

**Estrategias basadas en ciencia de datos para la predicción y reducción de la deserción en
instituciones públicas de educación superior en Medellín**

Camilo Andrés Losada Ule

Asesor

Lina Rocio Rivadeneira Munoz

Universidad Nacional Abierta y a Distancia UNAD
Escuela de Ciencias Básicas, Tecnología e Ingeniería ECBTI
Especialización en Ciencia de Datos y Analítica

2025

Nota de Aceptación

Nombre Director de Trabajo de Grado

Jurado

Jurado

Dedicatoria

Quiero dar gracias a Dios por permitirme la oportunidad de estudiar y alcanzar uno de mis más grandes sueños, ser especialista.

Dedico este trabajo a mi esposa Ana Victoria, mis abuelos Obdulio y Raquel, mi madre, Pricila y mi tía Marelvi , que han sido fuerza y motor para seguir alcanzando mis propósitos. Con todo mi cariño y amor les dedico este escrito.

Resumen

La presente monografía analiza el fenómeno de la deserción estudiantil en instituciones públicas de educación superior en la ciudad de Medellín entre los años 2014 y 2025, con un enfoque centrado en el uso de técnicas de ciencia de datos y analítica. A partir de la revisión crítica de investigaciones académicas, bases de datos oficiales como SPADIES y publicaciones institucionales, se identifican los factores clave asociados al abandono académico y se evalúa el grado de implementación de modelos predictivos en universidades públicas locales.

La investigación, desarrollada desde una mirada técnica y aplicada, permite evidenciar que si bien existen avances en el uso de herramientas como regresión logística, árboles de decisión y minería de datos educativa, aún persisten retos en términos de interoperabilidad de datos, capacitación técnica del personal y articulación entre las políticas públicas y las plataformas de monitoreo estudiantil.

Con base en el análisis de fuentes como los indicadores de SPADIES 3.0 y los sistemas de alerta temprana implementados por universidades como la UNAD y la Universidad de Antioquia, se construyó un compendio de metodologías utilizadas para identificar estudiantes en riesgo. Este trabajo también destaca el papel del Observatorio de Educación Superior de Medellín (SAPIÉNCIA) en la consolidación de buenas prácticas basadas en evidencia.

Como autor de esta monografía, sostengo que la integración de ciencia de datos en la gestión educativa no debe limitarse a fines estadísticos, sino que debe convertirse en una estrategia permanente para la toma de decisiones centrada en el bienestar estudiantil. Este estudio busca aportar tanto a la comprensión técnica del problema como a la formulación de soluciones replicables que fortalezcan la permanencia en la educación superior pública en Medellín.

Palabras claves: Deserción, datos, analítica, educación, predicción.

Abstract

This monograph explores the phenomenon of student dropout in public higher education institutions in Medellín, Colombia, from 2014 to 2025, focusing on the application of data science and analytics techniques. Through a comprehensive review of academic literature, official databases such as SPADIES, and institutional reports, key factors associated with academic withdrawal are identified and the extent to which predictive models are being implemented is assessed.

The study adopts a technical and applied perspective to examine how universities have incorporated tools such as logistic regression, decision trees, and educational data mining to detect students at risk of dropping out. Despite progress, significant challenges remain in terms of data integration, institutional coordination, and the development of analytical capabilities in the academic sector.

By analyzing predictive systems used by institutions such as UNAD and Universidad de Antioquia, as well as local policy instruments led by SAPIÉNCIA, this research compiles a methodological framework that facilitates understanding and comparison of analytical approaches. As the author, I argue that data science should not be limited to descriptive reporting but must serve as a continuous strategy to enhance decision-making and student retention in public higher education.

Keywords: Dropout, education, analytics, data, prediction.

Tabla de Contenido

Introducción	9
Justificación	10
Objetivos.....	12
Objetivo General	12
Objetivos Específicos.....	12
Fenómeno de la Deserción Estudiantil en Universidades Públicas de Medellín	14
Tipos de Deserción.....	16
Técnicas de Ciencia de Datos y Analítica para el Análisis de la Deserción.....	17
Modelos Predictivos y Desafíos Futuros con la Ciencia de Datos.....	18
Desarrollo de Objetivos	20
Recolectar y Compilar Información.....	20
Identificar los Enfoques Metodológicos y Herramientas.....	22
Aplicar Técnicas de Ciencia de Datos	24
Comparar los Resultados Obtenidos	26
Diseñar un Compendio Estructurado	29
Conclusiones.....	37
Recomendaciones	39
Referencias Bibliográficas	41
Apéndices.....	44

