

**CLASIFICACIÓN DE LAS MEDIDAS DE DISIMILARIDAD  $d$  ENTRE  
CONGLOMERADOS DE DATOS SEGÚN LA ESTRUCTURA MÉTRICA DEL  
ESPACIO MUESTRAL DE APLICACIÓN**

**ISADORE NABI**

<b>Espacio Muestral de <math>d</math></b>	<b>Estructura Métrica del Espacio Muestral de <math>d</math></b>	<b>Medida de Disimilaridad <math>d</math></b>	<b>Expresiones Funcionales de <math>d</math></b>
Espacios Métricos Completos (Espacios de Cauchy)	Según (Chiang, 2006, pág. 65), a) $d(x, y) = 0$ , para $x = y$ <sup>1</sup> b) $d(x, y) = d(y, x) > 0$ , para $x \neq y$ c) $d(x, y) \leq d(x, z) + d(z, y)$ , para $z \neq x, y$ <sup>2</sup>	<i>Distancia Euclidiana</i>	En $\mathbb{R}^2$ , según (Weisstein, Euclidian Metric, 2022): $d(x, y) = \sqrt{(x_1 - y_1)^2 + \dots + (x_n - y_n)^2}$ , con $x = x_1, \dots, x_n$ y $y = y_1, \dots, y_n$ .
	d) $d(x, y) \geq 0$ , positivada		<i>Coefficiente de Correlación de Pearson</i>

<sup>1</sup> También conocida como *identidad de los indiscernibles*.

<sup>2</sup> También conocida como desigualdad triangular, la cual no es más que una generalización del Teorema de Pitágoras para cualesquiera tipos de rectángulos.

	d, deducida de las tres propiedades anteriores.	<p><i>Coefficiente de Correlación de Spearman</i></p>	<p>Según (Barcelona Field Studies Center, 2022), cuando no hay rangos empatados (del mismo tamaño):</p> $d(x, y) = 1 - \left( \frac{6 \sum_{i=1}^n (X_i - Y_i)^2}{n^3 - n} \right)$ <p>Según (Laerd Statistics, 2022), cuando hay rangos empatados:</p> $d(x, y) = \frac{\sum_{i=1}^n (X_i - \bar{X}) \sum_{i=1}^n (Y_i - \bar{Y})}{\sqrt{(\sum_{i=1}^n (X_i - \bar{X})^2) (\sum_{i=1}^n (Y_i - \bar{Y})^2)}}$ <p>donde <math>\bar{X}</math> y <math>\bar{Y}</math> son las medias muestrales de <math>X</math> e <math>Y</math>, respectivamente.</p>
		<p><i>Coefficiente de Correlación de Kendall</i></p>	<p>Según (Statology, 2020):</p> $d(x, y) = \frac{C - D}{C + D}$ <p>donde <math>C</math> es el número de pares concordantes y <math>D</math> es el número de pares discordantes<sup>3</sup>.</p>

<sup>3</sup> Como se señala en (StatisticsHowTo, 2016), supóngase que dos entrevistadores calificaron a un grupo de doce solicitantes de empleo, datos recopilados en la siguiente tabla:

Espacios Semi- Métricos	Según (McAuley, 1956, pág. 315): a) $d(x, y) = 0$ , si $x = y$ b) La topología del espacio topológico $S$ es invariante con	Cuasi-Métrica	a) Según (Stojmirovic, 2005, pág. 23), las cuasi métricas en $\mathbb{R}$ que generan las denominadas topología superior ( $d^R$ ) e inferior ( $d^L$ ) son: $d^L(x, y) = u^L = \max\{x - y, 0\}$ $d^R(x, y) = u^R = \max\{y - x, 0\}$ En la expresión anterior, tanto $d^L$ como $d^R$ son mapas de la forma $\mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+$ . Además,
-------------------------------	---	---------------	--

Candidate	Interviewer 1	Interviewer 2
A	1	1
B	2	2
C	3	4
D	4	3
E	5	6
F	6	5
G	7	8
H	8	7
I	9	10
J	10	9
K	11	12
L	12	11

