tách chúng thành các phần cấu thành, hoặc tìm ra một số tính chất điển hình của chúng, không phụ thuộc vào cách biểu diễn chúng.

Ví dụ, các số nguyên có thể được phân tích thành các thừa số nguyên tố. Cách ta biểu diễn số 12 sẽ thay đổi tùy vào việc ta viết nó ở hệ thập phân hay nhị phân, nhưng biểu diễn $12 = 2 \times 2 \times 3$ luôn đúng. Từ cách biểu diễn này, ta có thể suy ra các tính chất hữu ích, chẳng hạn như 12 không chia hết cho 5, hoặc bất kỳ bội số nào của 12 cũng sẽ chia hết cho 3.

Giống như việc có thể khám phá bản chất thực sự của một số nguyên bằng cách phân tích nó thành các thừa số nguyên tố, ta cũng có thể phân tích ma trận theo cách cho thấy các thông tin về thuộc tính cơ bản của chúng, là điều không hiển nhiên do cách biểu diễn ma trận như một mảng các phần tử.

Một trong những loại phân tích ma trận được sử dụng rộng rãi nhất được gọi là **phân tích giá trị riêng**, trong đó ta phân tích ma trận thành một tập hợp các vectơ riêng và giá trị riêng.

Một **vectơ riêng** của ma trận vuông \mathbf{A} là một vectơ khác không \mathbf{v} sao cho phép nhân với \mathbf{A} chỉ thay đổi một tỷ lệ của \mathbf{v} :

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}. \quad (2.29)$$

Hệ số vô hướng λ được gọi là **giá trị riêng** tương ứng với vectơ riêng này. (Ta cũng có thể tìm **vectơ riêng trái** sao cho $\mathbf{v}^T \mathbf{A} = \lambda \mathbf{v}^T$, nhưng thường ta quan tâm đến vectơ riêng phải).

Nếu \mathbf{v} là một vectơ riêng của ma trận \mathbf{A} , thì bất kỳ vectơ nào tỷ lệ với nó, tức là $s\mathbf{v}$ với $s \in \mathbb{R}$, $s \neq 0$ cũng là một vectơ riêng của \mathbf{A} . Hơn nữa, $s\mathbf{v}$ vẫn có cùng giá trị riêng. Vì lý do này, ta thường chỉ tìm các vectơ riêng đơn vị.

Giả sử ma trận \mathbf{A} có n vectơ riêng độc lập tuyến tính, $\{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)}\}$, với các giá trị riêng tương ứng $\{\lambda_1, \dots, \lambda_n\}$. Ta có thể ghép tất cả các vectơ riêng lại để tạo

thành một ma trận \mathbf{V} , với mỗi vectơ riêng là một cột: $\mathbf{V} = [\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)}]$. Tương tự, ta có thể ghép các giá trị riêng lại để tạo thành một vectơ $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_n]^T$. Phân tích giá trị riêng của \mathbf{A} sau đó được cho bởi

$$\mathbf{A} = \mathbf{V} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^{-1}. \quad (2.30)$$

Ta đã thấy rằng việc xây dựng các ma trận với các giá trị riêng và vectơ riêng cụ thể cho phép ta kéo giãn không gian theo các hướng mong muốn. Tuy nhiên, ta thường muốn phân tích ma trận thành các giá trị riêng và vectơ riêng của nó. Việc làm này có thể giúp ta phân tích một số thuộc tính của ma trận, giống như việc phân tích một số nguyên thành các thừa số nguyên tố có thể giúp ta hiểu được đặc điểm của số nguyên đó.

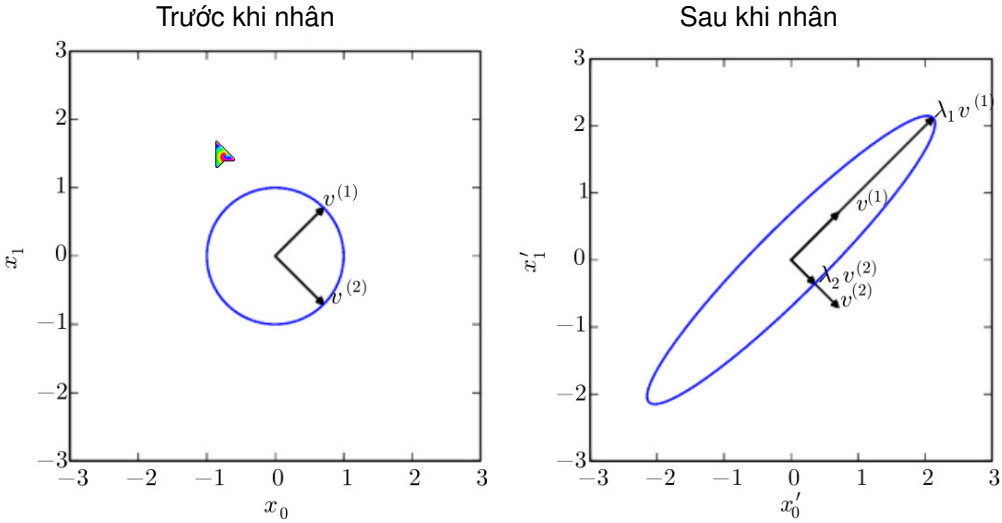
Không phải mọi ma trận đều có thể phân tích thành các giá trị riêng và vectơ riêng. Trong một số trường hợp, phép phân tích tồn tại, nhưng có thể liên quan đến số phức thay vì số thực. May mắn là, trong cuốn sách này, ta thường chỉ cần phân tích một lớp ma trận cụ thể có phân tích đơn giản. Cụ thể, mọi ma trận thực đối xứng đều có thể phân tích thành biểu thức chỉ sử dụng các vectơ riêng và giá trị riêng thực:

