

The background features four pairs of 3D mathematical manifolds. Each pair consists of a smooth, shaded surface and a corresponding wireframe mesh. The manifolds include a sphere with a grid, a torus with a grid, a complex multi-lobed surface, and a surface with a hole and a protrusion. A horizontal orange line is positioned below the title.

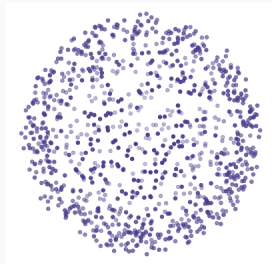
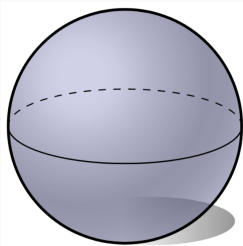
APRENDIZAJE EN VARIETADES PARA EL ANÁLISIS TOPOLÓGICO DE DATOS

XIMENA FERNÁNDEZ

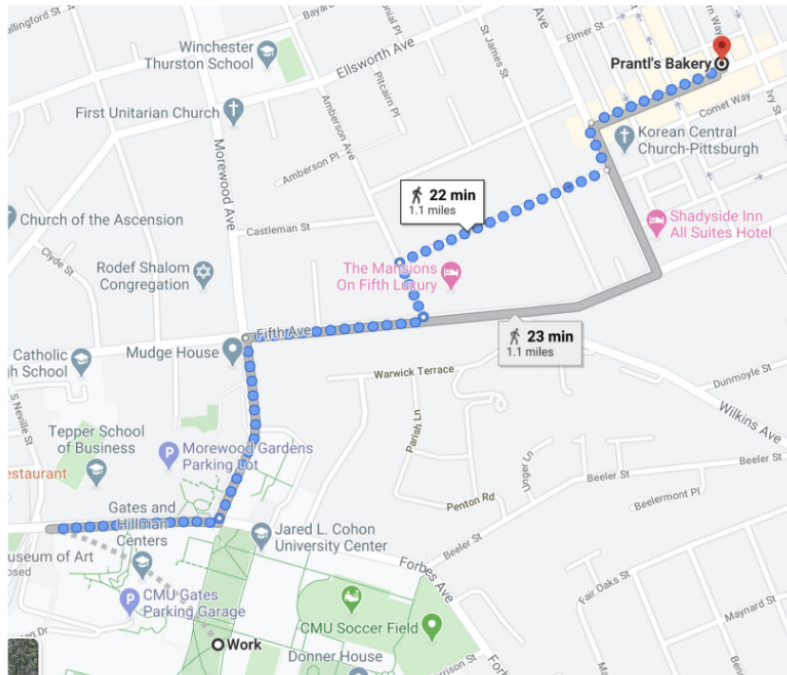
V ENCUENTRO DE JÓVENES TOPÓLOGOS
26 febrero 2021

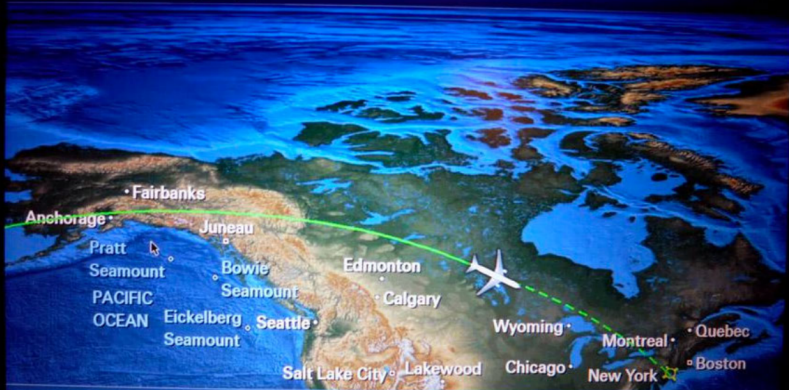
Liverpool-Oxford-Swansea Centre for Topological Data Analysis

Inferir **información** de una **variedad** a partir de una muestra de puntos.

