

# Proyecto de IA Aplicada: Predicción de Educación Primaria en Niños con Autismo

Ruben Dario Sabogal U, [rdsabogal@hotmail.com](mailto:rdsabogal@hotmail.com)  
Cristian Camilo Quebrada, [cristian\\_q\\_90@hotmail.com](mailto:cristian_q_90@hotmail.com)  
Edwin Pérez L., [edwinandperez@gmail.com](mailto:edwinandperez@gmail.com)

**Resumen**— Este proyecto aborda el subregistro y la inconsistencia en la categorización de estudiantes con Trastorno del Espectro Autista (TEA) en el Sistema de Matrícula (SIMAT) de Cali, así como la débil articulación con los Planes Individuales de Ajustes Razonables (PIAR) y los soportes del sector salud. Se plantea como objetivo general fortalecer la identificación, el registro y el seguimiento de la población con TEA para mejorar el acceso, la permanencia y la pertinencia de los apoyos educativos. La metodología combina: (i) evaluación de calidad de datos (completitud, coherencia y homologación de variables críticas), (ii) estandarización de criterios y evidencias para la categorización TEA y su alineación PIAR↔SIMAT, (iii) diseño de tableros e indicadores por año, comuna, IEO y grado, con panel de calidad, y (iv) un pipeline analítico opcional (EDA + métricas ROC/PR, sesgos, interpretabilidad) como apoyo al tamizaje, con resguardo ético y de privacidad. El plan de trabajo (septiembre–noviembre de 2025) incluye entregables de datos, análisis, documentación operativa, capacitación y cierre con plan de sostenibilidad y KPIs. Se espera disminuir tiempos de acceso a apoyos, reducir riesgo de deserción y mejorar la focalización territorial, mediante una gobernanza clara de roles, validaciones y plazos intersectoriales.

**Palabras clave**— Educación inclusiva, TEA, SIMAT, PIAR, calidad de datos, analítica educativa, ética de datos.

## I. INTRODUCCION

La identificación y el acompañamiento oportuno de estudiantes con Trastorno del Espectro Autista (TEA) constituyen un reto técnico y de gestión para los sistemas de información educativa. En el contexto de Cali, se ha detectado subregistro y falta de consistencia en la categorización de TEA dentro del Sistema Integrado de Matrícula (SIMAT), además de una débil articulación con los Planes Individuales de Ajustes Razonables (PIAR) y con los soportes del sector salud; esta situación afecta la planificación de apoyos, el seguimiento a la permanencia y la focalización territorial por comuna e institución (IEO).



Fig.1. Que es el autismo



Fig. 2. Sistema Integrado de matricula

Este informe propone un marco de intervención para mejorar la calidad del registro y del seguimiento de estudiantes con TEA en SIMAT, mediante: (i) evaluación de la calidad y coherencia de variables críticas; (ii) estandarización de criterios y evidencias para la categorización (incluida la clasificación como discapacidad psicosocial cuando aplique) y definición clara de responsabilidades entre IEO y entidad territorial certificada (ETC); y (iii) articulación operativa entre SIMAT y PIAR con respaldo del sector salud.

Desde el componente analítico-tecnológico, se plantea la implementación de análisis exploratorio de datos (EDA) sistemático, tableros con filtros por año, comuna, IEO y grado—incluyendo un panel de calidad—y, de manera

opcional, un pipeline de “screening” con métricas y validación ética para apoyar decisiones sin reemplazar criterios pedagógicos ni clínicos.



Fig. 3. ¿Qué es el PIAR?



Fig. 4. El PIAR

Las contribuciones de este trabajo son: (1) un diagnóstico estructurado del problema (árbol de causas-efectos) que explica demoras en apoyos, riesgo de deserción y desigualdades de inclusión; (2) un conjunto de objetivos y lineamientos operativos para alinear datos, procesos y responsabilidades intersectoriales; y (3) un plan de trabajo con entregables y cronograma (septiembre–noviembre de 2025) que integra datos, análisis, tableros, gestión del cambio y cierre con plan de sostenibilidad e indicadores clave.

El resto del documento se organiza así: la Sección 2 sintetiza el estado del arte y el marco normativo aplicable; la Sección 3 describe la metodología de evaluación y estandarización de datos; la Sección 4 presenta el diseño de tableros y el pipeline analítico opcional; la Sección 5 detalla el plan de trabajo y entregables; y la Sección 6 discute resultados esperados, consideraciones éticas y limitaciones.