$$\mathbf{A} = \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^T, \quad (2.31)$$

trong đó \mathbf{Q} là ma trận trực giao bao gồm các vectơ riêng của \mathbf{A} , và $\boldsymbol{\Lambda}$ là ma trận đường chéo. Giá trị riêng Λ_{ii} tương ứng với vectơ riêng ở cột i của \mathbf{Q} , ký hiệu là $\mathbf{Q}_{:,i}$. Vì \mathbf{Q} là ma trận trực giao, ta có thể xem \mathbf{A} như việc kéo giãn không gian với mỗi hướng $\mathbf{v}^{(i)}$ theo tỷ lệ λ_i . Hình 2.3 là một ví dụ về mô tả trên.

Mặc dù bất kỳ ma trận thực đối xứng thực \mathbf{A} nào cũng đều có phân tích giá trị riêng, nhưng phân tích giá trị riêng có thể không duy nhất. Nếu hai hoặc nhiều vectơ riêng có cùng giá trị riêng, thì bất kỳ tập hợp các vectơ trực giao nào nằm trong không gian con sinh bởi chúng cũng là các vectơ riêng với giá trị riêng đó, và ta có thể chọn một ma trận \mathbf{Q} sử dụng các vectơ riêng đó. Theo quy ước, ta thường sắp xếp các phần tử của $\boldsymbol{\Lambda}$ theo thứ tự giảm dần. Theo quy ước này, phân tích giá trị riêng là duy nhất khi tất cả các giá trị riêng đều là duy nhất.

Phân tích giá trị riêng của một ma trận cho ta biết nhiều thông tin hữu ích về ma trận đó. Ma trận là suy biến nếu và chỉ nếu ma trận đó có giá trị riêng nào đó bằng 0. Phân tích giá trị riêng của một ma trận thực đối xứng cũng có thể được sử dụng để tối ưu hóa các dạng toàn phương, tức là các biểu thức bậc hai có dạng $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$ với điều kiện $\|\mathbf{x}\|_2 = 1$. Mỗi khi \mathbf{x} bằng một vectơ riêng của \mathbf{A} , f sẽ có giá trị bằng giá trị riêng tương ứng. Giá trị lớn nhất của f trong miền ràng buộc



Hình 2.3: Ví dụ về tác động của vectơ riêng và giá trị riêng. Ở đây, ta có ma trận \mathbf{A} với hai vectơ riêng trực chuẩn, $\mathbf{v}^{(1)}$ với giá trị riêng λ_1 và $\mathbf{v}^{(2)}$ với giá trị riêng λ_2 . (Bên trái) Ta vẽ tập hợp tất cả các vectơ đơn vị $\mathbf{u} \in \mathbb{R}^2$ dưới dạng một đường tròn đơn vị. (Bên phải) Ta vẽ tập hợp tất cả các điểm $\mathbf{A}\mathbf{u}$. Bằng cách quan sát cách \mathbf{A} làm biến dạng đường tròn đơn vị, ta có thể thấy rằng nó kéo giãn không gian theo hướng $\mathbf{v}^{(i)}$ với tỷ lệ λ_i .

là giá trị riêng lớn nhất và giá trị nhỏ nhất trong miền ràng buộc là giá trị riêng nhỏ nhất.

Ma trận có tất cả các giá trị riêng dương được gọi là ma trận xác định dương. Ma trận có tất cả các giá trị riêng dương hoặc bằng 0 được gọi là ma trận nửa xác định dương. Tương tự, nếu tất cả các giá trị riêng âm, ma trận được gọi là ma trận xác định âm, và nếu tất cả các giá trị riêng âm hoặc bằng 0, ma trận là nửa xác định âm. Các ma trận nửa xác định dương rất thú vị vì chúng đảm bảo rằng $\forall \mathbf{x}, \mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$. Các ma trận xác định dương còn đảm bảo thêm rằng $\mathbf{x}^T \mathbf{A} \mathbf{x} = 0 \Rightarrow \mathbf{x} = \mathbf{0}$.

2.8 Phân tích giá trị kỳ dị

Trong Mục 2.7, ta đã thấy cách phân tích một ma trận thành các vectơ riêng và giá trị riêng. **Phân tích giá trị kỳ dị** (SVD) cung cấp một cách khác để phân tích ma trận thành các vectơ suy biến và giá trị suy biến. Phân tích giá trị kỳ dị cho phép ta khám phá một số thông tin tương tự như phép phân tích giá trị riêng. Tuy nhiên, phân tích giá trị kỳ dị có thể áp dụng tổng quát hơn. Mọi ma trận thực đều có

phân tích giá trị kỳ dị, nhưng điều này không đúng đối với phân tích giá trị riêng. Ví dụ, nếu một ma trận không phải là ma trận vuông, thì phép phân tích giá trị riêng không được định nghĩa và thay vào đó ta phải sử dụng phép phân tích giá trị kỳ dị.

