

Proyecto – Entrega final

Grupo 13: Freddy Rodrigo Mendoza Ticona, William Alexander Romero Bolívar, Maria Paula Salamanca Delgado, Jorge Oswaldo Suárez Rodríguez.

Detección e identificación de clústeres por afectación de dengue en Colombia: una aproximación desde el Aprendizaje no Supervisado

Resumen

El dengue es una enfermedad viral transmitida por la picadura de mosquitos infectados que se caracteriza por producir fiebre, dolor corporal, pérdida del apetito y, en casos graves, sangrado de mucosas. La población más vulnerable a esta enfermedad son niños y adultos mayores, sin embargo, puede afectar a cualquier grupo demográfico. Según el Instituto Nacional de Salud (INS), a julio de 2022, en Colombia se han presentado 34.017 casos, dentro de los cuales el 52,4% presentaron signos de alarma o graves (1,2).

Es posible prevenir y controlar su propagación concientizando a la población para evitar la proliferación de mosquitos que transmiten esta enfermedad. Así, a través de un análisis de clústeres se estratificaron 1.120 municipios de Colombia según el riesgo de incidencia de dengue de acuerdo con sus características demográficas (edad, población), socioeconómicas (desempleo, estratos), y climáticas (temperatura, precipitaciones), encontrando que existe un mayor riesgo de tasa de incidencia de dengue en municipios de clima cálido.

Los resultados de este proyecto permitirían dirigir recursos hacia la capacitación y concientización de la población para la prevención y control del dengue eficientemente en aquellas zonas de mayor riesgo de incidencia del dengue.

Introducción

Según la Organización Mundial de la Salud (OMS), se estima que anualmente el número de personas infectadas por dengue llega a 390 millones en el mundo y en riesgo de infección, a 3,9 billones (3). Adicionalmente, según la Organización Panamericana de la Salud (OPS), la circulación del dengue en Suramérica es hiperendémica, siendo Brasil y Colombia los que reportan el mayor número de casos; tan solo para el 2019 Colombia reportó 127.000 casos (4).

Los mosquitos *aegypti* son el principal vector de transmisión del dengue y están presentes en casi todos los países de Suramérica. El ciclo de transmisión se da por el contagio de la sangre de un humano infectado a un humano susceptible a través de la picadura del mosquito *Aedes* hembra que se reproduce cerca de las casas, poniendo huevos en recipientes de agua estancada. Se ha revisado factores asociados al aumento de la incidencia de infección por dengue, como el aumento global de la temperatura que aumenta el tiempo que el mosquito permanece infeccioso; el aumento estacional de las precipitaciones que contribuye a una mayor densidad de mosquitos; y condiciones de hacinamiento que incrementan la fracción de la población susceptible (5,6).

En estudios anteriores sobre el dengue se han utilizado algoritmos de aprendizaje no supervisado como los métodos de agrupamiento que son adecuados para la visualización de enfermedades, en especial los basados en la densidad para separar las regiones que tienen alta y baja densidad. En el estudio de Shaukat, et al, el análisis de datos de Pakistán por DBSCAN encontró que el dengue atacó principalmente la ciudad de Jhelum y Tehsil Jhelum, lo que permite enfocar estrategias de prevención de la enfermedad (7). Otro estudio realizado en la India encontró conglomerados de casos de dengue en Delhi usando un algoritmo de agrupamiento DBSCAN; estos puntos críticos se caracterizaron por ser grupos socioeconómicos bajos, estar cerca al río Yamuna, lagos y lugares con agua estancada (8).

Una revisión de la literatura realizada en Colombia encontró que, a nivel mundial en general, se ha usado modelos de aprendizaje supervisado con el objetivo de predicción de dengue, como de regresión logística con datos demográficos, clínicos y de laboratorio, así como modelos de regresión lineal, Random Forest y Support Vector Machine, con datos socioeconómicos, demográficos, climáticos y ambientales (9). También, estudios de análisis espacio-temporal como el realizado en Cali, Colombia, en el cual se encontró que el nivel socioeconómico, la densidad poblacional, proximidad a talleres con neumáticos, viveros de plantas y sistemas de alcantarillado, están relacionados con la enfermedad (10).