## 1. Diagnóstico del Problema:

Existe un subregistro y falta de consistencia en la categorización de estudiantes con Trastorno del Espectro Autista (TEA) en el sistema SIMAT, acompañado de una débil articulación con el PIAR y los soportes del sector salud. Esta situación afecta negativamente la planificación educativa, la asignación de apoyos, el seguimiento a la permanencia escolar y la focalización territorial (comuna/IEO).

## 2. Objetivos del Proyecto

### Objetivo general:

Mejorar la identificación, el registro y el seguimiento de estudiantes con TEA en SIMAT en la ciudad de Cali, fortaleciendo el acceso, la permanencia y la pertinencia de los apoyos educativos (PIAR), mediante una articulación efectiva entre los sectores de educación y salud.

### Objetivos específicos:

- Evaluar la calidad, completitud y coherencia de las variables relevantes en SIMAT (discapacidad/TEA, apoyos, IEO, comuna, grado, etc.) y su alineación con la ruta de registro.
- Estandarizar criterios y evidencias para la categorización de TEA (como discapacidad psicosocial cuando corresponda), definiendo responsabilidades y validaciones en IEO y ETC.

- Garantizar la articulación entre SIMAT y PIAR (ajustes razonables, ayudas técnicas) con respaldo del sector salud.
- Diseñar y validar un pipeline analítico (EDA + indicadores, con opción de ML) para el tamizaje y seguimiento, con enfoque ético y de privacidad.
- Implementar tableros y reportes periódicos por año, comuna, IEO y grado, incluyendo un panel de calidad de datos.

### 3. Árbol de Problemas

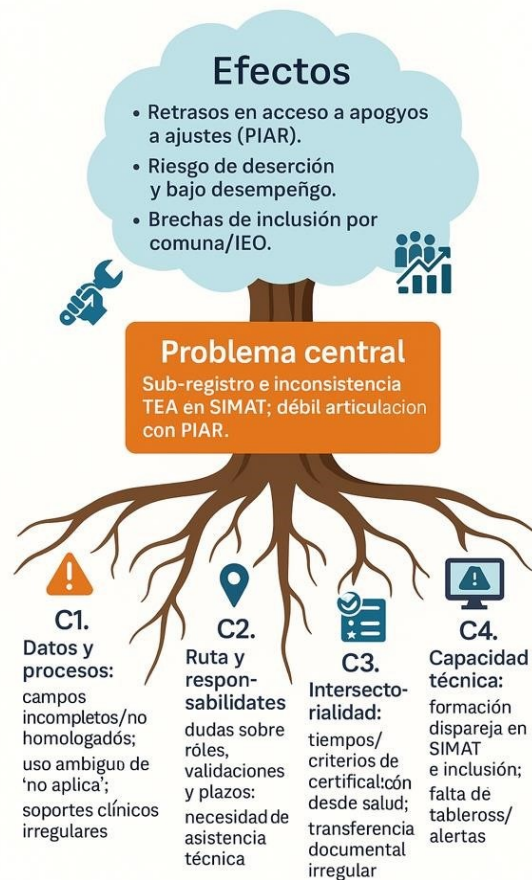


Fig. 5. árbol de Problemas

#### Efectos:

- Demoras en el acceso a apoyos y ajustes (PIAR).
- Riesgo de deserción escolar y bajo rendimiento académico.
- Desigualdades en inclusión educativa según comuna o IEO.

#### Causas principales:

- **C1. Datos y procesos:** Campos incompletos o no homologados, uso ambiguo de “no aplica”, soportes clínicos irregulares.

- **C2. Ruta y responsabilidades:** Falta de claridad en roles, validaciones y plazos; necesidad de asistencia técnica.
- **C3. Intersectorialidad:** Tiempos y criterios de certificación desde salud; transferencia documental inconsistente.
- **C4. Capacidad técnica:** Formación desigual en SIMAT e inclusión; ausencia de tableros y alertas.