Nhớ lại rằng phép phân tích giá trị riêng liên quan đến việc phân tích một ma trận \mathbf{A} để tìm ma trận \mathbf{V} gồm các vectơ riêng và một vectơ $\boldsymbol{\lambda}$ gồm các giá trị riêng sao cho có thể viết lại \mathbf{A} dưới dạng

$$\mathbf{A} = \mathbf{V} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^{-1}. \quad (\text{xem (2.30)})$$

Phân tích giá trị kỳ dị cũng tương tự, ngoại trừ lần này ta sẽ viết \mathbf{A} dưới dạng tích của ba ma trận:

$$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^T \quad (2.32)$$

Giả sử \mathbf{A} là ma trận cỡ $m \times n$. Khi đó \mathbf{U} được định nghĩa là ma trận cỡ $m \times m$, \mathbf{D} là ma trận cỡ $m \times n$, và \mathbf{V} là ma trận cỡ $n \times n$.

Mỗi ma trận trong số này được định nghĩa có một cấu trúc đặc biệt. Các ma trận \mathbf{U} và \mathbf{V} đều là ma trận trực giao. Ma trận \mathbf{D} là ma trận đường chéo. Lưu ý rằng \mathbf{D} không nhất thiết phải là ma trận vuông. Các phần tử dọc theo đường chéo của \mathbf{D} được gọi là các **giá trị kỳ dị** của ma trận \mathbf{A} . Các cột của \mathbf{U} được gọi là các **vectơ kỳ dị trái**. Các cột của \mathbf{V} được gọi là các **vectơ kỳ dị phải**.

Ta thực sự có thể diễn giải phân tích giá trị kỳ dị của ma trận \mathbf{A} dưới dạng phân tích giá trị riêng của các ma trận liên quan đến \mathbf{A} . Các vectơ kỳ dị trái của \mathbf{A} là các vectơ riêng của $\mathbf{A} \mathbf{A}^T$. Các vectơ kỳ dị phải của \mathbf{A} là các vectơ riêng của $\mathbf{A}^T \mathbf{A}$. Các giá trị kỳ dị khác 0 của \mathbf{A} là căn bậc hai của các giá trị riêng của $\mathbf{A}^T \mathbf{A}$. Điều này cũng đúng với $\mathbf{A} \mathbf{A}^T$.

Có lẽ tính năng hữu ích nhất của phân tích giá trị kỳ dị là ta có thể sử dụng nó để tổng quát hóa một phần phép nghịch đảo ma trận cho các ma trận không vuông, như sẽ thấy trong phần tiếp theo.

2.9 Ma trận giả nghịch đảo Moore – Penrose

Phép nghịch đảo ma trận không được định nghĩa cho các ma trận không vuông. Giả sử ta muốn xây dựng một ma trận nghịch đảo trái \mathbf{B} của ma trận \mathbf{A} , để có thể giải phương trình tuyến tính

$$\mathbf{A} \mathbf{x} = \mathbf{y} \quad (2.33)$$

bằng cách nhân bên trái mỗi vế để thu được

$$\mathbf{x} = \mathbf{B} \mathbf{y}. \quad (2.34)$$

Tùy thuộc vào cấu trúc của bài toán, có thể không thể thiết kế được một tương ứng duy nhất biến \mathbf{A} thành \mathbf{B} .

Nếu ma trận \mathbf{A} có nhiều hàng hơn cột, thì phương trình này có thể không có nghiệm. Nếu \mathbf{A} ít hàng hơn cột, thì có thể có nhiều nghiệm khác nhau. **Ma trận giả nghịch đảo Moore–Penrose** cho phép ta xử lý các trường hợp này. Giả nghịch đảo của \mathbf{A} được định nghĩa là ma trận

$$\mathbf{A}^+ = \lim_{\alpha \downarrow 0} (\mathbf{A}^T \mathbf{A} + \alpha \mathbf{I})^{-1} \mathbf{A}^T. \quad (2.35)$$

Các thuật toán thực tế để tính giả nghịch đảo không dựa trên định nghĩa này, mà dựa trên công thức

$$\mathbf{A}^+ = \mathbf{V} \mathbf{D}^+ \mathbf{U}^T, \quad (2.36)$$

trong đó \mathbf{U} , \mathbf{D} và \mathbf{V} là phân tích giá trị kỳ dị của \mathbf{A} , và giả nghịch đảo \mathbf{D}^+ của ma trận đường chéo \mathbf{D} được tính bằng cách lấy nghịch đảo của các phần tử khác không của nó, sau đó lấy chuyển vị của ma trận kết quả.

Khi ma trận \mathbf{A} có nhiều cột hơn hàng, việc giải một phương trình tuyến tính bằng ma trận giả nghịch đảo sẽ cung cấp một trong nhiều nghiệm có thể. Cụ thể, nó trả về nghiệm $\mathbf{x} = \mathbf{A}^+ \mathbf{y}$ với chuẩn Euclid tối thiểu $\|\mathbf{x}\|_2$ trong số tất cả các nghiệm có thể.