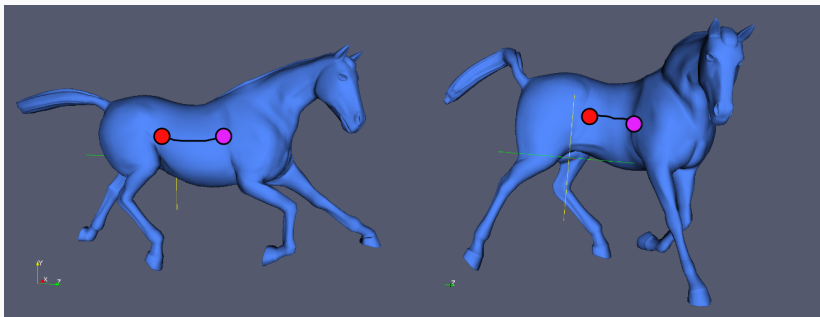


¿Cómo medir distancias?

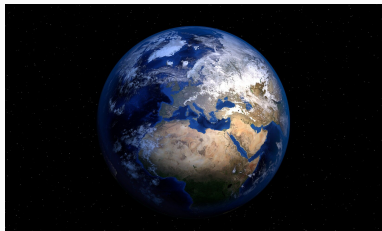




Home



¿Cómo medir distancias *intrínsecas*?



Bernhard Riemann (1826-1866)



Euclides (325 a.C.-265 a.C.)

¿Cómo medir distancias intrínsecas?

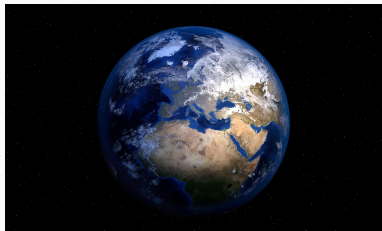
Dados x, y dos puntos de un *espacio topológico* \mathcal{M} .



$\gamma : I \rightarrow \mathcal{M}$ tal que $\gamma(0) = x$, $\gamma(1) = y$

$$d_{\mathcal{M}}(x, y) = \inf_{\gamma} (\text{len}(\gamma))$$

Terraplanistas vs Globistas



Propiedad

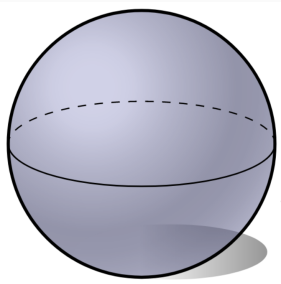
Dado $\delta > 0$, $\exists \epsilon > 0$ tq si $|x - y| < \epsilon$

$$|x - y| \leq d_{\mathcal{M}}(x, y) \leq (1 + \delta)|x - y|$$

Localmente, $d_{\mathcal{M}}$ se puede aproximar por la distancia Euclídea.

¿Cómo *inferir* distancias?

Aprendizaje de distancias



Dado $\epsilon > 0$, se construye el grafo G_ϵ con:

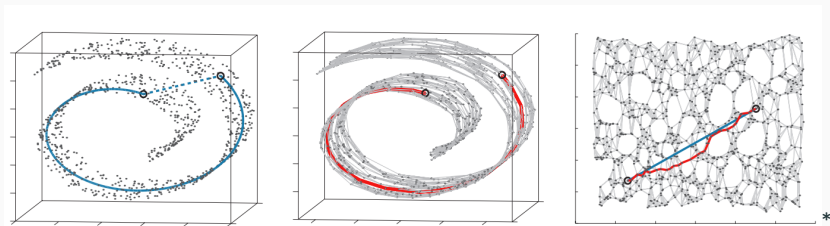
- vértices: $V = \mathbb{X}_n$,
- aristas: $E = \{(x_i, x_j) : |x_i - x_j| < \epsilon\}$.

Se define el **estimador**

$$d_{\mathbb{X}_n, \epsilon}(x, y) = \inf_{\gamma} \sum_{i=0}^r |x_{i+1} - x_i|$$

donde el ínfimo es sobre todos los caminos $\gamma = (x, x_1, \dots, x_r, y)$ con $(x_i, x_{i+1}) \in E(G_\epsilon)$ para todo $0 \leq i \leq r - 1$.

