

# Useful distance functions for machine learning

## Topics we'll cover

- ①  $L_p$  norms
- ② Metric spaces

## Measuring distance in $\mathbb{R}^m$

Usual choice: **Euclidean distance**:

$$\|x - z\|_2 = \sqrt{\sum_{i=1}^m (x_i - z_i)^2}.$$

For  $p \geq 1$ , here is  $\ell_p$  **distance**:

$$\|x - z\|_p = \left( \sum_{i=1}^m |x_i - z_i|^p \right)^{1/p}$$

- $p = 2$ : Euclidean distance
- $\ell_1$  distance:  $\|x - z\|_1 = \sum_{i=1}^m |x_i - z_i|$
- $\ell_\infty$  distance:  $\|x - z\|_\infty = \max_i |x_i - z_i|$

## Example 1

Consider the all-ones vector  $(1, 1, \dots, 1)$  in  $\mathbb{R}^d$ .  
What are its  $\ell_2$ ,  $\ell_1$ , and  $\ell_\infty$  length?

## Example 2

In  $\mathbb{R}^2$ , draw all points with:

- ①  $\ell_2$  length 1
- ②  $\ell_1$  length 1
- ③  $\ell_\infty$  length 1

## Metric spaces

Let  $\mathcal{X}$  be the space in which data lie.

A distance function  $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a **metric** if it satisfies these properties:

- $d(x, y) \geq 0$  (nonnegativity)
- $d(x, y) = 0$  if and only if  $x = y$
- $d(x, y) = d(y, x)$  (symmetry)
- $d(x, z) \leq d(x, y) + d(y, z)$  (triangle inequality)

## Example 1

$\mathcal{X} = \mathbb{R}^m$  and  $d(x, y) = \|x - y\|_p$

Check:

- $d(x, y) \geq 0$  (nonnegativity)
- $d(x, y) = 0$  if and only if  $x = y$
- $d(x, y) = d(y, x)$  (symmetry)
- $d(x, z) \leq d(x, y) + d(y, z)$  (triangle inequality)

## Example 2

$\mathcal{X} = \{\text{strings over some alphabet}\}$  and  $d = \text{edit distance}$

Check:

- $d(x, y) \geq 0$  (nonnegativity)
- $d(x, y) = 0$  if and only if  $x = y$
- $d(x, y) = d(y, x)$  (symmetry)
- $d(x, z) \leq d(x, y) + d(y, z)$  (triangle inequality)

## A non-metric distance function

Let  $p, q$  be probability distributions on some set  $\mathcal{X}$ .

The **Kullback-Leibler divergence** or **relative entropy** between  $p, q$  is:

$$d(p, q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)}.$$