Khi ma trận \mathbf{A} có nhiều hàng hơn cột, có thể phương trình không có nghiệm. Trong trường hợp này, việc sử dụng giả nghịch đảo sẽ cho ta giá trị \mathbf{x} sao cho $\mathbf{A}\mathbf{x}$ gần với \mathbf{y} nhất theo chuẩn Euclid $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2$.

2.10 Vết của ma trận

Vết của ma trận là tổng tất cả các phần tử trên đường chéo chính của ma trận:

$$\text{Tr}(\mathbf{A}) = \sum_i A_{ii}. \quad (2.37)$$

Toán tử vết hữu ích vì nhiều lý do. Một số phép toán khó biểu diễn mà không sử dụng ký hiệu tổng có thể được biểu diễn bằng cách sử dụng tích ma trận và toán tử vết. Ví dụ, toán tử vết cung cấp một cách viết khác cho chuẩn Frobenius của ma trận:

$$\|\mathbf{A}\|_F = \sqrt{\text{Tr}(\mathbf{A}\mathbf{A}^T)}. \quad (2.38)$$

Việc viết một biểu thức dưới dạng toán tử vết mở ra cơ hội để biến đổi biểu thức bằng cách sử dụng nhiều đồng nhất thức phổ biến. Ví dụ, toán tử vết bất biến

đối với phép chuyển vị ma trận:

$$\text{Tr}(\mathbf{A}) = \text{Tr}(\mathbf{A}^T). \quad (2.39)$$

Vết của ma trận vuông được tạo thành từ tích nhiều ma trận cũng bất biến khi di chuyển ma trận nhân tử cuối cùng lên vị trí đầu tiên, nếu cỡ của các ma trận là tương thích, tức là phép nhân ma trận khi viết ra có thể thực hiện được.

$$\text{Tr}(\mathbf{ABC}) = \text{Tr}(\mathbf{CAB}) = \text{Tr}(\mathbf{BCA}) \quad (2.40)$$

hay tổng quát hơn:

$$\text{Tr}\left(\prod_{i=1}^n \mathbf{F}^{(i)}\right) = \text{Tr}\left(\mathbf{F}^{(n)} \prod_{i=1}^{n-1} \mathbf{F}^{(i)}\right). \quad (2.41)$$

Tính bất biến đối với hoán vị vòng này vẫn đúng ngay cả khi các tích ma trận có cỡ khác nhau. Ví dụ, với $\mathbf{A} \in \mathbb{R}^{m \times n}$ và $\mathbf{B} \in \mathbb{R}^{n \times m}$, ta có

$$\text{Tr}(\mathbf{AB}) = \text{Tr}(\mathbf{BA}) \quad (2.42)$$

mặc dù $\mathbf{AB} \in \mathbb{R}^{m \times m}$ và $\mathbf{BA} \in \mathbb{R}^{n \times n}$.

Một kết quả khác cần nhớ là một số vô hướng chính là vết của nó: $a = \text{Tr}(a)$.

2.11 Định thức

Định thức của một ma trận vuông, ký hiệu là $\det(\mathbf{A})$, là một hàm ánh xạ mỗi ma trận thành một số thực. Định thức bằng tích của tất cả các giá trị riêng của ma trận. Giá trị tuyệt đối của định thức có thể được coi như một thước đo cho việc phép nhân với ma trận đó mở rộng hay thu hẹp không gian bao nhiêu. Nếu định thức bằng 0, thì không gian bị thu hẹp hoàn toàn dọc theo ít nhất một chiều, khiến nó mất toàn bộ thể tích. Nếu định thức bằng 1, thì phép biến đổi bảo toàn thể tích.

2.12 Ví dụ: phân tích thành phần chính

Một thuật toán học máy đơn giản, **phân tích thành phần chính** (PCA), có thể được suy ra chỉ bằng kiến thức cơ bản về đại số tuyến tính.

Cho tập gồm m điểm $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ trong không gian \mathbb{R}^n . Giả sử ta muốn áp dụng nén mất dữ liệu cho các điểm này. Nén mất dữ liệu có nghĩa là lưu trữ các điểm theo cách yêu cầu ít bộ nhớ hơn nhưng có thể mất một phần độ chính xác. Ta muốn mất càng ít độ chính xác càng tốt.

Một cách để mã hóa các điểm này là biểu diễn chúng dưới dạng phiên bản có số chiều thấp hơn. Đối với mỗi điểm $\mathbf{x}^{(i)} \in \mathbb{R}^n$, ta sẽ tìm một vectơ mã hóa tương ứng $\mathbf{c}^{(i)} \in \mathbb{R}^\ell$. Nếu $\ell < n$, sẽ cần ít bộ nhớ hơn để lưu trữ các điểm mã hóa so với dữ liệu gốc. Ta sẽ cần tìm một hàm mã hóa tạo ra mã cho đầu vào, $\mathbf{f}(\mathbf{x}) = \mathbf{c}$, và một hàm giải mã để tái tạo đầu vào từ mã hóa của nó, $\mathbf{x} \approx \mathbf{g}(\mathbf{f}(\mathbf{x}))$.