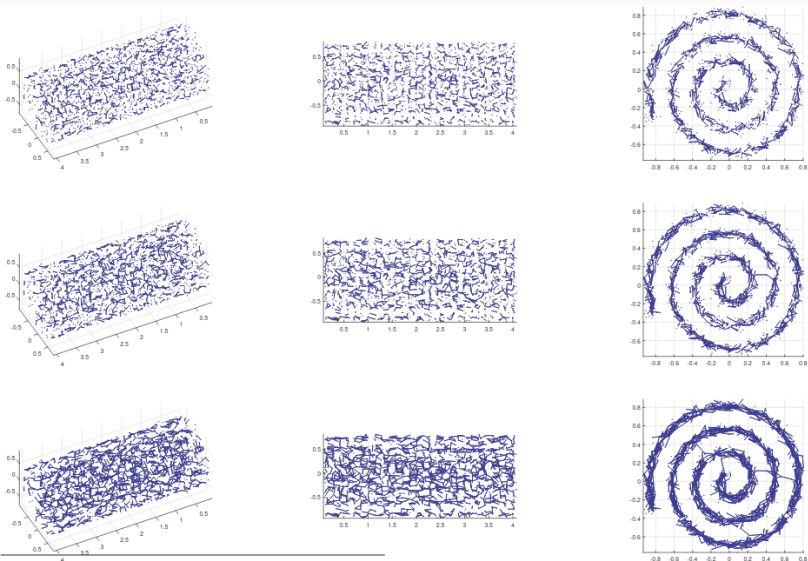
Distancia de ISOMAP



* Imagen: M. Bernstein, V. D. Silva, J. C. Langford, and J. B. Tenenbaum. *Graph approximations to geodesics on embedded manifolds*, 2000.

1. Calcular el grafo G_ϵ (o G_k).
2. Calcular la matriz de distancias.
3. Calcular el embedding.

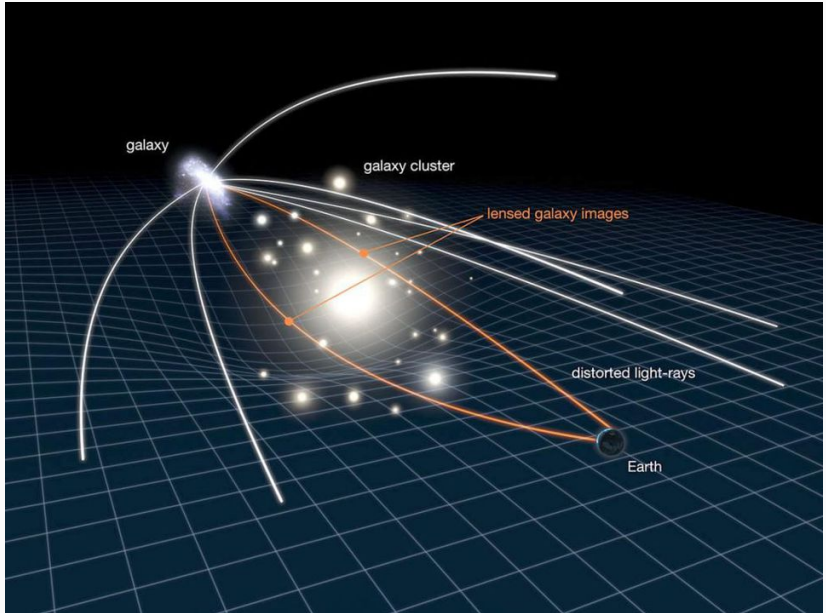
ISOMAP: Debilidades



† Imagen: B. Mahler. *Contagion Dynamics for Manifold Learning*, 2020.

**¿Cómo medir distancias
basadas en densidad?**

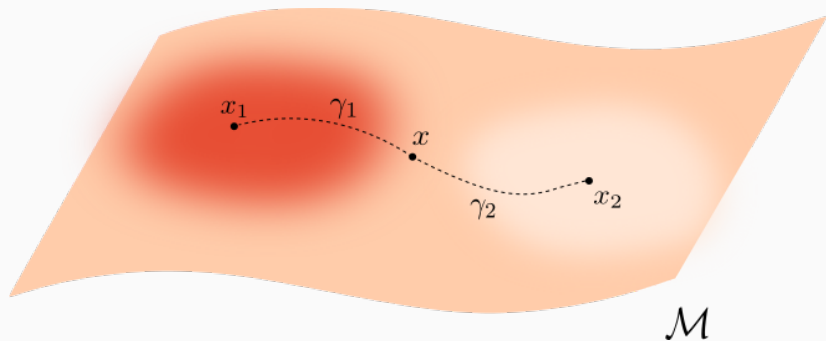
Distancias basadas en densidad



Distancias basadas en densidad



Distancias basadas en densidad



Distancia de Fermat

Sea \mathcal{M} una **variedad** suave de dimensión d embebida en \mathbb{R}^D .

