

MSBD5004 Mathematical Methods for Data Analysis

Homework 1

Due date: 30 September, Friday

1. Consider the vector space \mathbb{R}^n .

(a) Check that $\|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} |x_i|$ is indeed a norm on \mathbb{R}^n .

(b) Prove that: for any $\mathbf{x} \in \mathbb{R}^n$,

$$\|\mathbf{x}\|_\infty = \lim_{p \rightarrow \infty} \|\mathbf{x}\|_p.$$

(c) Prove the inequality

$$\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_1 \leq n\|\mathbf{x}\|_\infty, \quad \forall \mathbf{x} \in \mathbb{R}^n.$$

2. For any $\mathbf{A} \in \mathbb{R}^{m \times n}$, we have defined

$$\|\mathbf{A}\|_2 = \max_{\mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2$$

(a) Prove that $\|\cdot\|_2$ is a norm on $\mathbb{R}^{m \times n}$.

(b) Prove that $\|\mathbf{A}\mathbf{x}\|_2 \leq \|\mathbf{A}\|_2 \|\mathbf{x}\|_2$ for any $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{x} \in \mathbb{R}^n$.

(c) Prove that $\|\mathbf{A}\mathbf{B}\|_2 \leq \|\mathbf{A}\|_2 \|\mathbf{B}\|_2$ for all $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{n \times p}$.

3. Let a_1, a_2, \dots, a_m be m given real numbers. Prove that a median of a_1, a_2, \dots, a_m minimizes

$$\sum_{i=1}^m |a_i - b|$$

over all $b \in \mathbb{R}$.

4. Suppose that the vectors $\mathbf{x}_1, \dots, \mathbf{x}_N$ in \mathbb{R}^n are clustered using the K -means algorithm, with group representatives $\mathbf{z}_1, \dots, \mathbf{z}_k$.

(a) Suppose the original vectors \mathbf{x}_i are nonnegative, i.e., their entries are nonnegative. Explain why the representatives \mathbf{z}_j output by the K -means algorithm are also nonnegative.

(b) Suppose the original vectors \mathbf{x}_i represent proportions, i.e., their entries are nonnegative and sum to one. (This is the case when \mathbf{x}_i are word count histograms, for example.) Explain why the representatives \mathbf{z}_j output by the K -means algorithm also represent proportions (i.e., their entries are nonnegative and sum to one).

(c) Suppose the original vectors \mathbf{x}_i are Boolean, i.e., their entries are either 0 or 1. Give an interpretation of $(\mathbf{z}_j)_i$, the i -th entry of the j group representative.

5. Suppose that the vectors $\mathbf{x}_1, \dots, \mathbf{x}_N$ in \mathbb{R}^n are clustered using the K -means algorithm, with groups G_1, \dots, G_k and group representatives $\mathbf{z}_1, \dots, \mathbf{z}_k$. Prove that the sequence of objective function values

$$J = \sum_{j=1}^k \sum_{i \in G_j} \|\mathbf{x}_i - \mathbf{z}_j\|_2^2$$

produced by the K -means algorithm is non-increasing.

6. Consider the set of infinite sequences. We have defined the vector space $\ell_\infty := \{\mathbf{a} \mid \|\mathbf{a}\|_\infty < +\infty\}$ with norm $\|\cdot\|_\infty$. For each $k \geq 1$, let $\mathbf{a}_k = \{\underbrace{1, 1, \dots, 1}_{k \text{ times}}, 0, 0, \dots\} \in \ell_\infty$. Prove that $\mathbf{a}_k \not\rightarrow \mathbf{b} := (1, 1, 1, \dots)$ in ℓ_∞ with $\|\cdot\|_\infty$.
7. Let X be a vector space with norm $\|\cdot\|$. Prove the following results.
- (a) If $\{c_k\}$ is a convergent sequence in \mathbb{R} and $\{\mathbf{a}_k\}$ is a convergent sequence in X with limits c and \mathbf{a} respectively, then $\{c_k \cdot \mathbf{a}_k\}$ is a convergent sequence in X with limit $c \cdot \mathbf{a}$.
 - (b) If $\{\mathbf{a}_k\}$ and $\{\mathbf{b}_k\}$ are convergent sequences in X with limits \mathbf{a} and \mathbf{b} respectively, then $\{\mathbf{a}_k + \mathbf{b}_k\}$ is a convergent sequence in X with limit $\mathbf{a} + \mathbf{b}$.