Phân tích thành phần chính được xác định bởi cách ta chọn hàm giải mã. Cụ thể, một giải mã rất đơn giản, sử dụng phép nhân ma trận để ánh xạ mã trở lại không gian \mathbb{R}^n . Giả sử $\mathbf{g}(\mathbf{c}) = \mathbf{D}\mathbf{c}$, trong đó $\mathbf{D} \in \mathbb{R}^{n \times \ell}$ là ma trận xác định việc giải mã.

Việc tính toán mã tối ưu cho bộ giải mã này có thể là một vấn đề phức tạp. Để giữ cho việc mã hóa dễ dàng, phân tích thành phần chính yêu cầu các cột của \mathbf{D} phải trực giao với nhau. (Lưu ý rằng \mathbf{D} vẫn chưa phải là “ma trận trực giao” trừ khi $\ell = n$).

Với bài toán như mô tả ở trên, có nhiều nghiệm khả thi, vì ta có thể tăng tỷ lệ của $\mathbf{D}_{:,i}$ nếu giảm tỷ lệ tương ứng của c_i tại tất cả các điểm. Để đảm bảo bài toán có nghiệm duy nhất, ta ràng buộc tất cả các cột của \mathbf{D} phải có chuẩn đơn vị.

Để biến ý tưởng cơ bản này thành một thuật toán có thể triển khai, điều đầu tiên ta cần làm là tìm cách tạo ra điểm mã tối ưu \mathbf{c}^* cho mỗi điểm đầu vào \mathbf{x} . Một cách để làm điều này là cực tiểu hóa khoảng cách giữa điểm đầu vào \mathbf{x} và tái tạo của nó, $\mathbf{g}(\mathbf{c}^*)$. Ta có thể đo khoảng cách này bằng một chuẩn. Trong thuật toán phân tích thành phần chính, ta sử dụng chuẩn L^2 :

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{x} - \mathbf{g}(\mathbf{c})\|_2. \quad (2.43)$$

Ta có thể chuyển sang sử dụng chuẩn L^2 bình phương thay vì chuẩn L^2 , vì cả hai đều được cực tiểu hóa bởi cùng một giá trị của \mathbf{c} . Cả hai đều đạt giá trị nhỏ nhất với cùng một giá trị \mathbf{c} vì chuẩn L^2 không âm và phép bình phương là phép toán đơn điệu tăng đối với các đối số không âm.

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{x} - \mathbf{g}(\mathbf{c})\|_2^2. \quad (2.44)$$

Hàm đang được cực tiểu hóa đã được đơn giản hóa thành

$$[\mathbf{x} - \mathbf{g}(\mathbf{c})]^T [\mathbf{x} - \mathbf{g}(\mathbf{c})] \quad (2.45)$$

(theo định nghĩa của chuẩn L^2 ở công thức (2.20))

$$= \mathbf{x}^T \mathbf{x} - \mathbf{x}^T \mathbf{g}(\mathbf{c}) - \mathbf{g}(\mathbf{c})^T \mathbf{x} + \mathbf{g}(\mathbf{c})^T \mathbf{g}(\mathbf{c}) \quad (2.46)$$

(theo tính chất phân phối)

$$= \mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T \mathbf{g}(\mathbf{c}) + \mathbf{g}(\mathbf{c})^T \mathbf{g}(\mathbf{c}) \quad (2.47)$$

(vì đại lượng vô hướng $\mathbf{x}^T \mathbf{g}(\mathbf{c})$ bằng chuyển vị của chính nó).

Bây giờ ta có thể một lần nữa thay đổi hàm đang được cực tiểu hóa, bằng cách bỏ đi hạng tử đầu tiên, vì hạng tử này không phụ thuộc vào \mathbf{c} :

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} -2\mathbf{x}^T \mathbf{g}(\mathbf{c}) + \mathbf{g}(\mathbf{c})^T \mathbf{g}(\mathbf{c}). \quad (2.48)$$

Để tiện xa hơn, ta cần thay định nghĩa của $\mathbf{g}(\mathbf{c})$ vào:

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} -2\mathbf{x}^T \mathbf{D}\mathbf{c} + \mathbf{c}^T \mathbf{D}^T \mathbf{D}\mathbf{c} \quad (2.49)$$

$$= \arg \min_{\mathbf{c}} -2\mathbf{x}^T \mathbf{D}\mathbf{c} + \mathbf{c}^T \mathbf{I}_{\ell} \mathbf{c} \quad (2.50)$$

(do tính trực giao và chuẩn đơn vị ràng buộc trên \mathbf{D})

$$= \arg \min_{\mathbf{c}} -2\mathbf{x}^T \mathbf{D}\mathbf{c} + \mathbf{c}^T \mathbf{c} \quad (2.51)$$

Ta có thể giải bài toán tối ưu này bằng cách sử dụng phép tính vector (xem [Mục 4.1](#) nếu bạn chưa biết cách thực hiện điều này).