Sin embargo, no se encontró literatura donde se utilice el aprendizaje no supervisado para la estratificación del riesgo de dengue en Colombia, por lo cual, un análisis de clústeres por municipios de Colombia aportaría al caso de estudio. Así, este trabajo tiene como propósito estratificar el riesgo de infección por localización mediante una unidad de medida, como la tasa de incidencia de dengue del análisis de clústeres. Lo anterior sería de utilidad para las entidades gubernamentales y funcionarios de salud pública, al poder diferenciar zonas con mayor riesgo que requieran toma de decisiones sobre una intervención de control de propagación de esta enfermedad.

Metodología

Para el desarrollo de este proyecto se cuenta con una base de datos con información sobre 1.121 municipios y 1.017 características, como datos demográficos y socioeconómicos, casos acumulados de dengue del 2007 al 2019 y por semana epidemiológica e información climática como mediciones de temperatura y precipitaciones mensuales. Esta información se obtuvo de las bases de datos del Sistema de Vigilancia en Salud Pública (SIVIGILA) del Instituto Nacional de Salud (INS) y del Departamento Administrativo Nacional de Estadística (DANE).

Variable	Descripción	Tipo dato
Municipality code	Código del municipio	numérico
Municipality	Descripción de municipio	texto
Population (2007-2019)	Población del municipio por año	numérico
Cases (2007-2019)	Casos dengue reportados en el municipio x año	numérico
Age0-4(%)	Porcentaje de población menor de 4 años	numérico
Age5-14(%)	Porcentaje de población entre 5 y 14 años	numérico
Age15-29(%)	Porcentaje de población entre 15 y 29 años	numérico
Age>30(%)	Porcentaje de población mayor de 30 años	numérico
AfrocolombianPopulation(%)	Porcentaje de población afrocolombiana	numérico
IndianPopulation(%)	Porcentaje de población indígena	numérico
PeoplewithDisabilities(%)	Porcentaje de población con discapacidades (física, psicológica o mental)	numérico
Peoplewhocannotreadorwrite(%)	Porcentaje de población con que no puede leer/escribir	numérico
Secondary/HigherEducation(%)	Porcentaje de población que tiene educación secundaria	numérico
Employedpopulation(%)	Porcentaje de población empleada	numérico
Unemployedpopulation(%)	Porcentaje de población desempleada	numérico
Peopledoinghousework(%)	Porcentaje de población que realizan trabajo doméstico	numérico
Retiredpeople(%)	Porcentaje de población jubilada	numérico
Men(%)	Porcentaje de población masculina	numérico
Women(%)	Porcentaje de población femenina	numérico
Householdswithoutwateraccess (%)	Porcentaje de viviendas sin acceso a agua	numérico
Householdswithoutinternetaccess(%)	Porcentaje de viviendas sin acceso a internet	numérico
Buildingstratification1(%)	Porcentaje de viviendas estrato 1	numérico
Buildingstratification2(%)	Porcentaje de viviendas estrato 2	numérico
Buildingstratification3(%)	Porcentaje de viviendas estrato 3	numérico
Buildingstratification4(%)	Porcentaje de viviendas estrato 4	numérico
Buildingstratification5(%)	Porcentaje de viviendas estrato 5	numérico

A continuación, se encuentran algunas estadísticas de estas variables:

	2007/w1	2007/w2	2007/w3	2007/w4	2007/w5	2007/w6	2007/w7	2007/w8	2007/w9	2007/w10	...	2019/w3	2019/w4	2019/w5	2019/w6	2019/w7	2019/w8	2019/w9	2019/w10	2019/w11	2019/w12
count	1130	1130	1130	1130	1130	1130	1130	1130	1130	1130		1130	1130	1130	1130	1130	1130	1130	1130	1130	1130
mean	0.590179	0.39285	0.4375	0.521429	0.525893	0.609107	0.586607	0.554466	0.586607	0.545536		2.34643	2.436607	2.674107	2.58643	2.583929	2.496429	2.66875	2.703979	2.357879	1.030174
std	3.720587	2.62032	3.014061	3.332549	3.723851	4.42625	3.522125	3.647404	3.344343	2.57833		9.850514	10.502457	11.845824	11.093983	11.342507	10.709337	11.901746	11.982486	10.622129	4.85186
min	0	0	0	0	0	0	0	0	0	0		0	0	0	0	0	0	0	0	0	0
25%	0	0	0	0	0	0	0	0	0	0		0	0	0	0	0	0	0	0	0	0
50%	0	0	0	0	0	0	0	0	0	0		0	0	0	0	0	0	0	0	0	0
75%	0	0	0	0	0	0	0	0	0	0		0	0	0	0	0	0	0	0	0	0
max	82	63	63	70	100	107	66	73	59	43		146	174	174	182	187	186	212	223	242	106

Se encontró que el municipio de Pueblo Viejo no cuenta con información sobre temperatura y precipitaciones, por lo cual se decidió eliminar de la base este registro teniendo en cuenta que se consideran variables importantes para el análisis de incidencia del dengue. De esta manera, se cuenta con una base de datos limpia con información sobre 1.120 municipios y 1.017 características.

A esta base, se le aplicó el algoritmo K-medias para la serie temporal de casos y población (del 2007 al 2019) y se creó una única serie temporal que es la tasa de casos de dengue por población para realizar un K-medias utilizando Dynamic Time Warping (DTW) como medida de distancia, el resultado es la agrupación de municipios con similar tasa de casos de dengue en el tiempo (2007-2019). Además, se estimaron clústeres de la base con el último dato disponible de las series de tiempo (2019) a través de Clustering Jerárquico y DBSCAN.

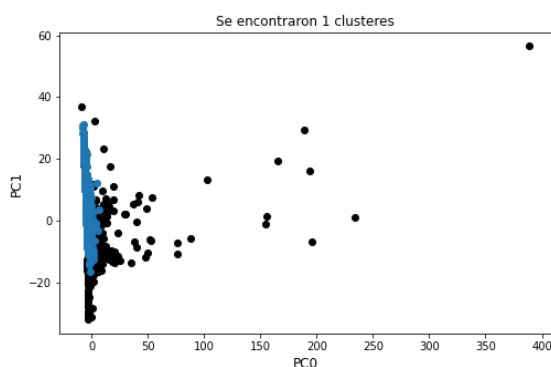
Resultados

Al aplicar K-medias sobre la serie de tiempo de contagios y población, se obtuvieron tres clústeres, entre los cuales se identificó uno que contiene los municipios con mayor número de contagios por población. Estos municipios se caracterizan por ser de clima cálido y húmedo, además de no ser pequeños en comparación con ciudades intermedias.