Sea $f : \mathcal{M} \rightarrow \mathbb{R}_{>0}$ una **función de densidad**.

Distancia de Fermat

Sea \mathcal{M} una **variedad** suave de dimensión d embebida en \mathbb{R}^D .

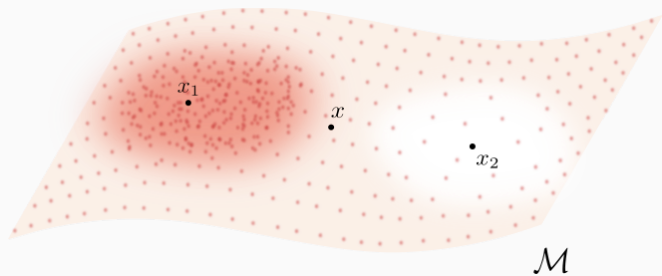
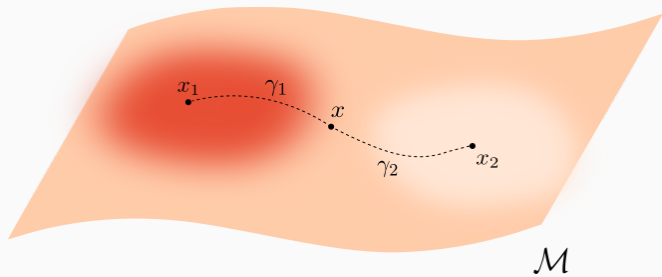
Sea $f : \mathcal{M} \rightarrow \mathbb{R}_{>0}$ una **función de densidad**.

Dado $p > 1$, se define la **distancia deformada**

$$d_{f,p}(x, y) = \inf_{\gamma} \int_I \frac{1}{f(\gamma(t))^{(p-1)/d}} \|\gamma'(t)\| dt.$$

$d_{f,p}$ se llama la **p -distancia de Fermat** en analogía con el Principio de Fermat en óptica.

Aprendizaje de distancias basadas en densidad

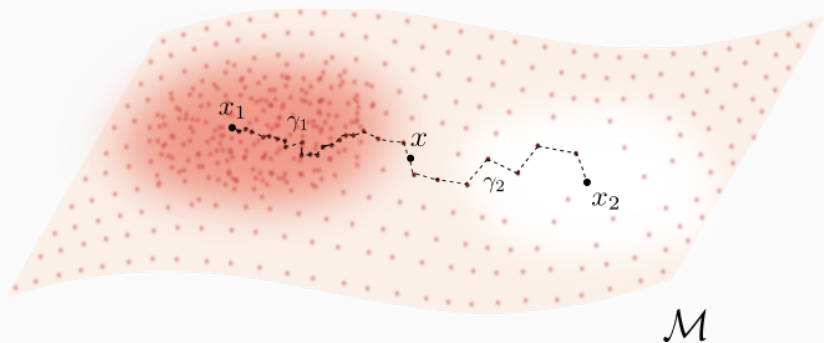


Distancia de Fermat muestral

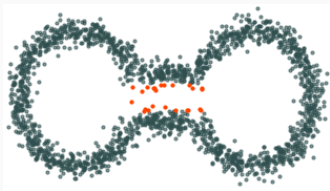
Dado $p > 1$, la **distancia de Fermat muestral** entre x, y se define por

$$d_{\mathbb{X}_n, p}(x, y) = \inf_{\gamma} \sum_{i=0}^r |x_{i+1} - x_i|^p$$

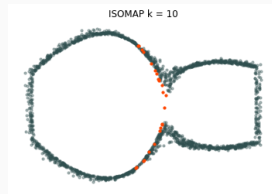
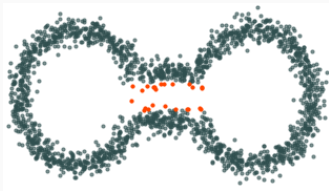
donde el ínfimo es sobre todos los caminos $\gamma = (x_0, \dots, x_{r+1})$ con $x_0 = x$, $x_{r+1} = y$ y $\{x_1, x_2, \dots, x_r\} \subseteq \mathbb{X}_n$.



