

# **MAC 0459 / 5865**

Topics in Data Science and Engineering

**R. Hirata Jr.** ([hirata@ime.usp.br](mailto:hirata@ime.usp.br))

Class 11 (2020)

# Distance

- a **positive** number that represents the separations between two objects
  - $S$  a set of objects
  - $D : S \times S \rightarrow \mathbb{R}^+$
- $D$  must satisfy four properties:
  - non-negative
  - identity
  - symmetry
  - triangle inequality

# Distance properties

$D : S \times S \rightarrow \mathbb{R}^+$  is a **distance function** (or **metric**) if it satisfies the properties bellow  $a, b, c \in S$ :

- **positive** -  $D(a, b) \geq 0$  for any  $a$  and  $b$
- **identity** -  $D(a, b) = 0$  if and only if  $a = b$
- **simmetry** -  $D(a, b) = D(b, a)$  for any  $a$  and  $b$
- **triangle inequality** - if  $D(a, c) \leq D(a, b) + D(b, c)$  for any  $a, b$  and  $c$

# Distance vs Topology

- Distance is usually defined from  $R^n$  to  $R^+$
- What if the objects “live” in a discrete space?

# Distance vs Topology

## Distance categories

- **Intrinsic** - obeys the topology and curvature of the space of objects
- **Extrinsic** - obeys the topology and curvature of the representation space

# Distance vs Topology

## Distance categories

- **Intrinsic** - obeys the topology and curvature of the space of objects  $d$ : connected digital path between two objects  $\rightarrow \mathbb{Z}^+$
- **Extrinsic** - obeys the topology and curvature of the representation space  $d$ : point coordinates that represent the objects  $\rightarrow \mathbb{R}^+$

# Distance vs Topology

## Tube map



### Check before you travel

† National Rail fares for stations should change at Terminals 1&2 for free rail travel to Terminals 3&4.

‡ Step-Free. Step-Free access for manual wheelchairs only.

§ Wheelchair. Manual train or bus will not stop until next platform.

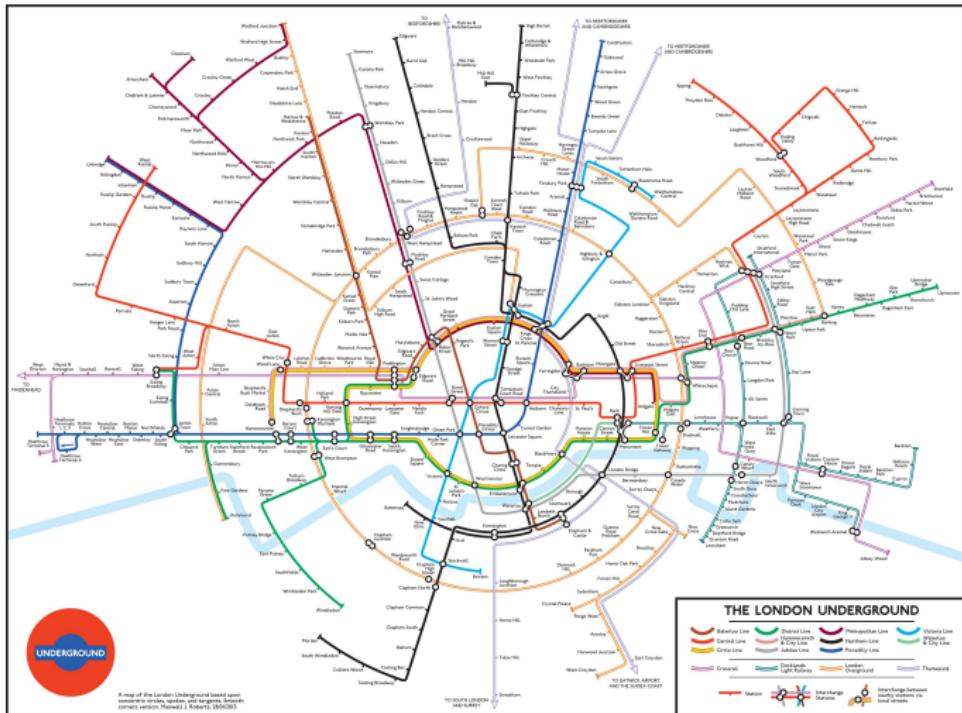
¶ Women. Step-Free access via the Central Line entrance.

\*\* Services or access at these stations are subject to variation.  
Please search TfL stations for full details.