Lista de Tablas

Tabla 1 <i>Técnicas de Ciencia de Datos y Analítica usadas</i>	32
Tabla 2 <i>Tasa de Deserción Anual IES Padre</i>	35
Tabla 3 <i>Tasa de Deserción por Nivel de Formación</i>	35

Lista de Apéndices

Apéndice A <i>Tabla Histórica de Tasas de Deserción 2014 –2025</i>	44
Apéndice B <i>Código en Python para el Modelo de Regresión Logística</i>	45
Apéndice C <i>Visualización de Riesgo de Deserción por Institución</i>	46
Apéndice D <i>Fragmento del Informe del ODES – SAPIENCIA</i>	47
Apéndice E <i>Correlación entre Deserción y Factores Sociales, Académicos, Económicos y Políticos</i>	48
Apéndice F <i>Código Python para la Matriz de Confusión</i>	50

Introducción

Esta monografía pretende abordar las investigaciones realizadas por estudiantes, docentes, investigadores, entes públicos y privados sobre los métodos de ciencias de datos y analítica usados para compilar información sobre la deserción estudiantil en Universidades públicas de la ciudad de Medellín entre el año 2014 hasta el presente.

Abordar esta problemática social desde el campo de la ciencia de datos y analítica, permitirá conocer cuáles son las estrategias y medidas que ha tomado gracias a la información recopilada, para ofrecer soluciones a los estudiantes y universidades que sufren por este tema. Pocas son las universidades públicas de Medellín, por tanto, los datos recopilados de estas instituciones permitirán obtener unos resultados más concretos. Adicional, al ser instituciones públicas podemos desde sus páginas web acceder a las bases de datos publicadas por ellos y contrastar la información con lo documentado hasta hoy.

En muchas investigaciones se ha abordado sobre el tema, pero más desde un ámbito sociológico, en la presente monografía lo abordaremos más desde las técnicas usadas en la ciencia de datos y analítica para dar conclusiones sobre la problemática, tratando de responder si los estudios han usado técnicas de ciencia de datos avanzadas y apropiadas o solo se ha desarrollado desde un ámbito convencional.

Justificación

La deserción estudiantil en las instituciones públicas de educación superior de Medellín es un problema crítico que trasciende el ámbito académico, afectando tanto a las personas involucradas como al desarrollo económico y social de la región. De acuerdo con informes recientes, las tasas de abandono académico han sido persistentes, lo que refleja la necesidad urgente de implementar estrategias innovadoras que aborden esta problemática desde sus raíces. Las condiciones sociales y económicas específicas de Medellín agravan el desafío, subrayando la importancia de enfoques basados en evidencia.

Este trabajo tiene como propósito dar a conocer los esfuerzos realizados por el gobierno local y nacional plasmado en diferentes investigaciones y publicaciones a lo largo de los últimos 10 años. Adicional, justificar el uso de la ciencia de datos como una herramienta clave para combatir la deserción estudiantil, dado su potencial para analizar grandes volúmenes de información y descubrir patrones subyacentes que influyen en el abandono académico.

La capacidad de los sistemas predictivos para procesar datos socioeconómicos, académicos y contextuales permite que las instituciones identifiquen de manera temprana a los estudiantes en riesgo de deserción, facilitando un uso más eficiente de los recursos disponibles. Además, estas tecnologías ofrecen la posibilidad de desarrollar plataformas interactivas que permitan monitorear el progreso de las intervenciones y ajustar estrategias en tiempo real.

La relevancia de esta monografía también radica en su alineación con los Objetivos de Desarrollo Sostenible (ODS), particularmente en lo que respecta a garantizar una educación inclusiva y equitativa (ODS 4). Reducir las tasas de deserción no solo mejora los índices de permanencia estudiantil, sino que también impulsa la movilidad social, promueve la equidad y fortalece las capacidades humanas. (Naciones Unidas, 2023)

En definitiva, este trabajo pretende evidenciar soluciones prácticas y replicables que integran ciencia de datos con enfoques educativos, optimizando la toma de decisiones y ofreciendo un impacto positivo tanto en el ámbito académico como en la sociedad en general. Su implementación representa un paso significativo hacia un sistema educativo más eficiente, justo y adaptado a las necesidades del entorno actual.

Objetivos

Objetivo General

Analizar las causas y factores asociados a la deserción estudiantil en instituciones públicas de educación superior en Medellín, mediante el uso de técnicas de ciencia de datos y analítica aplicada a diversas fuentes de información entre 2014 y la actualidad, con el fin de consolidar un compendio accesible que facilite la comprensión integral del fenómeno y difunda las metodologías analíticas empleadas, en coherencia con las evidencias reportadas en publicaciones científicas.