#### 4. Estado del Arte (Resumen)

- **Definición y prevalencia (OMS):** El autismo es un conjunto de condiciones del neurodesarrollo con una prevalencia global estimada de 1 en 127. Se recomiendan intervenciones psicosociales y comunitarias.
- **Marco normativo colombiano (Decreto 1421/2017):** Establece la atención educativa para personas con discapacidad, incluyendo ajustes razonables y responsabilidades del MEN, ETC e IEO.
- **Práctica local (Cali/SIMAT):** Ruta de registro, PIAR, responsabilidades y validaciones; el TEA se reporta como discapacidad psicosocial cuando aplica.
- **Tendencias analíticas:** Uso de EDA y tableros para focalización; aplicación de ML como herramienta de apoyo, con énfasis en desbalance, interpretabilidad y validación ética/clínica.

Referencias (formato IEEE)
[1] World Health Organization, “Autism,” Sep. 17, 2025. [Online]. Available: Organización Mundial de la Salud
[2] Centers for Disease Control and Prevention, “Data and Statistics on Autism Spectrum Disorder,” May 27, 2025. [Online]. Available: cdc.gov
[3] C. Lord et al., “The future of care and clinical research in autism—The Lancet Commission,” Lancet, 2021/2022. [Online]. Available: pmc.ncbi.nlm.nih.gov+1
[4] A. T. Wieckowski et al., “Sensitivity and Specificity of the M-CHAT(-R/F): Systematic Review and Meta-analysis,” JAMA Pediatr., 2023. [Online]. Available: JAMA Network
[5] Screening and diagnostic tools for ASD: systematic review and meta-analysis, 2023. [Online]. Available: PubMed
[6] K. A. Khowaja et al., “Single and Repeat Screening with M-CHAT-R in High-Likelihood Children,” 2024. [Online]. Available: pmc.ncbi.nlm.nih.gov
[7] Y. Ding, H. Zhang, and T. Qiu, “Deep learning approach to predict ASD: a systematic review and meta-analysis,” BMC Psychiatry, 2024. [Online]. Available: BioMed Central
[8] M. Briguglio et al., “ML approach to ASD diagnosis using ADOS-2,” Brain Sciences, 2023. [Online]. Available: MDPI
[9] X. Lian and M. S. Sunar, “Mobile AR Technologies for ASD Interventions: SLR,” Applied Sciences, 2021. [Online]. Available: MDPI
[10] “Exploring the Impact of AR in Children/Adolescents with ASD: Systematic Review,” IJERPH, 2020. [Online]. Available: MDPI
[11] “Augmented Reality and Learning-Cognitive Outcomes in ASD: Systematic Review,” Children, 2024. [Online]. Available: MDPI
[12] “Using AR Toward Improving Social Skills: Scoping Review,” JMIR Serious Games, 2023. [Online]. Available: games.jmir.org

[13] “Immersive Technology to Teach Social Skills to Students with ASD: Literature Review,” Review Journal of Autism and Developmental Disorders, 2021/2022. [Online]. Available: SpringerLink
[14] “School-Based Interventions for Increasing Autistic Pupils’ Social Inclusion,” Review Journal of Autism and Developmental Disorders, 2024. [Online]. Available: SpringerLink
[15] “School-based interventions targeting social communication,” Systematic Review, 2018. [Online]. Available: PubMed
[16] “School-based social skills interventions in inclusive settings,” 2021. [Online]. Available: PubMed
[17] Ministerio de Educación Nacional (Colombia), “Decreto 1421 de 2017,” 2017. [Online]. Available: Ministerio de Educación
[18] ICFES, “Compilación jurídica del Decreto 1421 de 2017,” 2019. [Online]. Available: normograma.icfes.gov.co
[19] MEN, “Orientaciones técnicas, administrativas y pedagógicas (PIAR),” 2017/2020. [Online]. Available: Ministerio de Educación+1
[20] MEN, “Seguimiento PIAR y PIP (reporte 2020),” 2020. [Online]. Available: Ministerio de Educación
[21] SITEAL–IIEP UNESCO, “Marco de educación inclusiva – Decreto 1421/2017,” 2017. [Online]. Available: siteal.iiep.unesco.org
[22] INCI, “El INCI y la educación inclusiva / datos SIMAT 2017,” 2024. [Online]. Available: inci.gov.co
[23] V. T. Badillo-Jiménez et al., “Percepción de inclusión escolar (Colombia),” Duazary, 2022. [Online]. Available: revistas.unimagdalena.edu.co
[24] “Los saberes para regular el autismo en Colombia,” Rev. Cienc. Salud, 2022. [Online]. Available: revistas.uurosario.edu.co
[25] A. Canal-Bedia et al., “Spanish validation of M-CHAT-R/F,” 2018. [Online]. Available: PubMed
[26] L. Gutierrez-Rojas et al., “Autism Assessment with English-Spanish Bilinguals,” JADD, 2025. [Online]. Available: PubMed+1
[27] A. Guerrero-Arias et al., “ADEC in Low-Income Spanish-Speaking Population,” JADD, 2024. [Online]. Available: SpringerLink