### Key to lines

Bakerloo
Central
Circle
District
Hammersmith & City
Jubilee
Metropolitan
Northern
Piccadilly
Victoria
Waterloo & City

# Distance vs Topology



# Distance vs Topology

Tube map vs real map

[https://www.reddit.com/r/dataisbeautiful/comments/b8ihhr/comparison\\_between\\_the\\_london\\_tube\\_map\\_and\\_its/](https://www.reddit.com/r/dataisbeautiful/comments/b8ihhr/comparison_between_the_london_tube_map_and_its/)

# Distance properties

$D : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is a **distance function** (or **metric**) if it satisfies the properties bellow  $x_i, x_j, x_k \in \mathbb{R}^d$ :

- **positive**, i.e.,  $D(x_i, x_j) \geq 0$  for any  $x_i$  and  $x_j$ .
- **identity**, i.e.,  $D(x_i, x_j) = 0$  if and only if  $x_i = x_j$ ,
- **simmetry**, i.e.,  $D(x_i, x_j) = D(x_j, x_i)$  for any  $x_i$  e  $x_j$ ,
- **triangle inequality**, i.e., if  $D(x_i, x_j) \leq D(x_i, x_k) + D(x_k, x_j)$  for any  $x_i, x_j$  and  $x_k$ ,

# Some distance functions

**Minkowski** or norm  $L_p$  :

$$D(x_i, x_j) = \left[ \sum_{i=k}^d (x_{ik} - x_{jk})^p \right]^{1/p}$$

# Some distance functions

**Minkowski** or norm  $L_p$  :

$$D(x_i, x_j) = \left[ \sum_{i=k}^d (x_{ik} - x_{jk})^p \right]^{1/p}$$

$p = 1$  **city-block, Manhattan**

# Some distance functions

**Minkowski** or norm  $L_p$  :

$$D(x_i, x_j) = \left[ \sum_{i=k}^d (x_{ik} - x_{jk})^p \right]^{1/p}$$

$p = 1$  **city-block, Manhattan**

$p = 2$  **Euclidean**

# Some distance functions

**Minkowski** or norm  $L_p$  :

$$D(x_i, x_j) = \left[ \sum_{i=k}^d (x_{ik} - x_{jk})^p \right]^{1/p}$$

$p = 1$  **city-block, Manhattan**

$p = 2$  **Euclidean**

$p \rightarrow \infty$  **sup-distance (Chebychev distance)**

# Some distance functions

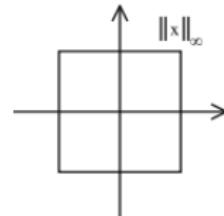
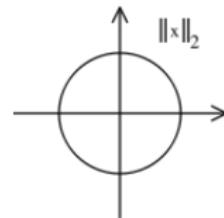
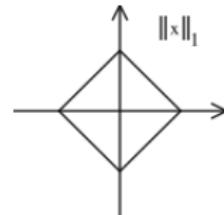
Minkowski or norm  $L_p$  :

$$D(x_i, x_j) = \left[ \sum_{i=k}^d (x_{ik} - x_{jk})^p \right]^{1/p}$$

$p = 1$  city-block, Manhattan

$p = 2$  Euclidean

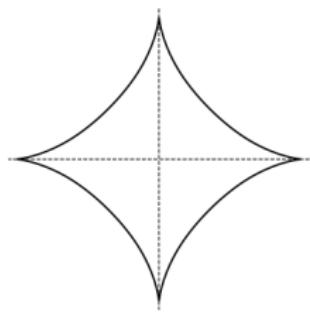
$p \rightarrow \infty$  sup-distance (Chebychev  
distance)



# Some distance functions

Norm  $L_p$

$p = 2/3$  (*subellipse*)

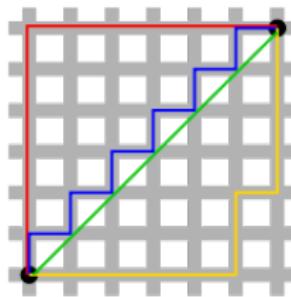
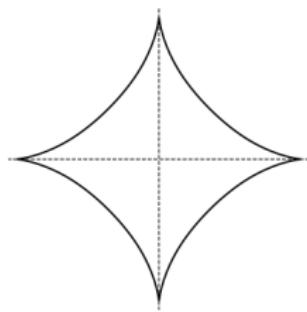


# Some distance functions

Norm  $L_p$

$p = 2/3$  (*subellipse*)

$p = 1$  **city-block, Manhattan**



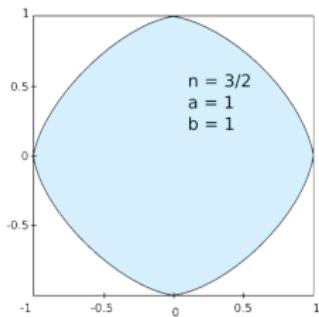
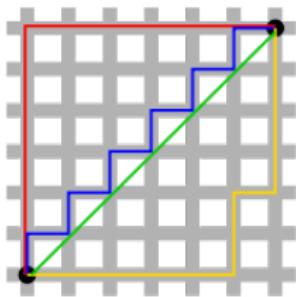
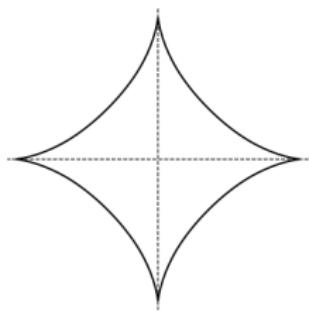
# Some distance functions

