# Useful distance functions for machine learning

## Topics we'll cover

1. $L_p$ norms

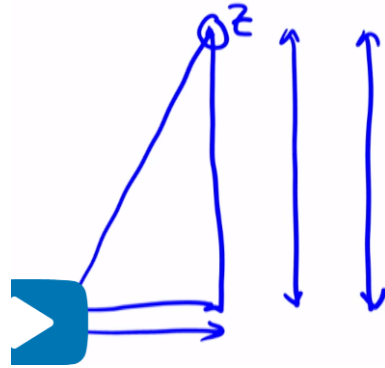2. Metric spaces

# 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 $(1, 1, \ldots, 1)$ in $\mathbb{R}^d$.
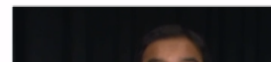What are its $\ell_2$, $\ell_1$, and $\ell_\infty$ length?

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

$\|x\|_2$

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

$= \sqrt{d}$

$\|x\|_1 =$

$|x_1| + \cdots + |x_d|$

$= d$

$\|x\|_\infty = 1$

# Example 2

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

&#10102; $\ell_2$ length 1

&#10103; $\ell_1$ length 1

&#10104; $\ell_\infty$ length 1

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

$\ell_2$

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

$\ell_\infty$



# Metric spaces

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 $d(x, y) = \|x - y\|_p$

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)

$$\mathcal{X} = \{A, C, G, T\}^*$$

$$x = A\ C\ C\ G\ T$$

$$y = C\ C\ G\ T$$

# 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)

$d(x,y) = $ # of insertions, deletions, substitutions to get from $x$ to $y$.

$d(x,y) \geq 0$

$d(y,y) = 0 \iff x = y$

$d(x,y) = d(y,x)$

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:

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

$p = (\tfrac{1}{2}, \tfrac{1}{4}, \tfrac{1}{8}, \tfrac{1}{8})$

$q = (\tfrac{1}{6}, \tfrac{1}{3}, \tfrac{1}{3}, \tfrac{1}{6})$

$d(p,q) = \tfrac{1}{2} \log \tfrac{1/2}{1/6} + \tfrac{1}{4} \log \tfrac{1/4}{1/3} + \tfrac{1}{8} \log \tfrac{1/8}{1/3}$
$\quad + \tfrac{1}{8} \log \tfrac{1/8}{1/6}$