

1.2 수열과 집합의 합과 곱

데이터를 분석하기 위해서는 많은 숫자의 합이나 곱을 계산해야 한다. 따라서 숫자의 합과 곱을 나타내는 수학 기호에 익숙해지는 것은 데이터 분석의 첫걸음이다. 이 절에 나온 기호와 수식들은 앞으로 계속 반복하여 나오므로 반드시 외우고 손으로 여러 번 쓰기를 바란다.

수열

수열(sequence)은 N 개 숫자 또는 변수가 순서대로 나열된 것이다. 다음은 수열의 예다.

$$1, 2, 3, 4 \quad (1.2.1)$$

$$x_1, x_2, x_3, x_4, x_5, x_6 \quad (1.2.2)$$

문자에 붙은 아래 첨자는 순서를 나타내는 숫자로서 인덱스(index)라고 부른다.

수열이 아주 길거나 수열의 길이가 숫자가 아닌 문자인 경우에는 ... (dots) 기호를 사용하여 다음처럼 가운데 부분을 생략할 수 있다.

$$x_1, x_2, \dots, x_N \quad (1.2.3)$$

집합

순서가 중요하지 않은 숫자들은 집합(set)으로 표시한다.

$$\{1, 2, 3, 4\} \quad (1.2.4)$$

$$\{x_1, x_2, x_3, x_4, x_5, x_6\} \quad (1.2.5)$$

집합에서도 원소가 많으면 가운데를 생략할 수 있다.

$$\{x_1, x_2, \dots, x_N\} \quad (1.2.6)$$

데이터 분석에서는 1부터 N 까지의 수열 또는 집합이 자주 나오기 때문에 위에서 사용한 기호 대신 다음과 같이 더 간단한 기호를 쓰는 경우도 많다.

$$x_{1:N} \quad (1.2.7)$$

$$\{x_i\}_N \quad (1.2.8)$$

집합에 알파벳 대문자로 이름을 붙일 수도 있다. 데이터 분석에서 자주 나오는 집합 중의 하나는 1, -2, 3.14와 같은 실수(real number) 전체의 집합이다. 실수 집합은 \mathbf{R} 이라는 이름을 가진다. 어떤 숫자 x 가 실수이면 집합 \mathbf{R} 에 포함되므로 다음과 같은 기호로 나타낸다.

$$x \in \mathbf{R} \quad (1.2.9)$$

만약 두 개의 숫자로 이루어진 숫자 쌍 (x_1, x_2) 가 있고 각각의 숫자 x_1, x_2 가 모두 실수라면 다음처럼 표시한다.

$$(x_1, x_2) \in \mathbf{R} \times \mathbf{R} \quad (1.2.10)$$

또는

$$(x_1, x_2) \in \mathbf{R}^2 \quad (1.2.11)$$

수열의 합과 곱

다음 기호는 수열을 더하거나 곱하는 연산을 짧게 줄여 쓴 것이다. 그리스 문자의 시그마(Σ)와 파이(Π)를 본 따서 만든 기호지만 시그마와 파이로 읽지 않고 영어로 썸(sum), 프로덕트(product)라고 읽는다. 합과 곱 기호 아래에는 인덱스의 시작값, 위에는 인덱스의 끝값을 표시한다. 곱셈은 알파벳 x와 혼동될 수 있기 때문에 $a \times b$ 가 아니라 $a \cdot b$ 와 같이 점(dot)으로 표시하거나 아예 생략한다.

$$\sum_{i=1}^N x_i = x_1 + x_2 + \dots + x_N \quad (1.2.12)$$

$$\prod_{i=1}^N x_i = x_1 \cdot x_2 \cdot \dots \cdot x_N \quad (1.2.13)$$

더하거나 곱하기를 반복해서 써야할 때는 합과 곱 기호를 사용하면 수식이 간결하고 명확해진다. 예를 들어 1부터 4까지 더해야 하는 경우에는 다음처럼 쓴다.

$$\sum_{i=1}^4 i = 1 + 2 + 3 + 4 \quad (1.2.14)$$

만약 10부터 90까지 10씩 증가하는 수열을 모두 더해야 한다면 다음처럼 쓴다.