Teniéndose en cuenta que, en la primera columna, las opciones del entrevistador 1 se ordenaron de menor a mayor, se puede hacer una comparación entre las opciones para el entrevistador 1 y 2. Con pares concordantes o discordantes, básicamente está respondiendo a la pregunta: ¿los jueces/evaluadores clasificaron los pares en el mismo orden? No necesariamente está buscando exactamente el mismo rango, sino más bien si un solicitante de empleo fue clasificado constantemente más alto por ambos entrevistadores. Tres posibles escenarios son posibles para estos pares ordenados: 1. pares empatados: ambos entrevistadores están de acuerdo. Por ejemplo, el candidato A se marcó como primera opción para ambos entrevistadores, por lo que están empatados. 2. Pares concordantes: ambos entrevistadores clasifican a ambos solicitantes en el mismo orden, es decir, ambos se mueven en la misma dirección. Si bien no tienen el mismo rango (es decir, ambos 1° o 2°), cada par se ordena igual de alto o igual de bajo. El entrevistador 1 clasificó a F como 6 y G como 7, mientras que el entrevistador 2 clasificó a F como 5 y G como 8. F y G son concordantes porque F siempre se clasificó más alto que G. 3. Pares discordantes: los candidatos E y F son discordantes porque los entrevistadores clasificaron en direcciones opuestas (uno dijo que E tenía una clasificación más alta que F, mientras que el otro dijo que F tenía una clasificación superior a 6).

	<p>respecto a la función distancia <math>d^4</math>.</p> <p>c) No se cumple que <math>d(x, y) = d(y, x)</math>, aunque las distancias siempre positivas cuando <math>x \neq y</math>.</p>		<p><math>d^L</math> es el conjugado de <math>d^R</math> y viceversa.</p> <p>b) Según (Stojmirovic, 2005, pág. 22), si <math>X</math> es un conjunto de datos y existe una medida <math>d</math> que describe un mapa de la forma <math>d: X \times X \rightarrow \mathbb{R}</math>, entonces la denominada <i>métrica discreta</i> tiene la forma funcional:</p> $d(x, y) = \begin{cases} 0, & \text{si } x = y \\ 1, & \text{si } x \neq y \end{cases}$ <p>c) Según (Stojmirovic, 2005, pág. 24), otra cuasi métrica en <math>\mathbb{R}_+</math> es:</p> $d(x, y) = \begin{cases} \min(1, y - x), & \text{si } x \leq y \\ 1, & \text{de otra forma} \end{cases}$
<i>Espacios Pseudo-Métricos</i>	Posee casi la misma estructura definida para los espacios de		Según (Planethmath, 2013), los siguientes son ejemplos de medidas cuasi-métricas:

<sup>4</sup> Es decir, si existe un punto  $p$  que es punto límite (una observación  $x$  será punto límite de un conjunto si cumple que para todo número arbitrariamente mayor que cero  $\epsilon > 0$  existe un miembro del conjunto  $y$  diferente de  $x$  tal que  $|y - x| < \epsilon$  (Weisstein, Limit Point, 2022); son conocidos también como puntos clúster o puntos de acumulación) de un subconjunto  $M \in S$ , entonces  $p$  es un punto límite distancia de  $M$  (i.e., un punto límite que puede ser considerado para calcular la distancia a  $M$ ) y viceversa. Lo anterior significa que la función distancia  $d$  es un invariante topológico (una función, en cuanto invariante topológico, es aquella función que para aplicarse depende únicamente de la topología del espacio, por lo que es aplicable a cualquier espacio homeomórfico, es decir, a cualquier espacio con la misma topología, con independencia de las consideraciones métricas) en el espacio muestral.

	<p>Cauchy, con la única diferencia que existe la posibilidad de <math>d(x, y) = 0</math> aun cuando <math>x \neq y</math>.</p>	<p><i>Pseudo-Métrica</i></p>	<p>a) Sea <math>X = \mathbb{R}^2</math> y considérese una función <math>d: X \times X \rightarrow \mathbb{R}_+</math> dada por <math>d[(x_1, x_2), (y_1, y_2)] =  x_1 - y_1 </math>, entonces <math>d(x, x) =  x_1 - x_1  = 0</math>, <math>d(x, y) =  x_1 - y_1  =  y_1 - x_1  = d(y, z)</math> y la desigualdad triangular se deriva naturalmente de la desigualdad triangular en <math>\mathbb{R}^1</math>, por lo que <math>(X, d)</math>, donde <math>d</math> es la medida de disimilaridad, satisface las condiciones de un espacio pseudo-métrico. En este espacio puede obtenerse una distancia nula entre dos observaciones diferentes, por ejemplo, <math>d[(2, 3), (2, 5)] =  2 - 2  = 0</math></p> <p>b) Sea <math>X</math> un conjunto, sea <math>x_0 \in X</math> y sea <math>F(X)</math> un conjunto de funciones tal que <math>X \rightarrow R</math>. Entonces, <math>d(f, g) =  f(x_0) - g(x_0) </math> es una pseudo-métrica de <math>F(X)</math>.</p> <p>c) Si <math>X</math> es un espacio vectorial y <math>p</math> es una</p>
--	--	------------------------------	--