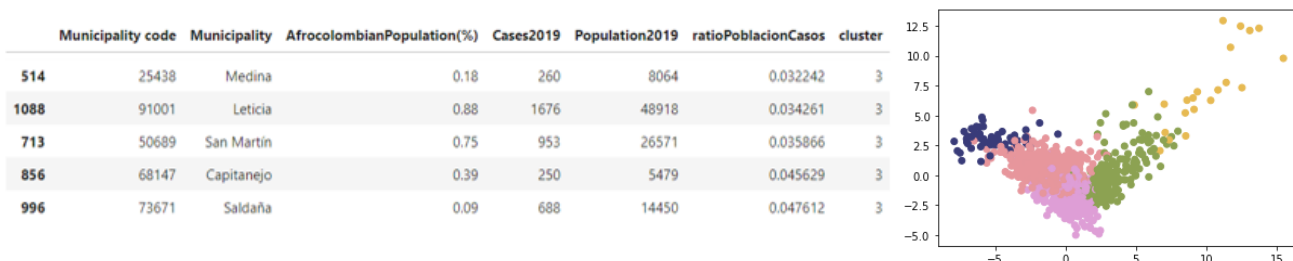
	Municipality	Municipality code	indice_total	cluster_ts
1050	Fortul	81300	0.020372	0
518	Nilo	25488	0.017301	0
690	Castilla La Nueva	50150	0.011041	0
1055	Aguazul	85010	0.010755	0
984	Meigar	73449	0.010042	0
687	Acacias	50006	0.008943	0
960	Alvarado	73026	0.008869	0
981	Lénida	73408	0.008670	0
607	Aipe	41016	0.008577	0
688	Barranca de Upiá	50110	0.008402	0
973	Espinal	73268	0.008287	0
706	Puerto López	50573	0.007699	0
710	San Carlos de Guaroa	50680	0.007435	0
617	Hobo	41349	0.007164	0
1107	San José del Guaviare	95001	0.007164	0
698	Guamal	50318	0.007114	0
1070	Tauramena	85410	0.007080	0
1054	Yopal	85001	0.007030	0
1109	El Retorno	95025	0.007016	0

Por otro lado, se estimó un modelo por DBSCAN a una base reducida de 73 componentes principales que explican el 98% de la varianza total de la base de datos de 2019 y se encontró que los 458 municipios que se clasificaron como ruido o outliers se caracterizan por tener tanto altas como bajas tasas de incidencia de dengue.

	Municipality	Municipality code	Casos_total	Temperatura_promedio	Precipitacion_promedio
9	Anorí	5040	75	22.745777	266.232825
599	San José del Palmar	27660	21	22.383884	260.546915
82	Remedios	5604	252	25.867615	260.301663
5	Amalfi	5031	32	22.061463	257.702365
117	Vegachí	5858	208	24.055471	254.445760
...
297	Susacón	15774	1	10.410217	80.896387
865	Concepción	68207	35	8.622600	77.888286
195	Tunja	15001	23	11.981617	77.742075
288	Siachoque	15740	0	10.423938	77.338384
227	El Cocuy	15244	1	8.628353	74.024666



Por último, se estimó un modelo por Clustering Jerárquico a una base reducida de 21 componentes principales que explican el 95% de la varianza total, obteniendo cinco clústeres, entre los cuales se identificó uno que contiene los municipios con mayor tasa de contagios por población. Estos municipios se caracterizan por ser de clima cálido y húmedo, resultados muy similares a los obtenidos por K-medias sobre series de tiempo.



De acuerdo con lo anterior, se encuentra que a través de los algoritmos K-medias en series de tiempo y Clustering Jerárquico se logran identificar clústeres de municipios que tienen mayor riesgo de presentar brotes de dengue, a diferencia de DBSCAN, con el cual se clasificaron estos municipios como atípicos. Teniendo en cuenta que el interés del proyecto es la identificación de las zonas en las que se deberían concentrar las campañas de concientización de prevención y control contra el dengue, los algoritmos K-medias sobre series de tiempo y Clustering Jerárquico proporcionan los resultados esperados.

Una limitación de la implementación de K-medias sobre series de tiempo es que requiere mayores recursos computacionales si la serie es de gran longitud. Para futuras implementaciones, se podrían tener en cuenta las series de tiempo de temperatura y precipitación para la estimación de un modelo de K-medias de series de tiempo multidimensional. Adicionalmente, estos resultados podrían implementarse para un análisis de geolocalización que permitiera la identificación de puntos calientes en Colombia, además de servir como ruta de priorización por zonas y puntos de mayor concentración.