$$\nabla_{\mathbf{c}} (-2\mathbf{x}^T \mathbf{D}\mathbf{c} + \mathbf{c}^T \mathbf{c}) = \mathbf{0} \quad (2.52)$$

$$-2\mathbf{D}^T \mathbf{x} + 2\mathbf{c} = \mathbf{0} \quad (2.53)$$

$$\mathbf{c} = \mathbf{D}^T \mathbf{x}. \quad (2.54)$$

Điều này làm cho thuật toán trở nên hiệu quả: ta có thể mã hóa \mathbf{x} một cách tối ưu chỉ bằng một phép toán ma trận-véc tơ. Để mã hóa một véc tơ, ta dùng hàm mã hóa

$$\mathbf{f}(\mathbf{x}) = \mathbf{D}^T \mathbf{x}. \quad (2.55)$$

Sử dụng thêm phép nhân ma trận nữa, ta cũng có thể định nghĩa phép tái tạo của phương pháp phân tích thành phần chính:

$$\mathbf{r}(\mathbf{x}) = \mathbf{g}(\mathbf{f}(\mathbf{x})) = \mathbf{D}\mathbf{D}^T \mathbf{x}. \quad (2.56)$$

Tiếp theo, ta cần chọn ma trận mã hóa \mathbf{D} . Để làm điều này, ta quay lại ý tưởng về việc cực tiểu khoảng cách L^2 giữa các đầu vào và các điểm tái tạo. Vì ta sẽ sử dụng cùng một ma trận \mathbf{D} để giải mã tất cả các điểm, nên không thể xem xét các

điểm một cách riêng lẻ nữa. Thay vào đó, ta phải giảm thiểu chuẩn Frobenius của ma trận sai số được tính trên tất cả các chiều và tất cả các điểm.

$$\mathbf{D}^* = \arg \min_{\mathbf{D}} \sqrt{\sum_{i=1}^m \|\mathbf{x}^{(i)} - \mathbf{r}(\mathbf{x}^{(i)})\|_2^2} = \arg \min_{\mathbf{D}} \sqrt{\sum_{i,j} \left(x_j^{(i)} - \mathbf{r}(\mathbf{x}^{(i)})_j\right)^2} \quad (2.57)$$

sao cho $\mathbf{D}^T \mathbf{D} = \mathbf{I}_\ell$.

Để xây dựng thuật toán tìm \mathbf{D}^* , ta sẽ bắt đầu bằng cách xét trường hợp $\ell = 1$. Trong trường hợp này, \mathbf{D} chỉ là một vectơ, ký hiệu \mathbf{d} . Thay phương trình (2.56) vào phương trình (2.57) và xét đơn giản hóa \mathbf{D} thành \mathbf{d} , bài toán được rút gọn thành

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \sum_{i=1}^m \|\mathbf{x}^{(i)} - \mathbf{d}\mathbf{d}^T \mathbf{x}^{(i)}\|_2^2 \quad \text{sao cho} \quad \|\mathbf{d}\|_2 = 1. \quad (2.58)$$

Vì chuẩn L^2 của vectơ chính là chuẩn Frobenius, và chuẩn Frobenius của ma trận không đổi qua phép chuyển vị, nên

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \sum_{i=1}^m \left\| (\mathbf{x}^{(i)})^T - (\mathbf{x}^{(i)})^T \mathbf{d}\mathbf{d}^T \right\|_F^2 \quad \text{sao cho} \quad \|\mathbf{d}\|_2 = 1. \quad (2.59)$$

Bây giờ, ta viết lại bài toán dưới dạng hữu ích hơn, thông qua ma trận thiết kế của các mẫu, thay vì lấy tổng qua các vectơ mẫu riêng lẻ. Điều này sẽ cho phép ta sử dụng ký hiệu ngắn gọn hơn. Gọi $\mathbf{X} \in \mathbb{R}^{m \times n}$ là ma trận được định nghĩa bằng cách xếp chồng tất cả các vectơ mô tả các điểm dữ liệu, sao cho $\mathbf{X}_{i,:} = (\mathbf{x}^{(i)})^T$. Ta có thể viết lại bài toán dưới dạng

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \|\mathbf{X} - \mathbf{X}\mathbf{d}\mathbf{d}^T\|_F^2 \quad \text{sao cho} \quad \mathbf{d}^T \mathbf{d} = 1. \quad (2.60)$$

Tạm thời bỏ qua ràng buộc, ta có thể đơn giản hóa phần chuẩn Frobenius như sau:

$$\arg \min_{\mathbf{d}} \|\mathbf{X} - \mathbf{X}\mathbf{d}\mathbf{d}^T\|_F^2 \quad (2.61)$$

$$= \arg \min_{\mathbf{d}} \text{Tr} \left((\mathbf{X} - \mathbf{X}\mathbf{d}\mathbf{d}^T)^T (\mathbf{X} - \mathbf{X}\mathbf{d}\mathbf{d}^T) \right) \quad (2.62)$$

(do phương trình (2.38))

$$= \arg \min_{\mathbf{d}} \text{Tr} (\mathbf{X}^T \mathbf{X} - \mathbf{X}^T \mathbf{X}\mathbf{d}\mathbf{d}^T - \mathbf{d}\mathbf{d}^T \mathbf{X}^T \mathbf{X} + \mathbf{d}\mathbf{d}^T \mathbf{X}^T \mathbf{X}\mathbf{d}\mathbf{d}^T) \quad (2.63)$$

$$= \arg \min_{\mathbf{d}} \text{Tr} (\mathbf{X}^T \mathbf{X}) - \text{Tr} (\mathbf{X}^T \mathbf{X}\mathbf{d}\mathbf{d}^T) - \text{Tr} (\mathbf{d}\mathbf{d}^T \mathbf{X}^T \mathbf{X}) + \text{Tr} (\mathbf{d}\mathbf{d}^T \mathbf{X}^T \mathbf{X}\mathbf{d}\mathbf{d}^T) \quad (2.64)$$

$$= \arg \min_{\mathbf{d}} -\text{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T) - \text{Tr}(\mathbf{d} \mathbf{d}^T \mathbf{X}^T \mathbf{X}) + \text{Tr}(\mathbf{d} \mathbf{d}^T \mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T) \quad (2.65)$$