Fig. 6 Bibliografía- Estado del Arte

#	Enlace / Referencia	Qué estudiaron	Hallazgos clave / utilidad para tamizaje
1	<a href="#">JAMA Netw Open 2024 — ML con 28 variables mínimas (n=30.660). (JAMA Network)</a>	EHR/antecedentes	Alta sensibilidad/especificidad; hitos del desarrollo y alimentación como <i>features</i> informativas; <b>generaliza</b> a cohortes externas.
2	<a href="#">JAMA Netw Open 2024 — vigilancia del desarrollo vs M-CHAT. (JAMA Network)</a>	Datos de vigilancia rutinaria	Predicción de ASD <b>superando M-CHAT</b> ; integrable en flujo clínico.
3	<a href="#">Sci Rep 2025 — BORN Ontario (nacimientos + admin). (Nature)</a>	Registros poblacionales	<b>Transformers/ensembles</b> identifican 18m–5a con mayor probabilidad ASD; viabilidad <b>poblacional</b> .
4	<a href="#">Angell 2025 — equidad en EHR ML. (PubMed)</a>	EHR; fairness	Diferencias por <b>sexo</b> ; recomienda métricas de justicia (Equal Opportunity/Eq. Odds) y auditorías antes de uso real.
5	<a href="#">Pan 2025 — abrev. M-CHAT-R con ML. (ScienceDirect)</a>	Cuestionario	Subconjuntos óptimos mantienen poder de cribado con menor carga (útil en primaria).
6	<a href="#">Sci Rep 2025 — Q-CHAT-10 subsets + ML. (Nature)</a>	Cuestionario	Ítems <b>compactos</b> generalizan a <b>diagnóstico clínico</b> externo (NZ/SA → Polonia).
7	<a href="#">npj Digit Med / Nature 2025 — videos caseros automatizados. (Nature)</a>	Video DL	<b>Factible</b> : 3 tareas cortas, extracción de conducta, buen desempeño para cribado remoto.
8	<a href="#">Metaanálisis 2025 — tele-video screening. (PMC)</a>	Revisión + meta	Evidencia acumulada de <b>buena exactitud</b> y <b>conveniencia</b> en tele-salud.
9	<a href="#">Frontiers Neuroinformatics 2025 — prosodia/voz. (Frontiers)</a>	Audio/voz	Rasgos acústicos y prosódicos clasifican subgrupos; utilidad <b>complementaria</b> en cribado digital.
10	<a href="#">CIHI (case study) — datos reales + AI. (cihi.ca)</a>	Implementación	Caso de uso institucional sobre cómo escalar tamizaje temprano con datos administrativos.

11	<a href="#">Revisión de sesgos (J Clin Epidemiol 2025). (jclinepi.com)</a>	Fairness	Catálogo de <b>sesgos</b> clínicos en ML y grupos PROGRESS; guía para <b>auditorías</b> .
12	<a href="#">Estudio en prensa/medios sobre el de JAMA 2024 (contexto público). (The Guardian)</a>	Divulgación	Resume el hallazgo de <b>28 medidas</b> y advierte <b>no sustituir</b> métodos clínicos; útil para stakeholders.

**Fig. 7 Bibliografía - Estado del Arte -Modelos de machine learning específicos para tamizaje de TEA**

## 5. Análisis de Soluciones Existentes

### Político-institucional:

- Aplicación del Decreto 1421/2017 y lineamientos de educación inclusiva.
- Definición de indicadores clave (KPIs) sobre acceso, permanencia y calidad del registro TEA.

### Operativo (SIMAT/PIAR):

- Formalización de la ruta, roles y validaciones (actas, soportes, plazos).
- Homologación de IEO, comunas y grados; guías para el uso de “no aplica” y actualización diagnóstica.
- Alineación de registros SIMAT con PIAR y apoyos pedagógicos.