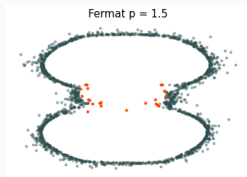
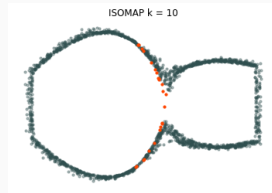
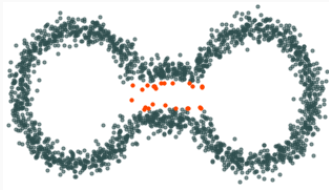
Ejemplo



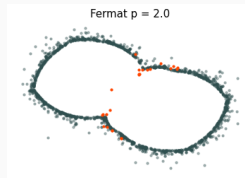
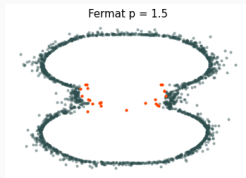
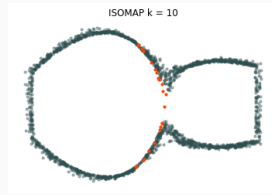
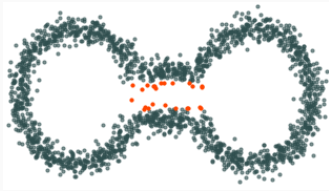
Ejemplo



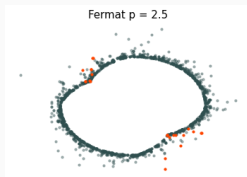
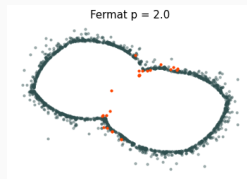
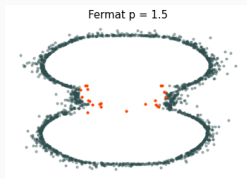
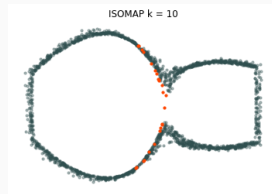
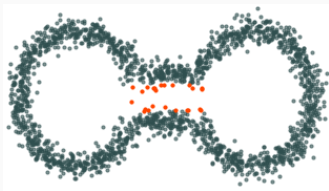
Ejemplo



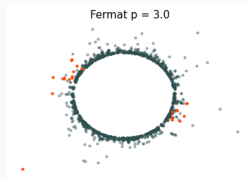
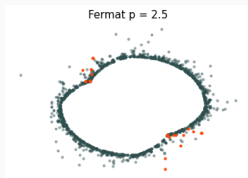
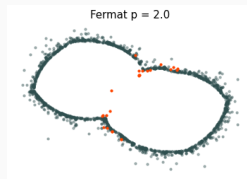
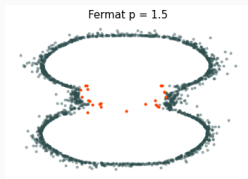
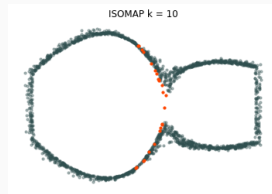
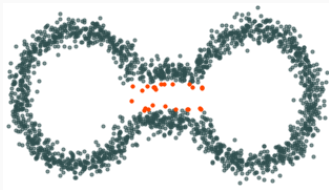
Ejemplo



