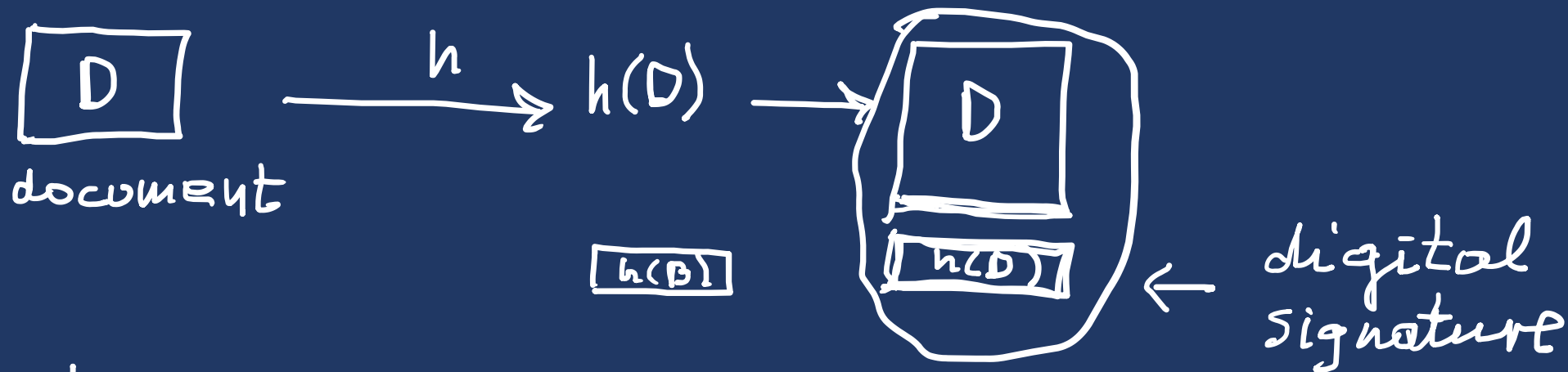


# LOCALITY SENSITIVE HASH FUNCTIONS FAMILIES

Standard hash functions, e.g.  $h \sim \text{MD5}$



$D'$  ← small modif. of  $D$

$$\text{rng}(h) \subseteq \{0,1\}^{256}$$

$h(D) \quad h(D')$  ← they are very different

$$\sum_{i=1}^{256} \mathbb{P}[h(D)_i \neq h(D')_i] \approx \frac{256}{2}$$

$D$  I will pay 100 \$  
to John Smith  
at 31.12.2020  
JCI  $\xrightarrow{h}$   $(1, 0, 1, 1, 1, 0, 0, 0, \dots)$

$D'$  I will pay 100 \$  
to J.S  
:  
:  $\xrightarrow{h}$   $(1, 1, 0, 0, 1, 0, 1, \dots)$

OUR GOAL :  $D \underset{\text{similar}}{\approx} D' \longrightarrow h(D) \underset{\substack{\uparrow \\ \text{sketch of } D}}{\text{similar}} h(D')$

# Metric space

$(X, d) \leftarrow$  metric space if

1)  $d: X \times X \longrightarrow [0, \infty)$  ( $d(x, y) \geq 0$ )

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

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

4)  $d(x, z) \leq d(x, y) + d(y, z)$

triangle inequality



EX 1  $X = \mathbb{R}^n$ ,  $p \geq 1$ .