### Tecnológico-analítico:

- Implementación de EDA sistemático y tableros con filtros por año, comuna, IEO y grado, incluyendo panel de calidad.
- Desarrollo de pipeline de screening (opcional), con validación ética y reporte de sesgos.

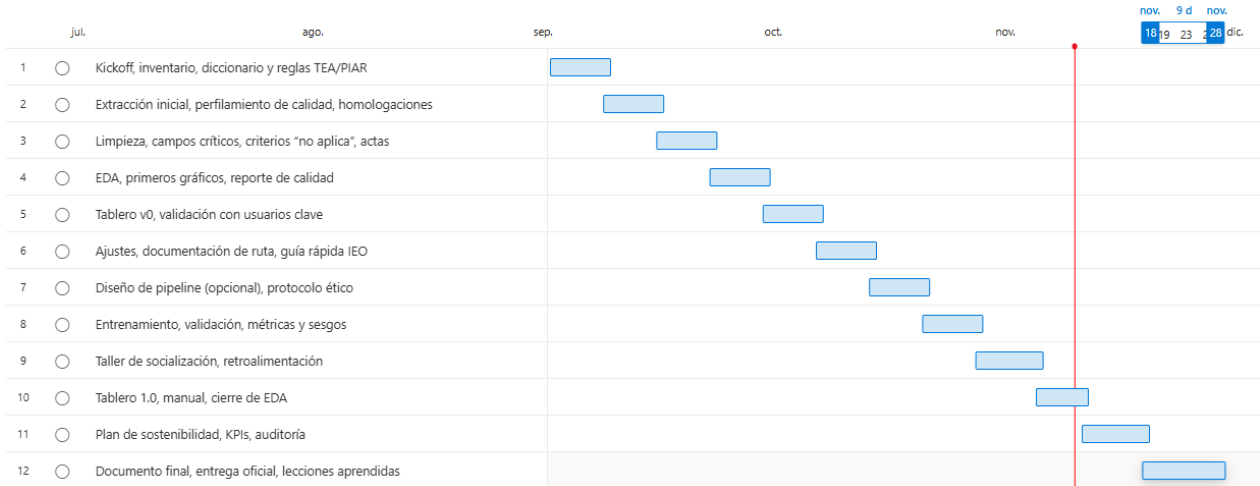
## 6. Planeación - Plan de Trabajo (1 Sep – 30 Nov 2025)

### Tareas y Entregables:

- **Datos:** Inventario de fuentes SIMAT, diccionario de datos, reglas de categorización TEA, protocolos PIAR ↔ SIMAT.
- **Análisis:** EDA de calidad (faltantes, duplicados, homologaciones), visualización de tendencias.
- **Modelo de apoyo (opcional):** Diseño y validación de pipeline, métricas (ROC/PR), reporte de sesgos.
- **Tableros:** Dashboard con filtros, manual de usuario, plan de actualización.
- **Gestión del cambio:** Talleres con responsables SIMAT/Orientación/Apoyos, guía de “no aplica”, plazos de actualización.
- **Cierre:** Documento final con diagnóstico, resultados, tableros y plan de sostenibilidad; acta de lecciones aprendidas y plan 2026.



## Cronograma por Semana:



## II. Análisis Exploratorio de Datos

### Propósito

Evaluar la calidad, consistencia y completitud de los datos disponibles, con el fin de garantizar que las variables empleadas para el modelado reflejen información relevante, representativa y libre de sesgos o ruido estadístico.

### Descripción de la Exploración

El dataset analizado está compuesto por 520 registros y 61 variables, vinculadas a instituciones educativas y características sociodemográficas de los estudiantes. El análisis exploratorio incluyó: - Identificación de tipos de variables (numéricas, categóricas, mixtas) y su distribución general. - Detección de valores faltantes y su impacto en la calidad del conjunto de datos. Por ejemplo, *num\_convenio* presenta un 80.19% de faltantes, y *nombre2* un 38.08%. - Revisión de valores atípicos y de consistencia interna (por ejemplo, variables con baja varianza o sin dispersión significativa). Clasificación preliminar de las variables según su posible utilidad analítica (predictiva o descriptiva) y eliminación de redundancias.

### Herramientas Utilizadas

Se implementó un enfoque reproducible en Python, empleando librerías como pandas, numpy, matplotlib, seaborn y scikit-learn. El flujo de trabajo incluyó: - Imputación de datos mediante media (numéricos) y moda (categóricos) para minimizar pérdida de información. - Codificación de variables categóricas con LabelEncoder para su integración al proceso de modelado. - Estandarización de variables con StandardScaler, asegurando comparabilidad entre magnitudes y evitando sesgos de escala.