			<p>semi-norma<sup>5</sup> sobre el espacio <math>X</math>, entonces <math>d(x, y) = p(x - y)</math> es una pseudo-métrica sobre <math>X</math>.</p> <p>d) La <i>pseudo-métrica trivial</i> <math>d(x, y) = 0</math> para todo <math>x, y \in X</math> es una pseudo-métrica.</p>
<p><i>Variedades Pseudo-Riemannianas</i><sup>6</sup></p>	<p>Una variedad <math>n - dimensional</math> es un espacio topológico con la propiedad de que cada observación dentro de dicho espacio tiene un vecindario (un entorno formado partiendo de un centro y una distancia fija) que</p>	<p><i>Métrica de Minkowski</i></p>	<p>Según (Kolmogórov &amp; Fomin, 1978, pág. 56):</p> $\left( \sum_{k=1}^n  x_k + y_k ^p \right)^{\frac{1}{p}} \leq \left( \sum_{k=1}^n  x_k ^p \right)^{\frac{1}{p}} + \left( \sum_{k=1}^n  y_k ^p \right)^{\frac{1}{p}},$ <p>donde <math>p</math> expresa el orden de la potencia a la cual las funciones que conforman el espacio son integrables. Un espacio <math>L_p</math> es aquel espacio de funciones <math>f</math> (donde cada coordenada u observación puede ser representada mediante una función) conformado por aquellas <math>f</math> que, elevadas a una</p>

<sup>5</sup> Como señala (Weisstein, Seminorm, 2022), es una función en un espacio vectorial  $V$  denotada como  $\|v\|$ , tal que cumple las siguientes condiciones para todos los vectores  $v$  y  $w$  que pertenecen a  $V$ , así como para todo escalar  $c$ :

1.  $\|v\| \geq 0$ ,
2.  $\|cv\| = |c|\|v\|$ ,
3.  $\|v + w\| \leq \|v\| + \|w\|$ .

Nótese que es posible que  $\|v\| = 0$  para vectores  $v$  no nulos.

<sup>6</sup> Es una variedad con un tensor métrico no-degenerado en todas partes, lo cual se explicará a continuación.

	<p>es homeomórfico (topológicamente equivalente) a algún conjunto abierto de un espacio euclidiano <math>n - dimensional^7</math>. Una variedad pseudo-riemmaniana posee casi la misma estructura definida para las variedades, con la salvedad de que además es posible obtener un tensor métrico<sup>8</sup> no-</p>		<p>potencia <math>p</math>, son integrables. Por ejemplo, los espacios <math>L_2</math> son los espacios de funciones que son cuadrado-integrables, es decir, que la integral de <math>f^2</math> existe (<i>i.e.</i>, tiene un valor finito).</p>
--	--	--	--

<sup>7</sup> Que se diferencia de un espacio métrico (o de un espacio métrico completo) en que el espacio métrico completo es la conjunción del espacio euclidiano y de la métrica euclidiana. Véase (Kolmogórov & Fomin, 1978, pág. 52).

<sup>8</sup> Como se señala en (Wikipedia, 2021), un tensor métrico es un mapa bilineal (función lineal en cada uno de sus argumentos que combina elementos de dos espacios vectoriales -que pueden definirse como matrices de datos- y producen un elemento de un tercer espacio vectorial -que puede definirse también como una matriz de datos-) no degenerado (es decir, el mapa bilineal y su mapa conjugado son isomórficos, *i.e.*, topológicamente equivalentes), suave (sus  $n - ésimas$  derivadas evaluados en un punto existen) y simétrico (la distancia de  $x$  a  $y$  es igual a la distancia de  $y$  a  $x$ ) que asigna un número real a un par de vectores tangentes (relacionados entre sí a través de una estructura diferencial -de derivadas- que se mostrará al definir la estructura de la medida de disimilitud de las variedades pseudo-riemannianas) en cada espacio tangente (como señala (Weisstein, Tangent Space, 2022), es una copia de un espacio  $\mathbb{R}^n$  tangencial a una variedad compacta  $M$ , es decir, a una variedad que es compacta en tanto espacio topológico; un espacio