Conclusión

Con una base de 1.120 municipios de Colombia y sus características demográficas (edad, población), socioeconómicas (desempleo, estratos), y climáticas (temperatura, precipitaciones), se estimaron clústeres por K-medias para series de tiempo y Clustering Jerárquico. Con el primero, se encontró que los municipios como Fortul, Nilo, Castilla La Nueva, entre otros, son aquellos que tienen un mayor riesgo de brote del dengue y se caracterizan en especial por su clima cálido y baja población.

Por otro lado, a través de Clustering Jerárquico se encontró que municipios como Saldaña, Capitanejo, San Martin, Leticia, Medina, entre otros, son los municipios que tienen mayor riesgo de incidencia de dengue, en relación con la proporción de contagios sobre el total de su población. Asimismo, estos municipios se caracterizan por su clima húmedo y cálido. Por último, aunque a través de DBSCAN se esperaba encontrar zonas densas de municipios cuyas características permitieran la identificación de clústeres con mayor tasa de incidencia de dengue, no se obtuvieron resultados que estuvieran alineados con el objetivo de este proyecto, ya que se agruparon como atípicos tanto municipios con alta tasa de incidencia como de baja incidencia.

De esta manera, se encontró que los municipios de clima cálido y mayor humedad hacen parte de las zonas de mayor riesgo y, por tanto, son a estos municipios a los que se deberían dirigir recursos para la capacitación y concientización de la población para la prevención y control del dengue eficientemente.

Bibliografía

- 1 . Instituto Nacional de Salud. Protocolo de Vigilancia de Dengue [Internet]. INS 2022. [citado el 20 de agosto de 2022. Disponible en: https://www.ins.gov.co/buscadoreventos/Lineamientos/Pro_Dengue.pdf
2. Instituto Nacional de Salud. Informe de Evento Dengue [Internet]. INS 2022. Citado el 20 de agosto de 2022. Disponible en: <https://www.ins.gov.co/buscadoreventos/Informesdeevento/DENGUE%20PE%20VII%202022.pdf>
3. World Health Organization. Dengue and severe dengue [Internet]. WHO 2022. Citado el 04 de septiembre de 2022. Disponible en: <https://www.who.int/news-room/fact-sheets/detail/dengue-and-severe-dengue#:~:text=The%20number%20of%20dengue%20cases,affecting%20mostly%20younger%20age%20group.>
4. Organización Panamericana de la Salud. Dengue [Internet]. OPS. 2020. Citado el 20 de agosto de 2022. Disponible en: <https://www.paho.org/es/temas/dengue>
5. Thomas SJ, Rothman AL. Dengue virus infection: Epidemiology. In: UpToDate, Shefner JM (Ed), UpToDate, Waltham, MA. (Accessed on September 04, 2022.)
6. Bhatt S, Gething PW, Brady OJ, et al. The global distribution and burden of dengue. *Nature*. 2013;496(7446):504-507. doi:10.1038/nature12060
7. Shaukat K, Masood N, Shafaat AB, Jabbar K, Shabbir H, Shabbir S. Dengue fever in perspective of clustering algorithms. arXiv preprint arXiv:1511.07353. 2015 Nov 23.
8. G. M. Nandana, S. Mala and A. Rawat, "Hotspot Detection of Dengue Fever Outbreaks Using DBSCAN Algorithm," *2019 9th International Conference on Cloud Computing, Data Science & Engineering (Confluence)*, 2019, pp. 158-161, doi: 10.1109/CONFLUENCE.2019.8776916.
9. Hoyos W, Aguilar J, Toro M. Dengue models based on machine learning techniques: A systematic literature review. *Artif Intell Med*. 2021;119:102157. doi:10.1016/j.artmed.2021.102157
10. Delmelle E, Hagenlocher M, Kienberger S, Casas I. A spatial model of socioeconomic and environmental determinants of dengue fever in Cali, Colombia. *Acta Trop*. 2016;164:169-176. doi:10.1016/j.actatropica.2016.08.028