$$\|\bar{x}\|_p = \sqrt[p]{\sum_{l=1}^n \|x_l\|^p}$$

THM. For each  $p \geq 1$ ,  $d_p$  is a metric

~~THM~~  $d_p(\bar{x}, \bar{y}) = \|\bar{x} - \bar{y}\|_p$

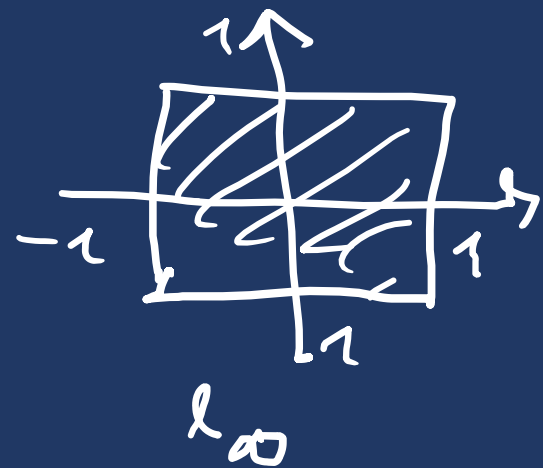
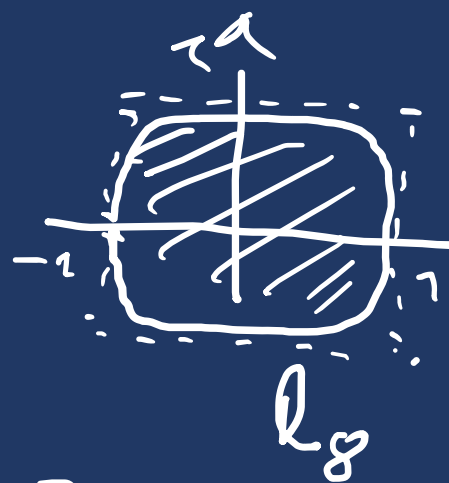
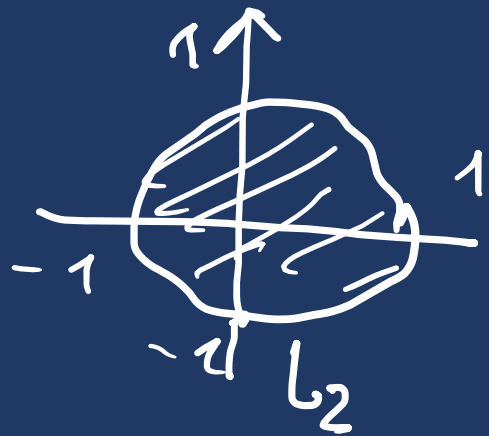
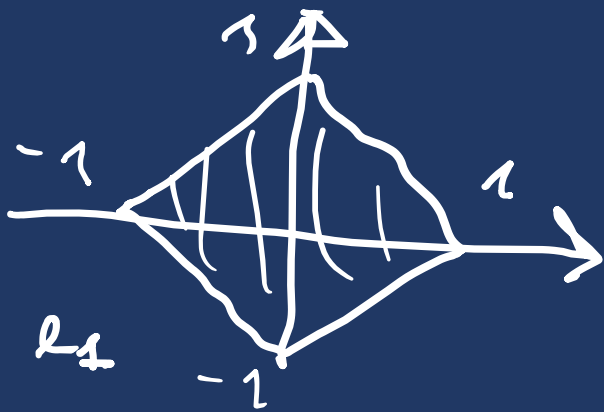
$p=2$ :  $d_2(\bar{x}, \bar{y}) = \sqrt{(x_1 - y_1)^2 + \dots + (x_n - y_n)^2}$  ← euclidean distance

$p=1$ :  $d_1(\bar{x}, \bar{y}) = \sum_{l=1}^n |x_l - y_l|$

$\lim_{p \rightarrow \infty} \|\bar{x}\|_p = \max \{ |x_l| : l=1, \dots, n \}$

$$\overline{B}(x, r) = \{y \in X : d(x, y) \leq r\}$$

$$\overline{B}_p(0, 1) = \{\bar{y} \in \mathbb{R}^n : \|\bar{y}\|_p \leq 1\}$$



$$(n=2) \overline{B}_1(0, 1) = \{(y_1, y_2) : |y_1| + |y_2| \leq 1\}$$

**EX 2**  $X = \{0, 1\}^n$ ;  $\bar{x}, \bar{y} \in X$

$$d_H(\bar{x}, \bar{y}) = \sum_{l=1}^n |\bar{x}_l - \bar{y}_l| \quad (\ell_1\text{-distance on } \mathbb{R}\{0, 1\}^n)$$