	degenerado En Todas Partes <sup>9</sup>		
--	--	--	--

## REFERENCIAS

- Barcelona Field Studies Center. (21 de Abril de 2022). *Spearman's Rank Correlation Coefficient*. Obtenido de Spearman's Rank:  
<https://geographyfieldwork.com/SpearmansRank.htm>
- Chiang, A. (2006). *Métodos Fundamentales de Economía Matemática*. México, D.F.: McGraw-Hill. Obtenido de  
<https://elvisjgblog.files.wordpress.com/2018/02/mc3a9todos-fundamentales-de-economc3ada-matemc3a1tica-4ta-edicic3b3n-alpha-c-chiang-freelibros-org.pdf>
- Gujarati, D. N., & Porter, D. C. (2010). *Econometría*. México D.F.: McGraw-Hill Educación. Obtenido de  
<https://fvela.files.wordpress.com/2012/10/econometria-damodar-n-gujarati-5ta-ed.pdf>
- Halabiski, A. (23 de Febrero de 2022). *Euclidean Distance In 'n'-Dimensional Space*. Obtenido de Stanford University:  
[https://hlab.stanford.edu/brian/euclidean\\_distance\\_in.html](https://hlab.stanford.edu/brian/euclidean_distance_in.html)
- Kaplan, W. (1985). *CÁLCULO AVANZADO*. MÉXICO, D.F.: COMPAÑÍA EDITORIAL CONTINENTAL, S.A. DE C.V., MÉXICO.
- Kolmogórov, A. N., & Fomin, S. V. (1978). *Elementos de la Teoría de Funciones y del Análisis Funcional* (Tercera ed.). (q. e.-m. Traducido del ruso por Carlos Vega, Trad.) Moscú: MIR.

---

topológico es compacto si todo cubrimiento abierto -colección de conjuntos abiertos de un espacio topológico cuya unión contiene algún conjunto dado- posee un subcubrimiento finito -un cubrimiento que es subconjunto de otro cubrimiento-) de la variedad generada por el conjunto de datos. Como señala (Kaplan, 1985, pág. 305), un tensor de orden dos (bidimensional) posee la forma  $g_{ij} = \frac{\partial \xi^r}{\partial x^i} \frac{\partial \xi^r}{\partial x^j}$ , donde  $\xi$  expresa el sistema de coordenadas estándar (del cual se generan todos los nuevos sistemas de coordenadas),  $x^i$  representa la  $i$  – ésima variable,  $x^j$  la  $j$  – ésima variable de la matriz de datos (puesto que una matriz bidimensional como las utilizadas en el álgebra lineal son equivalentes a tensores de orden 2) y el símbolo  $\partial$  denota la derivada parcial del numerador respecto del denominador.

<sup>9</sup> Es decir, en cualquier localidad (subregión) del espacio en que se tomen las medidas de disimilaridad.



- Laerd Statistics. (1 de Febrero de 2022). *Spearman's Rank-Order Correlation*. Obtenido de Statistical Guides: <https://statistics.laerd.com/statistical-guides/spearmans-rank-order-correlation-statistical-guide.php>
- McAuley, L. F. (Diciembre de 1956). A Relation Between Perfect Separability, Completeness, and Normality in Semi-Metric Spaces. *Pacific Journal of Mathematics*, VI(2), 315-326.
- Planetmath. (22 de Marzo de 2013). *example of pseudometric space*. Obtenido de planetmath.org: <https://planetmath.org/exampleofpseudometricspace>
- StatisticsHowTo. (20 de Agosto de 2016). *Concordant Pairs and Discordant Pairs*. Obtenido de Statistics Definitions: <https://www.statisticshowto.com/concordant-pairs-discordant-pairs/>
- Statology. (26 de Febrero de 2020). *Kendall's Tau: Definition + Example*. Obtenido de Statistics. Simplified: <https://www.statology.org/kendalls-tau/>
- Stojmirovic, A. (2005). *Quasi-metrics, Similarities and Searches: aspects of geometry of protein datasets*. Wellington: Victoria University of Wellington. Obtenido de <https://arxiv.org/pdf/0810.5407.pdf>
- Weisstein, E. W. (19 de Abril de 2022). *Euclidian Metric*. Obtenido de MathWorld - A Wolfram Web Resource: <https://mathworld.wolfram.com/EuclideanMetric.html>
- Weisstein, E. W. (19 de Abril de 2022). *Limit Point*. Obtenido de MathWorld - A Wolfram Web Resource: <https://mathworld.wolfram.com/LimitPoint.html>
- Weisstein, E. W. (19 de Abril de 2022). *Seminorm*. Obtenido de MathWorld - A Wolfram Web Resource: <https://mathworld.wolfram.com/Seminorm.html>
- Weisstein, E. W. (19 de Abril de 2022). *Tangent Space*. Obtenido de MathWorld - A Wolfram Web Resource: <https://mathworld.wolfram.com/TangentSpace.html>
- Wikipedia. (10 de Noviembre de 2021). *Pseudo-Riemannian manifold*. Obtenido de Smooth Manifolds: [https://en.wikipedia.org/wiki/Pseudo-Riemannian\\_manifold](https://en.wikipedia.org/wiki/Pseudo-Riemannian_manifold)