Objetivos Específicos

Recolectar y compilar información relevante sobre la deserción estudiantil en instituciones públicas de educación superior en Medellín, mediante la búsqueda sistemática en libros, artículos académicos y bases de datos científicas publicadas desde el 2014 hasta el presente, haciendo énfasis en aquellas que emplean ciencia de datos y técnicas analíticas.

Identificar los enfoques metodológicos y herramientas de ciencia de datos utilizados en estudios académicos y científicos sobre deserción estudiantil, destacando patrones comunes y variaciones significativas.

Aplicar técnicas de ciencia de datos para analizar los factores clave de la deserción estudiantil, identificando patrones, tendencias y correlaciones significativas en la información recopilada.

Comparar los resultados obtenidos en distintas investigaciones, para establecer tendencias, hallazgos coincidentes o divergentes, y vacíos de conocimiento en el uso de la ciencia de datos aplicada al fenómeno.

Diseñar un compendio estructurado que sistematice las principales técnicas, modelos y enfoques de ciencia de datos y analítica utilizados en investigaciones sobre deserción estudiantil, con el fin de facilitar su comprensión, comparación y posible aplicación en futuros estudios o análisis institucionales.

Fenómeno de la Deserción Estudiantil en Universidades Públicas de Medellín

Esta monografía se enfoca en evidenciar estrategias innovadoras que empleen ciencia de datos y analítica para abordar la problemática de la deserción estudiantil en las instituciones públicas de educación superior de Medellín. A partir de un análisis profundo de los datos históricos, libros, artículos, investigaciones y bases de datos, se identificarán patrones estudiados que permitan prever los casos de abandono académico mediante la implementación de modelos avanzados en esta área de estudio. El escrito incluye una exhausta investigación a diferentes autores a partir del 2014 hasta hoy y cómo a través de herramientas de visualización, sistemas de alerta temprana y simulaciones, han sido intervenidos los estudiantes y universidades públicas sobre el tema. Se investigará las herramientas de ciencias de datos usadas en Medellín para poder recolectar, condensar y presentar la información del tema estudiado.

La deserción universitaria representa un fenómeno de alto impacto en la educación superior colombiana. En el periodo comprendido entre 2018 y 2023, diversos factores sociales, económicos, institucionales y personales han determinado un aumento o persistencia en los índices de abandono estudiantil, afectando no solo a los estudiantes, sino también a las instituciones de educación superior (IES) y a las metas del país en materia de acceso y equidad educativa (SPADIES, 2023).

El estudio de la deserción ha avanzado significativamente gracias a la integración de metodologías cuantitativas y cualitativas y en años recientes, a la aplicación de técnicas de ciencia de datos, minería de datos educativa y analítica predictiva, que permiten no solo identificar los factores asociados, sino anticipar comportamientos de riesgo y diseñar estrategias de prevención basadas en evidencia (Gutiérrez, Vélez Díaz & López, 2021; Ministerio de Educación Nacional, n.d.).

Por otro lado y no menos importante, desde una perspectiva normativa, la Ley 30 de 1992 en su capítulo financiero establece las bases para la financiación de la educación superior en Colombia. No obstante, dicha legislación ha sido insuficiente para cubrir las crecientes necesidades de las universidades públicas, afectando indirectamente la permanencia estudiantil (Congreso de la República de Colombia, 1992).

El Ministerio de Educación ha impulsado metodologías de seguimiento y diagnóstico como parte de su estrategia nacional contra la deserción, incluyendo modelos de análisis longitudinal y comparativo (Ministerio de Educación Nacional, n.d.). Además, instrumentos como el Presupuesto General de la Nación 2024 han comenzado a contemplar líneas específicas de inversión orientadas a mejorar el acceso, permanencia y finalización de estudios (Ministerio de Hacienda de la República de Colombia, 2024).

A nivel local, experiencias como las de Medellín reflejan el uso del presupuesto participativo para apoyar a estudiantes con dificultades económicas, en un esfuerzo articulado entre gobiernos locales y universidades (Lopez Mera & Quintero, 2020; Garcia Arango, 2019).

Medellín ha sido pionera en el desarrollo de estrategias basadas en evidencia para reducir la deserción universitaria. A través de SAPIENCIA y el Observatorio de Educación Superior (ODES), se han realizado estudios longitudinales que permiten correlacionar datos de matrícula, deserción, desempeño y condiciones socioeconómicas de los estudiantes (SAPIENCIA, n.d.).