$$\sum_{k=1}^9 10k = 10 \cdot 1 + 10 \cdot 2 + \dots + 10 \cdot 9 = 10 + 20 + \dots + 90 \quad (1.2.15)$$

곱셈도 마찬가지이다. 다음은 10부터 20까지의 수를 모두 곱하는 식이다.

$$\prod_{i=10}^{20} i = (10) \cdot (11) \cdot \dots \cdot (20) \quad (1.2.16)$$

합이나 곱을 중첩하여 여러 번 쓰는 경우도 있다. 합과 곱을 중첩하여 쓸 때는 다음처럼 괄호를 생략할 수 있다. 합이나 곱이 중첩된 경우에는 인덱스가 여러 개가 된다.

$$\sum_{i=1}^N \left(\sum_{j=1}^M x_{ij} \right) = \sum_{i=1}^N \sum_{j=1}^M x_{ij} \quad (1.2.17)$$

다음은 합과 곱을 중첩한 수식의 예이다.

$$\begin{aligned} \sum_{i=1}^2 \sum_{j=1}^3 (i+j) &= \sum_{i=1}^2 \left(\sum_{j=1}^3 (i+j) \right) \\ &= \sum_{i=1}^2 ((i+1) + (i+2) + (i+3)) \\ &= ((1+1) + (1+2) + (1+3)) + ((2+1) + (2+2) + (2+3)) \\ \prod_{m=1}^3 \prod_{n=1}^2 (m+2n) &= \prod_{m=1}^3 \left(\prod_{n=1}^2 (m+2n) \right) \\ &= \prod_{m=1}^3 ((m+2 \cdot 1) \cdot (m+2 \cdot 2)) \\ &= ((1+2 \cdot 1) \cdot (1+2 \cdot 2)) \cdot ((2+2 \cdot 1) \cdot (2+2 \cdot 2)) \cdot ((3+2 \cdot 1) \cdot (3+2 \cdot 2)) \end{aligned} \quad (1.2.18)$$

거듭 강조하지만 수학 공부는 눈으로 읽기만 하고 손으로 쓰지 않으면 아무런 의미가 없다. 지금까지 나온 수식을 꼭 손으로 반복하여 쓰면서 의미를 익히기 바란다.

연습 문제 1.2.1

다음 수식을 풀어 써라. 이 수식들은 이후에 머신러닝 모형에 등장할 수식이다.

(1) 이 식은 분류 모형 중의 하나인 서포트 벡터 머신(support vector machine) 모형에 나온다.

$$\sum_{i=1}^3 \sum_{j=1}^3 a_i a_j y_i y_j x_i x_j \quad (1.2.20)$$

(2) 이 식은 특잇값 분해(singular value decomposition)에 나온다.

$$\sum_{k=1}^3 \sum_{i=1}^3 \sigma_i^2 (v_i w_k)^2 \quad (1.2.21)$$

(3) 이 식은 카테고리 분포(categorical distribution)의 추정에 사용된다.

$$\prod_{i=1}^4 \prod_{k=1}^4 \theta_k^{x_{i,k}} \quad (1.2.22)$$

(4) 가우시안 혼합 모형(Gaussian mixture model)에 다음과 비슷한 수식이 나온다.

$$\prod_{i=1}^4 \sum_{k=1}^2 \pi_k x_i \mu_k \quad (1.2.23)$$

연습 문제 1.2.2

수열의 곱은 여러개의 값 중 하나를 선택하는 경우에도 쓰일 수 있다. 수열 x_i 가 다음과 같다고 하자.

$$x_i : x_1, x_2, x_3, x_4 \quad (1.2.24)$$

이 값 중 하나만 선택하고 싶다면 다음처럼 모두 0이고 하나만 1인 수열 y_i 를 사용하면 된다.

$$y_i : 0, 1, 0, 0 \quad (1.2.25)$$

(1) x_i 와 y_i 가 위와 같을 때 다음 값을 계산하라.

$$\prod_i x_i^{y_i} \quad (1.2.26)$$

(2) 만약 수열 y_i 에서 $y_3 = 1$ 이고 나머지값이 0이라면 답이 어떻게 달라지는가?

수열의 합과 곱 연산은 다음과 같은 성질을 가지고 있다.

