# 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**:

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

$\ell_p$ distance between 2 vectors $x$ & $z$,

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: $\|x - z\|_1 = \sum_{i=1}^{m}|x_i - z_i|$   $V + y$
- $\ell_\infty$ distance: $\|x - z\|_\infty = \max_i |x_i - z_i|$   $y$

differ abs
each coordi

$$x = (1, 1, \cdots, 1)$$

# Example 1

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

$\ell_2$ m

$\|x\|_2$

$= \sqrt{1^2 + 1^2 + \ldots + 1^2}$

$= \sqrt{d}$

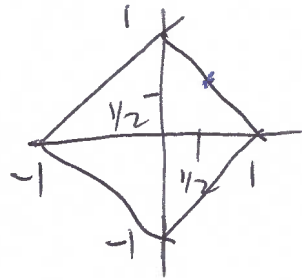$\ell_1$ m

$\|x\|_1 =$

$|x_1| + \ldots + |x_d|$
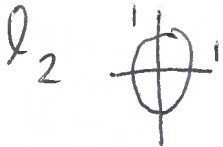
$= d$

$0$

$\|x\|_\alpha = 1$

$0$

# Example 2

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

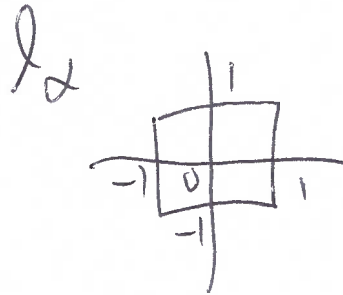① $\ell_2$ length 1
② $\ell_1$ length 1
③ $\ell_\infty$ length 1

$$\ell_1 : \{(x_1, x_2) : |x_1| + |x_2| = 1\}$$

unit ball for $\ell_1$

$\ell_2$

$$\{(x_1, x_2) : \sqrt{x_1^2 + x_2^2} = 1\}$$

$\ell_\infty$

# Metric spaces
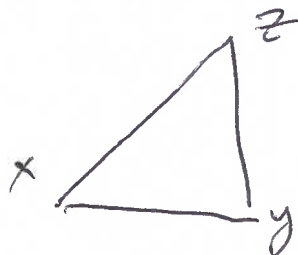
$$d(x, x') = 3.6$$

data

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

A distance function $d : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ 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 $d(x, y) = \|x - y\|_p$

$$d(x,y) = \sum_{i=1}^{m} |x_i - y_i|$$

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) ✓

$$|x_i - z_i| \leq |x_i - y_i| + |y_i - z_i|$$

Sum over all $i$

$\ell_1$ satisfies metric distance 4 properties

# Example 2

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

$$x = \{A, C, G, T\}^*$$
$$x = ACCGT$$
$$y = CCGT$$

edit distance

$$d(x,y) = \# \text{ of insertion, deletion, substitutes}$$
$$\text{to get from } x \text{ to } y$$

$$d(x,y) \geq 0$$
$$d(x,y) = 0 \iff x = y$$
$$d(x,y) = d(y,x)$$
$$\text{triangle inequality}$$

## A non-metric distance function

$$d \to \text{distance}$$
$$d(p,q) \text{ distance between } p \text{ \& } q$$

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

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

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

$$p = \left(\tfrac{1}{2}, \tfrac{1}{4}, \tfrac{1}{8}, \tfrac{1}{8}\right)$$
$$q = \left(\tfrac{1}{6}, \tfrac{1}{3}, \tfrac{1}{3}, \tfrac{1}{6}\right)$$

we can compute $l_1, l_2$ but

the choice
is changed $d(p,q) = \frac{1}{2} \log \frac{1/2}{1/6} + \frac{1}{4} \log \frac{1/4}{1/3} + \frac{1}{8} \log \frac{1/8}{1/3} + \frac{1}{8} \log \frac{1/8}{1/6}$

KL div $\not\Rightarrow$

$$w/c \text{ is } d(p,q)$$