Hamming distance

$$= \sum_{l=1}^n \mathbb{I}[x_l \neq y_l] = \sum_i \mathbb{I}[i \in A_x \Delta A_y] = |A_x \Delta A_y|$$

$$\mathbb{I}[\varphi] = \begin{cases} 1 : \varphi \\ 0 : \neg\varphi \end{cases}$$

$$\begin{aligned} x_l \neq y_l &\equiv (x_l=1, y_l=0) \vee (x_l=0, y_l=1) \\ &\equiv (i \in A_x \setminus A_y) \vee (i \in A_y \setminus A_x) \equiv \end{aligned}$$

$$\{0, 1\}^n \ni \bar{x} \longrightarrow A_{\bar{x}} = \{i \in \{1, \dots, n\} : \bar{x}_i = 1\} \subseteq \{1, \dots, n\}$$

$$\rightarrow i \in A_x \Delta A_y$$

•  $(\mathcal{P}(\{1, \dots, n\}), d)$  ← metric space

$$d(A, B) = |A \Delta B|$$

• Fix any set  $X$ .

$$\mathcal{P}_{\text{fin}}(X) = \{D \subseteq X : |D| < \infty\}$$

$$d(A, B) = |A \Delta B|$$



$X = \text{ASCII}^*$

$A \subseteq X$  : documents  
+4

Main problem: our metrics are unbounded.

We want use such metrics  $d$  that

$$0 \leq d \leq 1.$$

standard way.

above line ✓

THM. Suppose  $f: [0, \infty) \rightarrow [0, \infty)$  s.t.

1)  $f(0) = 0$

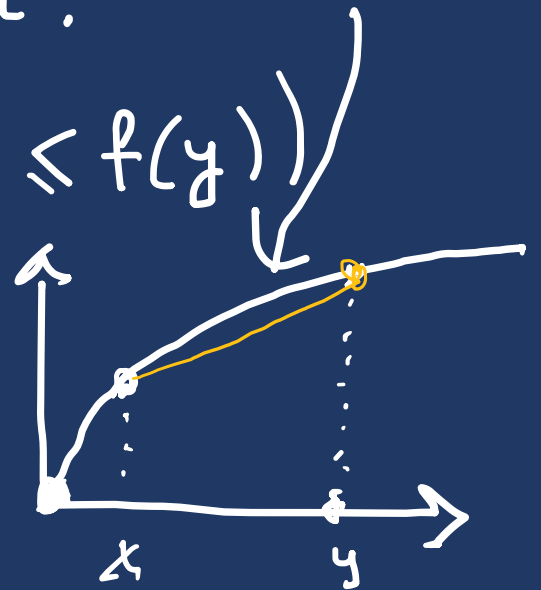
2)  $f$  is monotonic ( $0 \leq x < y \rightarrow f(x) \leq f(y)$ )

3)  $f$  is concave

Let  $(X, d)$  be a metric space. Then

$$\rho(x, y) = f(d(x, y))$$

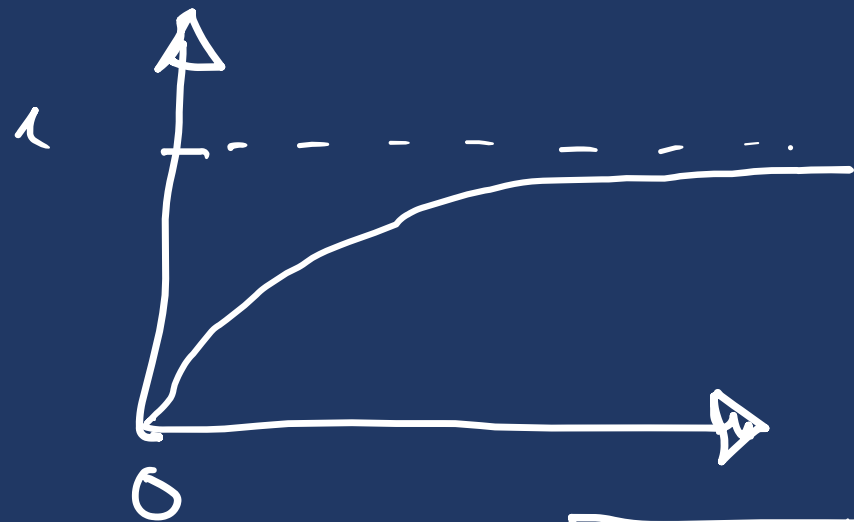
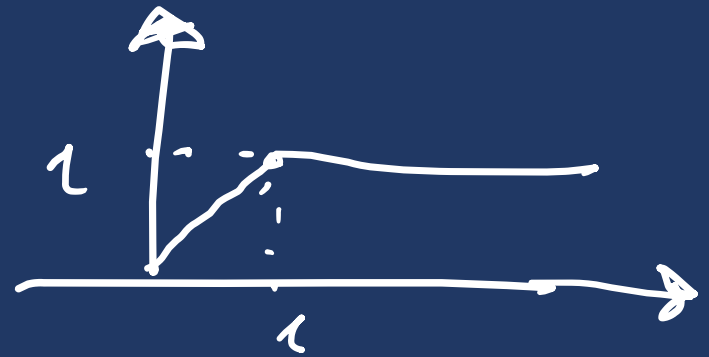
is a metric on  $X$ .