(vì số hạng không phụ thuộc vào \mathbf{d} thì không ảnh hưởng tới $\arg \min$) (2.66)

$$= \arg \min_{\mathbf{d}} -2\text{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T) + \text{Tr}(\mathbf{d} \mathbf{d}^T \mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T) \quad (2.67)$$

(vì theo phương trình (2.41), có thể hoán vị vòng các ma trận trong toán tử vết)

$$= \arg \min_{\mathbf{d}} -2\text{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T) + \text{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T \mathbf{d} \mathbf{d}^T) \quad (2.68)$$

(dùng tiếp tính chất trên)

Bây giờ ta nhắc lại bài toán có ràng buộc:

$$\arg \min_{\mathbf{d}} -2\text{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T) + \text{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T \mathbf{d} \mathbf{d}^T) \quad \text{sao cho} \quad \mathbf{d}^T \mathbf{d} = 1 \quad (2.69)$$

$$= \arg \min_{\mathbf{d}} -2\text{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T) + \text{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T) \quad \text{sao cho} \quad \mathbf{d}^T \mathbf{d} = 1 \quad (2.70)$$

(theo ràng buộc trên)

$$= \arg \min_{\mathbf{d}} -\text{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T) \quad \text{sao cho} \quad \mathbf{d}^T \mathbf{d} = 1 \quad (2.71)$$

$$= \arg \max_{\mathbf{d}} \text{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T) \quad \text{sao cho} \quad \mathbf{d}^T \mathbf{d} = 1 \quad (2.72)$$

$$= \arg \max_{\mathbf{d}} \text{Tr}(\mathbf{d}^T \mathbf{X}^T \mathbf{X} \mathbf{d}) \quad \text{sao cho} \quad \mathbf{d}^T \mathbf{d} = 1 \quad (2.73)$$

Bài toán tối ưu này có thể giải bằng phương pháp phân tích giá trị riêng. Cụ thể, vectơ tối ưu \mathbf{d} là vectơ riêng của ma trận $\mathbf{X}^T \mathbf{X}$ ứng với giá trị riêng lớn nhất.

Lập luận này chỉ áp dụng cho trường hợp $\ell = 1$ và chỉ tìm được thành phần chính thứ nhất. Tổng quát hơn, khi ta muốn tìm một cơ sở của các thành phần chính, ma trận \mathbf{D} được cho bởi ℓ vector riêng ứng với các giá trị riêng lớn nhất. Điều này có thể được chứng minh bằng phương pháp quy nạp. Bạn đọc có thể tự chứng minh xem như một bài tập.

2.13 Không gian Euclid

Cho không gian vectơ V trên trường số thực. Xét hàm f gán tương ứng mỗi cặp phần tử (\mathbf{x}, \mathbf{y}) , trong đó $\mathbf{x}, \mathbf{y} \in V$, với một số thực $f(\mathbf{x}, \mathbf{y})$. Hàm f được gọi là **song tuyến tính**, nếu khi cố định thành phần này, thì nó là ánh xạ tuyến tính theo thành phần kia, tức là với mọi $\mathbf{x}, \mathbf{x}', \mathbf{y}, \mathbf{y}' \in V$ và $\alpha \in \mathbb{R}$, điều kiện sau phải thỏa mãn:

$$\begin{aligned} f(\mathbf{x} + \mathbf{x}', \mathbf{y}) &= f(\mathbf{x}, \mathbf{y}) + f(\mathbf{x}', \mathbf{y}) \quad \text{và} \quad f(\alpha \mathbf{x}, \mathbf{y}) = \alpha f(\mathbf{x}, \mathbf{y}) \\ f(\mathbf{x}, \mathbf{y} + \mathbf{y}') &= f(\mathbf{x}, \mathbf{y}) + f(\mathbf{x}, \mathbf{y}') \quad \text{và} \quad f(\mathbf{x}, \alpha \mathbf{y}) = \alpha f(\mathbf{x}, \mathbf{y}) \end{aligned} \quad (2.74)$$

Dạng song tuyến tính f trên V được gọi là

- **đối xứng** nếu

$$f(\mathbf{x}, \mathbf{y}) = f(\mathbf{y}, \mathbf{x}) \quad \text{với mọi } \mathbf{x}, \mathbf{y} \in V. \quad (2.75)$$

- **xác định dương** nếu

$$f(\mathbf{x}, \mathbf{x}) \geq 0 \quad \text{với mọi } \mathbf{x} \in V, \quad \text{và} \quad f(\mathbf{x}, \mathbf{x}) = 0 \Rightarrow \mathbf{x} = \mathbf{0}. \quad (2.76)$$