### Criterios de Selección de Variables

- Variables con menos del 50% de valores faltantes se conservaron para modelado.
- Eliminación de variables constantes o con coeficiente de variación inferior al 1%.
- Exclusión de campos identificadores, direcciones y nombres sin valor analítico.
- Descarte de variables con una categoría dominante (>95%).

Como resultado, se seleccionaron 41 variables (27 numéricas y 14 categóricas) que ofrecen un balance adecuado entre calidad y diversidad informativa. El análisis permitió depurar la base de datos conservando aquellas variables que aportan valor al entendimiento de los patrones de comportamiento. El resultado es un conjunto más limpio, consistente y representativo, ideal para la aplicación de técnicas no supervisadas de agrupamiento.

### III. Modelos de Referencia y Experimentos

#### Hipótesis de Trabajo

- H1: Es posible identificar grupos homogéneos dentro del conjunto de datos, reflejando características comunes entre los individuos.
- H2: Las variables numéricas aportan mayor capacidad de diferenciación que las categóricas.
- H3: Existen relaciones no lineales que justifican el uso de algoritmos basados en densidad o jerarquía.
- H4: El tratamiento de valores faltantes y la estandarización afectan directamente la estabilidad de los clusters.

#### Modelos Explorados

Se implementaron tres técnicas de agrupamiento no supervisado:

Modelo	Descripción	Nº Clusters	Silhouette	Davies-Bouldin	Calinski-Harabasz
<b>K-Means</b>	Agrupamiento particional con distancia euclídea	8	0.260	1.325	41.86
<b>DBSCAN</b>	Basado en densidad, sin necesidad de definir K	2	0.410	0.818	14.99
<b>Agglomerative</b>	Jerárquico, método Ward	8	0.240	1.302	38.81

=====

DETERMINACIÓN DEL NÚMERO ÓPTIMO DE CLUSTERS

=====

Evaluando diferentes números de clusters...

2 clusters: Silhouette=0.237, Inertia=18893.63

3 clusters: Silhouette=0.242, Inertia=17515.19

4 clusters: Silhouette=0.249, Inertia=16387.52

5 clusters: Silhouette=0.242, Inertia=15334.87

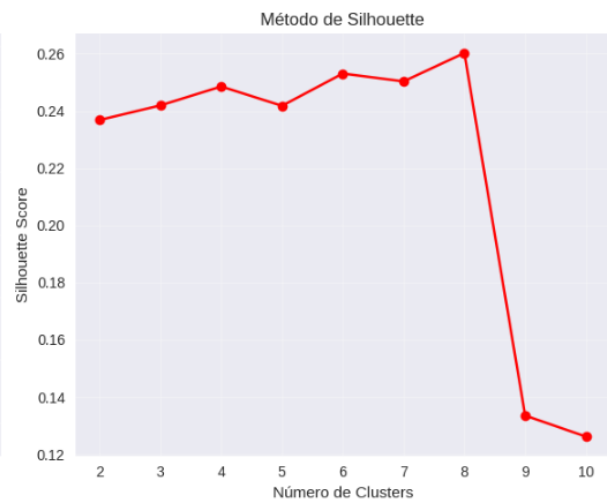
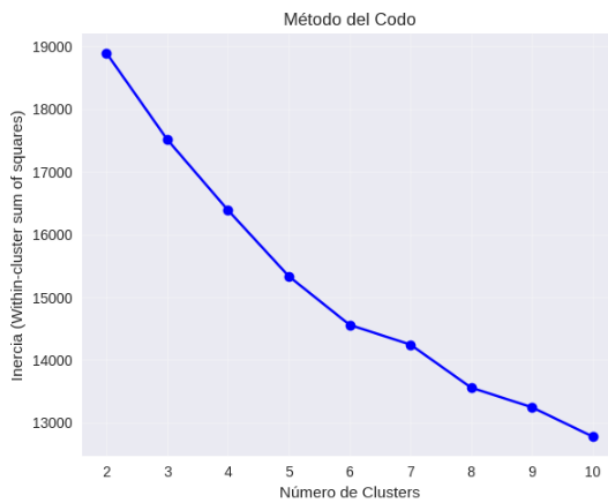
6 clusters: Silhouette=0.253, Inertia=14560.44