E1.  $f(x) = \min\{x, 1\}$

E2.  $f(x) = \frac{x}{1+x}$



$$f'(x) = \frac{1}{(1+x)^2} > 0$$

$$f''(x) = \frac{-2}{(1+x)^3} < 0$$

$f$  is concave

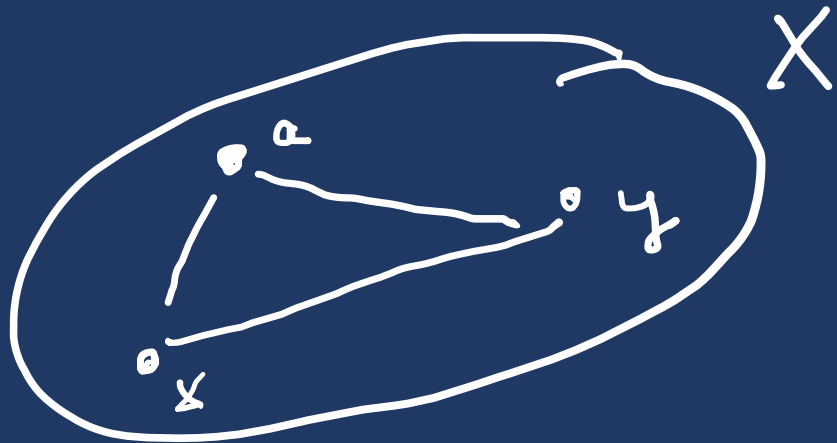
$$\phi(\bar{x}, \bar{y}) = \frac{\sqrt{\sum (x_i - y_i)^2}}{1 + \sqrt{\sum (x_i - y_i)^2}}$$

$\mathbb{R}^n$

THM (Steinhaus) Let  $(X, d)$  be metric space, let  $a \in X$ . Let

$$\rho(x, y) = \frac{2 \cdot d(x, y)}{d(x, a) + d(y, a) + d(x, y)}$$

Then  $(X, \rho)$  is a metric space.



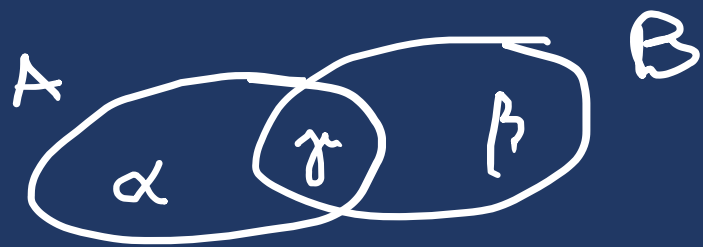
$$\underline{d(x, a) + d(a, y) + d(x, y)} \geq \underline{d(x, y) + d(x, y)}$$

$$\rho(x, y) \leq \frac{2 \cdot d(x, y)}{2 \cdot d(x, y)} = 1$$

EX. Suppose  $x \neq a$ . Calculate  $\rho(x, a)$  !!!