Nếu f là một dạng song tuyến tính đối xứng xác định dương trên V thì f được gọi là một **tích vô hướng** trên V , và V được gọi là **không gian Euclid**. Khi đó ta ký hiệu $f(\mathbf{x}, \mathbf{y})$ bởi $\langle \mathbf{x}, \mathbf{y} \rangle$ và đọc là “tích vô hướng \mathbf{x} với \mathbf{y} ”. Ta cũng viết lại các điều kiện của tích vô hướng theo ký hiệu mới, lưu ý rằng do tính đối xứng, nên chỉ cần đảm bảo tính tuyến tính theo một biến là đủ:

$$\begin{aligned} \langle \mathbf{x} + \mathbf{x}', \mathbf{y} \rangle &= \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{x}', \mathbf{y} \rangle \quad \text{và} \quad \langle \alpha \mathbf{x}, \mathbf{y} \rangle = \alpha \langle \mathbf{x}, \mathbf{y} \rangle \\ \langle \mathbf{x}, \mathbf{y} \rangle &= \langle \mathbf{y}, \mathbf{x} \rangle \\ \langle \mathbf{x}, \mathbf{x} \rangle &\geq 0 \\ \langle \mathbf{x}, \mathbf{x} \rangle &= 0 \Rightarrow \mathbf{x} = \mathbf{0} \end{aligned}$$

Trong không gian Euclid, tích vô hướng $\langle \mathbf{x}, \mathbf{y} \rangle$ sẽ sinh ra một chuẩn:

$$\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}. \quad (2.77)$$

Một ví dụ đơn giản, trong không gian \mathbb{R}^n , với

$$\mathbf{x} = (x_1, x_2, \dots, x_n)^T, \quad \mathbf{y} = (y_1, y_2, \dots, y_n)^T \in \mathbb{R}^n$$

có tích vô hướng

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^n x_i y_i \quad (2.78)$$

và chuẩn sinh bởi tích vô hướng

$$\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{\sum_{i=1}^n |x_i|^2} \quad (2.79)$$

chính là chuẩn L^2 trong \mathbb{R}^n .

... đang soạn!

Đại số tuyến tính là một trong những ngành toán học cơ bản cần thiết để hiểu về học máy. Một lĩnh vực toán học quan trọng khác thường xuyên xuất hiện trong học máy là lý thuyết xác suất, sẽ được trình bày tiếp theo.

Bài tập Chương 2

1. Cho hai hàm vectơ m thành phần với n biến $\mathbf{f}, \mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ xác định bởi: với $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$, $\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_m(\mathbf{x})]^T$ trong đó $f_i(\mathbf{x}) = f_i(x_1, x_2, \dots, x_n)$ là các hàm số n biến. Ta quy ước ký hiệu $\mathbf{f}(\mathbf{x})$ bởi \mathbf{f} , và tương tự đối với các hàm khác. Ký hiệu $\mathbf{J}(\mathbf{f}) = \left(\frac{\partial f_i}{\partial x_j} \right)$ là **ma trận Jacobi** của \mathbf{f} tại \mathbf{x} , $i = \overline{1, m}$ là chỉ số hàng, $j = \overline{1, n}$ là chỉ số cột. Chứng minh

a) $\mathbf{J}(\mathbf{f} + \mathbf{g}) = \mathbf{J}(\mathbf{f}) + \mathbf{J}(\mathbf{g})$

b) $\mathbf{J}(\alpha \mathbf{f}) = \alpha \mathbf{J}(\mathbf{f})$, trong đó $\alpha \in \mathbb{R}$

2. Cho trước ma trận \mathbf{A} cỡ $m \times n$, vectơ \mathbf{b} có chiều m , và vectơ \mathbf{x} có n biến. Chứng minh

a) $\mathbf{J}(\mathbf{Ax}) = \mathbf{A}$

b) $\mathbf{J}(\mathbf{b}) = \mathbf{0}$, từ đó suy ra $\mathbf{J}(\mathbf{Ax} - \mathbf{b}) = \mathbf{A}$

3 (*). Định nghĩa **gradient của hàm số** là vectơ mà mỗi thành phần lần lượt là đạo hàm riêng của hàm số đó theo mỗi biến. Cho các hàm vectơ \mathbf{f}, \mathbf{g} xác định như trong bài 1. Khi đó $\mathbf{f}^T \mathbf{g}$ là một hàm số. Chứng minh

$$\nabla (\mathbf{f}^T \mathbf{g}) = \mathbf{J}(\mathbf{f})^T \mathbf{g} + \mathbf{J}(\mathbf{g})^T \mathbf{f}.$$

4. Cho trước ma trận \mathbf{A}, \mathbf{B} cỡ $m \times n$, vectơ \mathbf{c}, \mathbf{d} có chiều m , và vectơ \mathbf{x} có n biến. Chứng minh

a) $\nabla [(\mathbf{Ax} - \mathbf{b})^T (\mathbf{Cx} - \mathbf{d})] = \mathbf{A}^T (\mathbf{Cx} - \mathbf{d}) + \mathbf{C}^T (\mathbf{Ax} - \mathbf{b})$

b) $\nabla [\|\mathbf{Ax} - \mathbf{b}\|_2^2] = 2\mathbf{A}^T (\mathbf{Ax} - \mathbf{b})$