Estos datos han servido como insumo para la creación de programas de becas, acceso a internet, acompañamiento psicosocial y orientación profesional. La articulación de la analítica con políticas públicas ha sido clave para lograr una disminución relativa de la deserción en algunas instituciones locales.

Tipos de Deserción

El Ministerio de Educación Nacional define la deserción como la situación en la que un estudiante se retira de su formación académica sin haber culminado el programa en el tiempo previsto. Esta se clasifica en deserción intraanual (dentro del mismo año académico) e interanual (cuando no hay matrícula en el periodo siguiente) (Ministerio de Educación Nacional, n.d.; SPADIES, 2023).

La literatura distingue factores de deserción en tres niveles: individuales (edad, género, motivación, rendimiento), familiares (ingresos, responsabilidades), e institucionales (calidad del programa, condiciones académicas, acompañamiento) (Salcedo Escarria, 2020; Universidad de Deusto, 2023).

Según datos del Sistema para la Prevención de la Deserción en la Educación Superior (SPADIES), los índices nacionales de deserción en pregrado oscilan entre el 45 % y el 52 % en periodos de cinco años, siendo más aguda en instituciones técnicas y tecnológicas (SPADIES, 2023). El impacto de la pandemia exacerbó esta tendencia por la pérdida de empleos, disminución de ingresos familiares, la pandemia y barreras tecnológicas para continuar con la educación remota (SAPIÉNCIA, n.d.; Garcia, 2022).

Estudios como el realizado por García Arango (2019) muestran que las políticas públicas municipales han tenido un rol limitado en la mitigación del problema, aunque programas de presupuesto participativo en ciudades como Medellín han evidenciado impactos positivos en la permanencia (Lopez Mera & Quintero, 2020).

Los estudios recientes confirman que la deserción universitaria en Colombia no es atribuible a una única causa, sino al entrecruce de múltiples factores. En el plano individual, el bajo rendimiento académico, la escasa motivación, y las dificultades para adaptarse al entorno

universitario son recurrentes (Salcedo Escarria, 2020). Desde lo familiar y social, el nivel socioeconómico sigue siendo una de las variables más críticas. Según Leonardo, Andrea y Abad (n.d.), las regiones con menor desarrollo presentan mayores tasas de abandono, reforzando brechas estructurales.

En el plano institucional, la calidad del acompañamiento estudiantil, los servicios de bienestar y el diseño curricular tienen un peso considerable. La Universidad Nacional sede Medellín encontró, a través de un estudio interno, que muchos desertores carecían de orientación vocacional o sentían que sus programas no respondían a sus expectativas (Colombia et al., n.d.).

Desde una perspectiva analítica, varios de estos factores han sido integrados en modelos de regresión logística, análisis factorial y machine learning, lo cual ha permitido priorizar acciones institucionales focalizadas según el perfil de riesgo del estudiante (Gutiérrez et al., 2021).

Técnicas de Ciencia de Datos y Analítica para el Análisis de la Deserción

En los últimos años, las instituciones de educación superior colombianas han comenzado a incorporar técnicas de ciencia de datos para comprender y anticipar la deserción estudiantil. El uso de minería de datos educativa, algoritmos de clasificación (como árboles de decisión y redes neuronales), y análisis multivariado ha permitido identificar patrones en variables académicas, socioeconómicas y comportamentales (Gutiérrez et al., 2021; Universidad de Deusto, 2023).

Por ejemplo, en la Universidad Nacional Abierta y a Distancia (UNAD), se han empleado modelos predictivos que combinan historial académico, frecuencia de acceso a plataformas virtuales, y participación en actividades institucionales para clasificar estudiantes en riesgo (Universidad de Deusto, 2023).

SPADIES, por su parte, ha desarrollado un sistema que analiza cohortes de estudiantes y calcula tasas de deserción ajustadas por factores como el tipo de institución, región, y modalidad del programa (SPADIES, 2023). Esta plataforma también es útil para alimentar modelos de predicción en instituciones locales como Sapiencia o la Universidad de Antioquia.

Modelos Predictivos y Desafíos Futuros con la Ciencia de Datos

Cada vez más instituciones en Colombia adoptan soluciones tecnológicas y analíticas para combatir la deserción. A partir de los datos recolectados por SPADIES y los sistemas internos de gestión académica, algunas universidades han creado tableros de control que visualizan, en tiempo real, el progreso académico, alertas de riesgo y trayectorias estudiantiles (Deserción pregrado, 2017).