Norm  $L_p$

$p = 2/3$  (*subellipse*)

$p = 1$  **city-block, Manhattan**

$p = 3/2$  (*superellipse*)



# Some distance functions

**Camberra** (weighted version of Manhattan):

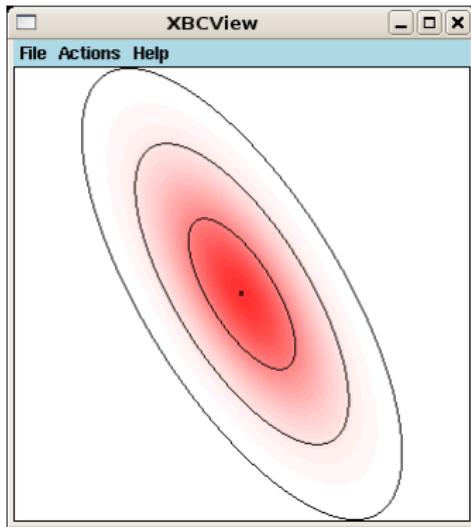
$$D(x_i, x_j) = \sum_{i=k}^d \frac{|x_{ik} - x_{jk}|}{|x_{ik}| + |x_{jk}|}$$

# Some distance functions

## Mahalanobis

$$D(x_i, x_j) = (x_i - x_j)^t \Sigma^{-1} (x_i - x_j)$$

$\Sigma$  symmetric positive definite.



# Similarity functions

A **similarity function**  $S$  satisfies:

- $0 \leq S(x_i, x_j) \leq 1$  **for any**  $x_i$  **and**  $x_j$ ,
- $S(x_i, x_j) = S(x_j, x_i)$  **for any**  $x_i$  **and**  $x_j$ ,
- $S(x_i, x_j) = 1$  **if and only if**  $x_i = x_j$ ,
- $S(x_i, x_j)S(x_j, x_k) \leq [S(x_i, x_j) + S(x_j, x_k)] S(x_i, x_k)$

# Some similarity functions

## Cosine similarity

$$S(x_i, x_j) = \cos\alpha = \frac{x_i \cdot x_j}{\|x_i\| \|x_j\|}$$

## Similarities from distances:

$$S(x_i, x_j) = \frac{1}{1 + D(x_i, x_j)}$$

if  $D(x_i, x_j) = 0$  then  $S(x_i, x_j) = 1$ .

# Distances between objects with categorical characteristics

## Binary data (absent/present characteristic)

$x = (x_1, x_2, \dots, x_d)$  : object with  $d$  characteristics

$x_k = 1$ : characteristic  $k$  present in the object

$x_k = 0$ : characteristic  $k$  absent in the object

$x_i$  and  $x_j$ : any two objects

$n_{11}$ : characteristics with value 1 in both objects

$n_{00}$ : characteristics with value 0 in both objects

$n_{10}$ : characteristics with value 1 in  $x_i$  and value 0 in  $x_j$

$n_{01}$ : characteristics with value 0 in  $x_i$  and value 1 in  $x_j$

$$d = n_{11} + n_{00} + n_{10} + n_{01}$$

# Distances between objects with binary characteristics

Measures that consider some ratio between the matches:

$$\frac{n_{11} + n_{00}}{n_{11} + n_{00} + \omega(n_{10} + n_{01})}$$

$\omega = 1$  simple

$\omega = 2$  Roger & Tanimoto

$\omega = 1/2$  Gower & Legendre

# Distances between objects with binary characteristics

Measures that consider some ratio only between the 1 – 1 matches  
(0 – 0 are not considered):

$$\frac{n_{11}}{n_{11} + \omega(n_{10} + n_{01})}$$

$\omega = 1$       Jaccard index

$\omega = 2$       Sokal & Sneath

$\omega = 1/2$       Gower & Legendre

# Distances between objects with categorical characteristics

A possible similarity measure when the characteristics are not binary is:

$$S(x_i, x_j) = \frac{1}{d} \sum_{l=1}^d S_{ijl}$$

which

$$S_{ijl} = \begin{cases} 0 & \text{if } x_i \text{ and } x_j \text{ do not match in } l, \\ 1 & \text{if } x_i \text{ and } x_j \text{ match in } l. \end{cases}$$

# Value “Normalization”

- Distinct characteristics can have different scales of representation
- This can introduce some bias to the distance/similarity measure
- At the end, to all the analysis

# “Normalization”

## Value “Normalization”

$$\hat{x} = \frac{x - \bar{x}}{\sigma}$$

$\bar{x}$  sample mean

$\sigma$  sample standard deviation

## Linear “normalization” to $[0, 1]$ or $[-1, 1]$

## Non-linear “normalization”

$$\hat{x} = \frac{x - \bar{x}}{r\sigma}$$

$$\hat{x} = \frac{1}{1 + e^{-x}}$$