7 clusters: Silhouette=0.250, Inertia=14245.57

8 clusters: Silhouette=0.260, Inertia=13559.06

9 clusters: Silhouette=0.133, Inertia=13246.14

10 clusters: Silhouette=0.126, Inertia=12778.60



Número óptimo de clusters (según Silhouette): 8

Silhouette Score: 0.2602

Número sugerido (método del codo): 6



# MODELO 1: K-MEANS CLUSTERING

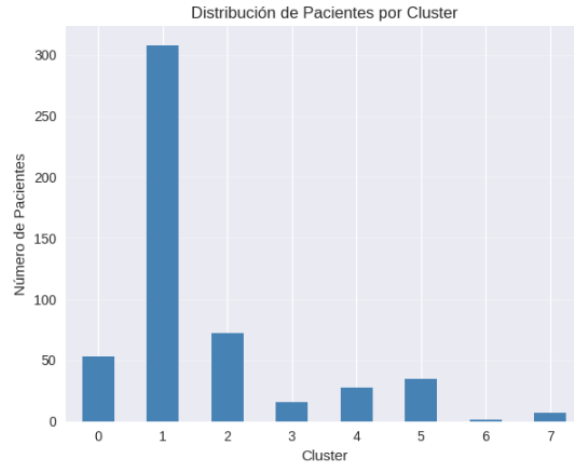
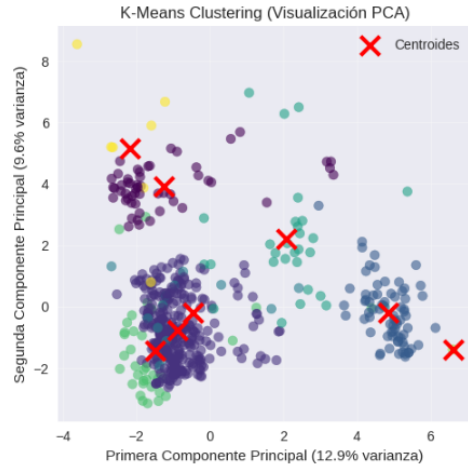
Usando número óptimo de clusters: 8

## Métricas de Rendimiento:

Silhouette Score: 0.2602 (mayor es mejor, rango: -1 a 1)  
 Davies-Bouldin Index: 1.3253 (menor es mejor)  
 Calinski-Harabasz: 41.87 (mayor es mejor)  
 Inercia: 13559.06 (menor es mejor)

## Distribución de pacientes por cluster:

Cluster 0: 53 pacientes (10.2%)  
 Cluster 1: 308 pacientes (59.2%)  
 Cluster 2: 72 pacientes (13.8%)  
 Cluster 3: 16 pacientes (3.1%)  
 Cluster 4: 28 pacientes (5.4%)  
 Cluster 5: 35 pacientes (6.7%)  
 Cluster 6: 1 pacientes (0.2%)  
 Cluster 7: 7 pacientes (1.3%)



# MODELO 2: DBSCAN CLUSTERING

## Parámetros DBSCAN:

eps estimado: 7.088  
 min\_samples: 5

## Resultados DBSCAN:

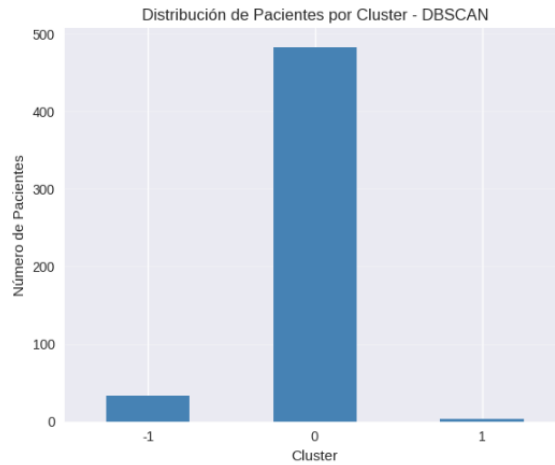
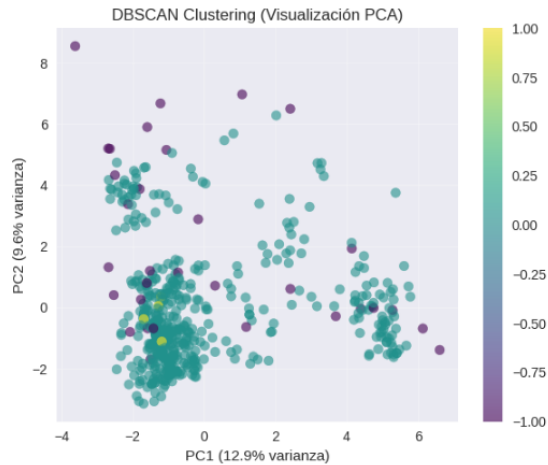
Número de clusters: 2  
 Puntos de ruido (outliers): 33 (6.3%)