Ejemplo



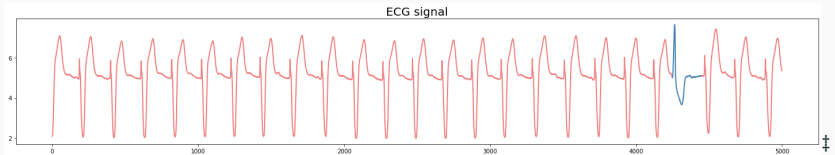
Ejemplo



Aplicaciones

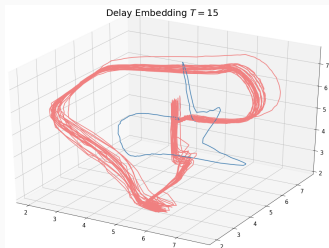
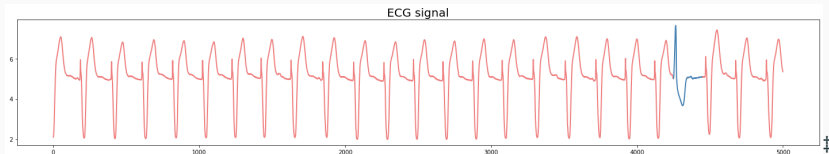
Detección de anomalías en series temporales

Señal de electrocardiograma con un latido anómalo (arritmia).



Detección de anomalías en series temporales

Señal de electrocardiograma con un latido anómalo (arritmia).



Detección de anomalías en series temporales

Señal de electrocardiograma con un latido anómalo (arritmia).

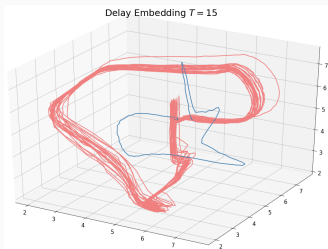
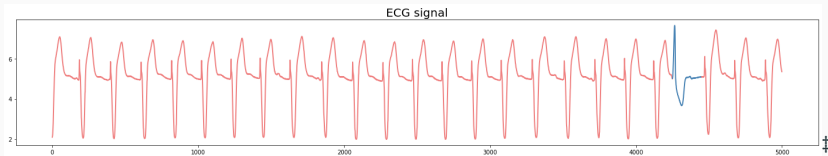
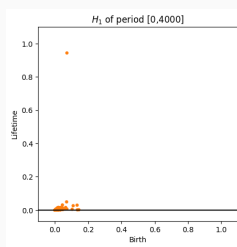


Diagrama de persistencia con distancia de Fermat $\rho = 2$.



‡Datos de la base de datos Physionet, MIT Laboratorio de Fisiología Computacional.

Detección de anomalías en series temporales

Señal de electrocardiograma con un latido anómalo (arritmia).

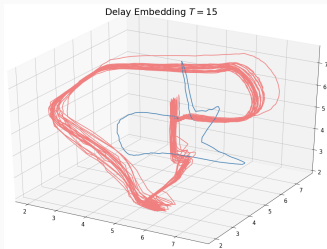
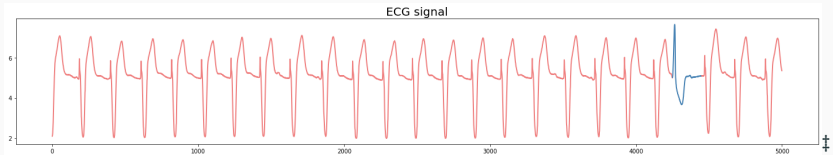
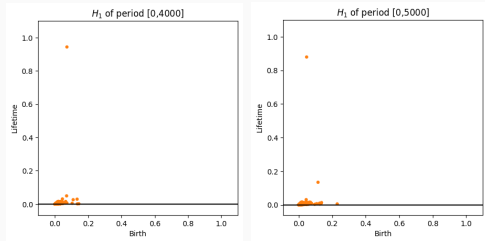


Diagrama de persistencia con distancia de Fermat $\rho = 2$.



‡Datos de la base de datos Physionet, MIT Laboratorio de Fisiología Computacional.

Referencias

- J. B. Tenenbaum. *A Global Geometric Framework for Non-linear Dimensionality Reduction* (2000).
- M. Bernstein, V. D. Silva, J. C. Langford, J. B. Tenenbaum. *Graph approximations to geodesics on embedded manifolds* (2000).
- E. Borghini, X. Fernández, P. Groisman, G. Mindlin. *Intrinsic persistent homology via density-based distance learning*. arXiv:2012.07621 (2020).

web: <https://ximenafernandez.github.io/>

github: ximenafernandez

email: x.l.fernandez@swansea.ac.uk

GRACIAS!