(1) 인덱스 문자가 바뀌어도 실제 수식은 달라지지 않는다.

$$\sum_{i=1}^N x_i = \sum_{j=1}^N x_j \quad (1.2.27)$$

(2) 상수 c 를 곱한 후에 합을 한 결과는 먼저 합을 구하고 상수를 곱한 것과 같다.

$$\sum_{i=1}^N cx_i = c \sum_{i=1}^N x_i \quad (1.2.28)$$

(3) 더해야 하는 값들이 여러 항의 합으로 되어 있으면 각각의 합을 먼저 구한 후에 더해도 된다.

$$\sum_{i=1}^N (x_i + y_i) = \sum_{i=1}^N x_i + \sum_{i=1}^N y_i \quad (1.2.29)$$

(4) 합이나 곱을 중첩하는 경우에는 중첩의 순서를 바꾸어도 결과가 같다.

$$\sum_{i=1}^N \sum_{j=1}^M = \sum_{j=1}^M \sum_{i=1}^N \quad (1.2.30)$$

$$\prod_{i=1}^N \prod_{j=1}^M = \prod_{j=1}^M \prod_{i=1}^N \quad (1.2.31)$$

예를 들어 다음 두 식은 항들의 순서만 바뀌었고 그 합은 같다는 것을 알 수 있다.

$$\sum_{i=1}^2 \sum_{j=1}^3 x_{ij} = (x_{11} + x_{12} + x_{13}) + (x_{21} + x_{22} + x_{23}) \quad (1.2.32)$$

$$\sum_{j=1}^3 \sum_{i=1}^2 x_{ij} = (x_{11} + x_{21}) + (x_{12} + x_{22}) + (x_{13} + x_{23}) \quad (1.2.33)$$

연습 문제 1.2.3

다음 두 식의 좌변과 우변이 같음을 증명하라. (힌트: 등호의 왼쪽과 오른쪽 각각의 식을 풀어서 같아짐을 보인다.) 이 수식들은 선형대수에서 벡터 및 행렬의 곱에 유용하게 사용된다.

(1)

$$\left(\sum_{i=1}^3 x_i \right)^2 = \sum_{i=1}^3 \sum_{j=1}^3 x_i x_j \quad (1.2.34)$$

(2)

$$\sum_{i=1}^3 \sum_{j=1}^3 x_i y_{ij} = \sum_{i=1}^3 \left(x_i \sum_{j=1}^3 y_{ij} \right) \quad (1.2.35)$$

집합의 합과 곱

수열이 아니라 집합의 원소들의 합과 곱을 구할 때는 인덱스 대신 집합 기호를 사용한다.

만약 집합 X 의 원소가 다음과 같다면,

$$X = \{x_1, x_2, x_3\} \quad (1.2.36)$$

이 집합의 원소의 합과 곱은 다음처럼 표시한다. 이 때는 합과 곱 기호 안에 인덱스가 없다.

$$\sum_X x = x_1 + x_2 + x_3 \quad (1.2.37)$$

$$\prod_X x = x_1 \cdot x_2 \cdot x_3 \quad (1.2.38)$$

원소 중에서 특정한 조건을 가진 원소만 포함시키거나 제외하여 합과 곱을 구하는 경우도 있다. 이 때는 인덱스 위치에 조건을 표시한다. 예를 들어 다음 식은 집합 X 의 원소 중 0이 아닌 것만 곱한 값을 뜻한다.

$$\prod_{x \in X, x \neq 0} x \quad (1.2.39)$$

연습 문제 1.2.4

두 집합 X_1, X_2 가 있고 x_1 은 X_1 의 원소들, x_2 은 X_2 의 원소들을 가리킬 때 다음 두 식의 좌변과 우변이 같음을 증명하라. 문제를 간단하게 하기 위해 여기에서는 각각의 집합이 3개의 원소만 가지고 있다고 가정하자. 이 식의 확장된 버전은 추후 베이지안 네트워크의 합-곱(sum-product) 알고리즘에 사용된다.

$$\prod_{i=1}^2 \sum_{X_i} x_i = \sum_{X_1 \times X_2} \prod_{i=1}^2 x_i \quad (1.2.40)$$