## Métricas de Rendimiento (sin ruido):

Silhouette Score: 0.4107  
 Davies-Bouldin Index: 0.8177  
 Calinski-Harabasz: 14.99

## Distribución de pacientes por cluster:

Ruido (outliers): 33 pacientes (6.3%)  
 Cluster 0: 483 pacientes (92.9%)  
 Cluster 1: 4 pacientes (0.8%)



#### MODELO 3: AGGLOMERATIVE CLUSTERING (JERÁRQUICO)

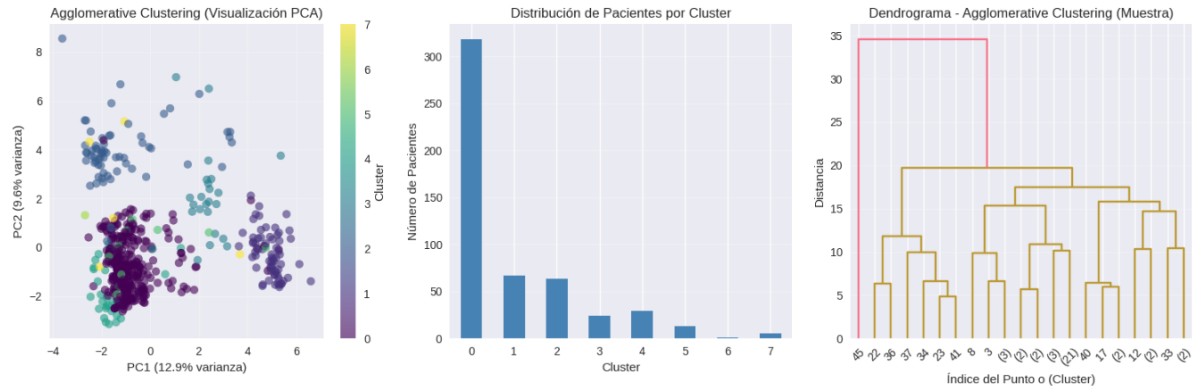
Usando número óptimo de clusters: 8

##### Métricas de Rendimiento:

Silhouette Score: 0.2397 (mayor es mejor, rango: -1 a 1)  
 Davies-Bouldin Index: 1.3016 (menor es mejor)  
 Calinski-Harabasz: 38.81 (mayor es mejor)

##### Distribución de pacientes por cluster:

Cluster 0: 318 pacientes (61.2%)  
 Cluster 1: 67 pacientes (12.9%)  
 Cluster 2: 63 pacientes (12.1%)  
 Cluster 3: 24 pacientes (4.6%)  
 Cluster 4: 29 pacientes (5.6%)  
 Cluster 5: 13 pacientes (2.5%)  
 Cluster 6: 1 pacientes (0.2%)  
 Cluster 7: 5 pacientes (1.0%)



## Resultados y Validación

- Los métodos del codo y Silhouette determinaron que el rango óptimo de agrupación se encuentra entre 6 y 8 clusters.
- DBSCAN obtuvo el mejor rendimiento global (Silhouette = 0.4107), destacando por su capacidad para detectar estructuras no lineales y manejar ruido.
- El análisis comparativo confirmó H1 y parcialmente H3, evidenciando patrones latentes que reflejan diferencias sustanciales entre grupos.

## Interpretación de Resultados

Los tres modelos ofrecen perspectivas complementarias: mientras K-Means permite segmentaciones claras y replicables, DBSCAN revela zonas densas de comportamiento atípico, y Agglomerative facilita una lectura jerárquica útil para interpretar relaciones entre grupos. La combinación de métricas y visualizaciones sugiere que el dataset presenta estructuras complejas que pueden aprovecharse en etapas de predicción supervisada.