Application :  $|A \Delta B| = |A \Delta \phi| + |B \Delta \phi| + |A \Delta B|$ ,  $a = \phi$

$$\begin{aligned} P(A, B) &= \frac{2 \cdot |A \Delta B|}{|A \Delta \phi| + |B \Delta \phi| + |A \Delta B|} = \\ &= \frac{2 \cdot |A \Delta B|}{|A| + |B| + |A \Delta B|} = \frac{2(\alpha + \beta)}{(\alpha + \gamma) + (\beta + \gamma) + (\alpha + \beta)} \end{aligned}$$



$$= \frac{2 \cdot (\alpha + \beta)}{2(\alpha + \beta + \gamma)} = \frac{|A \Delta B|}{|A \cup B|}$$

$$d_J(A, B) = \frac{|A \triangle B|}{|A \cup B|}$$

Jaccard  
distance

- $d_J(A, B) \leq 1$

- $(d_J(A, B) = 1) \equiv (A \cap B = \emptyset)$

SIMILARITY:  $s: X \times X \rightarrow [0, 1]$

- $s(x, x) = 1$

- $s(x, y) = s(y, x)$

PERFECT SIMILARITY:  $s(x, y) = 1 - d(x, y)$

where  $d$  is a metric s.t.  $0 \leq d \leq 1$ .

$$1 - d_J(A, B) = 1 - \frac{|A \triangle B|}{|A \cup B|} =$$

$$= \frac{|A \cup B| - |A \triangle B|}{|A \cup B|} = \frac{(\alpha + \beta + \gamma) - (\alpha + \beta)}{|A \cup B|}$$



$$= \frac{\gamma}{|A \cup B|} = \frac{|A \cap B|}{|A \cup B|}$$

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

Jaccard  
similarity

# Estimation of $\tau(A, B)$

$$\Omega = \{1, 2, \dots, n\}$$

$$A, B \subseteq \Omega$$

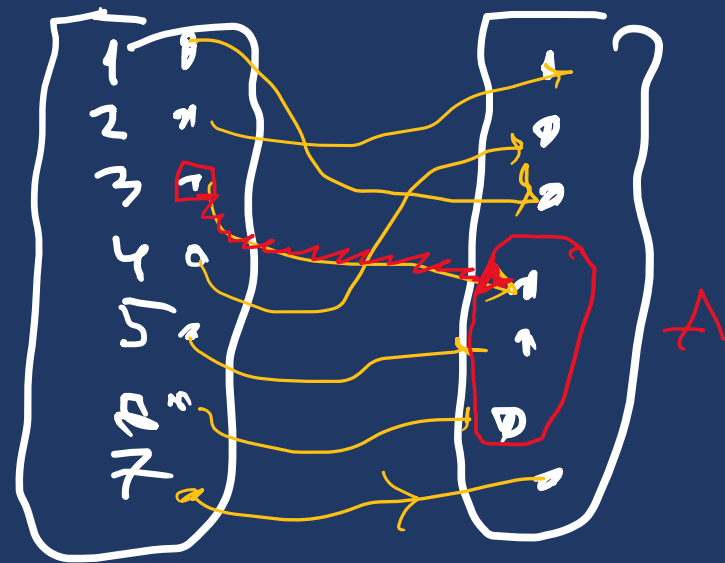
$S_n$  = all permutations of  $\Omega$  ;  $|S_n| = n!$

$\sigma \in S_n$ ,  $A \subseteq \Omega$  ;  $A \neq \emptyset$

$$h_A(\sigma) = \min \{k : \sigma_k \in A\}$$

$$h_A(\sigma) \in \{1, 2, \dots, n\}$$

$$\sigma : \{1 \dots n\} \xrightarrow[\text{output}]{1-1} \{1 \dots n\}$$



Look at  $S_n$  (all perms) as a probab. space:

$$P(\{\sigma\}) = \frac{1}{n!}$$

$$\Pr_{\sigma}(h_A(\sigma) = h_B(\sigma)) = \sum_{k=1}^n \Pr_{\sigma}[h_A(\sigma) = h_B(\sigma) | h_A(\sigma) = k] \cdot \Pr[h_A(\sigma) = k]$$

$$= \sum_{k=1}^n \Pr_{\sigma}[h_B(\sigma) = k | h_A(\sigma) = k] \cdot \Pr[h_A(\sigma) = k]$$

$$= \sum_{k=1}^n \frac{|A \cap B|}{|A \cup B|} \cdot \Pr[h_A(\sigma) = k]$$

$$= \frac{|A \cap B|}{|A \cup B|}$$



$$\{\sigma_1, \dots, \sigma_{k-1}\} \cap (A \cup B) = \emptyset$$

!!!

$$\Pr_{\sigma} [h_A(\sigma) = h_B(\sigma)] = \frac{1}{2^{|A \oplus B|}}$$