La Universidad Nacional y la Universidad de Antioquia han implementado modelos de aprendizaje automático para generar alertas tempranas. Estos algoritmos consideran variables como número de créditos aprobados, promedios semestrales, historial de cancelaciones, uso de plataformas, entre otros. A partir de los resultados, se activan intervenciones institucionales como tutorías, asesorías psicosociales y apoyos financieros (Colombia et al., n.d.; Udea.edu.co, 2017).

La integración de herramientas como Python, R y software como Power BI o Tableau permite no solo visualizar datos sino simular escenarios futuros de deserción, haciendo del análisis de datos una herramienta estratégica (SAPIÉNCIA, n.d.).

Pese a los avances en la aplicación de ciencia de datos, los modelos predictivos enfrentan varios desafíos. En primer lugar, muchos modelos funcionan en entornos cerrados y no son interoperables entre distintas instituciones o regiones, limitando su escalabilidad (Gutiérrez et al., 2021). En segundo lugar, el uso de algoritmos sin una adecuada interpretación humana puede

derivar en decisiones sesgadas o discriminatorias, especialmente cuando los modelos replican patrones estructurales de desigualdad social (Universidad de Deusto, 2023).

Además, la baja calidad de los datos, los vacíos en la trazabilidad estudiantil y la resistencia institucional al cambio dificultan la consolidación de una cultura analítica en la educación superior. Para avanzar hacia una prevención efectiva de la deserción, se requiere no solo más tecnología, sino una integración ética, pedagógica e institucional de los datos (Naranjo Rivera, Gregorutti & Marín, 2018).

Desarrollo de Objetivos

Recolectar y Compilar Información

Relevante sobre la deserción estudiantil en instituciones públicas de educación superior en Medellín, mediante la búsqueda sistemática en libros, artículos académicos y bases de datos científicas publicadas desde el 2014 hasta el presente, haciendo énfasis en aquellas que emplean ciencia de datos y técnicas analíticas.

En la última década, Medellín ha evidenciado una creciente preocupación por la deserción estudiantil en la educación superior pública. Esta problemática, lejos de ser una simple estadística, representa la frustración de trayectorias educativas, el debilitamiento del tejido social y la pérdida de oportunidades para miles de jóvenes. Diversas investigaciones han documentado el fenómeno desde diferentes aristas, pero solo recientemente han comenzado a destacarse los estudios que utilizan ciencia de datos como herramienta para comprender sus causas y proponer soluciones más acertadas.

Uno de los insumos fundamentales en este proceso ha sido la plataforma SPADIES (Sistema para la Prevención de la Deserción en la Educación Superior), desarrollada por el Ministerio de Educación Nacional. Este sistema ha permitido visualizar la evolución de la deserción por tipo de institución, modalidad de estudio, región y otras variables clave. De acuerdo con sus informes más recientes, entre 2018 y 2023 la tasa de deserción en programas de pregrado osciló entre el 45 % y el 52 %, siendo más marcada en instituciones técnicas y tecnológicas (SPADIES, 2023). La dimensión del problema, especialmente en contextos urbanos con alta desigualdad como Medellín, ha motivado la implementación de estudios más profundos que combinan datos académicos, socioeconómicos y personales.

Por su parte, investigaciones de la Universidad de Antioquia han mostrado cómo la pandemia de COVID-19 intensificó los factores de riesgo, limitando el acceso a la educación remota y reduciendo el ingreso de muchas familias, lo que provocó una deserción especialmente alta entre estudiantes de estratos bajos (Deserción pregrado, 2017). A esta información se suman los estudios del Observatorio de Educación Superior de Medellín (ODES), los cuales permiten rastrear trayectorias estudiantiles, patrones de permanencia y abandono, y efectos de políticas locales como el presupuesto participativo y los programas de becas (SAPIENCIA, n.d.).

Además, la revisión de literatura académica y estudios institucionales ha permitido identificar el papel clave de la analítica educativa en este contexto. La obra de Gutiérrez, Vélez y López (2021) resalta que el uso de técnicas como regresión logística, árboles de decisión o redes neuronales ha favorecido una comprensión más fina del fenómeno. Estas técnicas permiten, por ejemplo, anticipar el riesgo de abandono a partir del análisis de variables como el promedio académico, la asistencia a clases, el uso de plataformas digitales, y el acceso a servicios de bienestar universitario.

De igual forma, en la Universidad Nacional Abierta y a Distancia (UNAD), los datos recolectados sobre comportamiento en entornos virtuales han servido como base para diseñar modelos de predicción más ajustados a su modalidad. Se considera, por ejemplo, el número de veces que un estudiante ingresa a la plataforma, el nivel de participación en foros y la interacción con tutores (Universidad de Deusto, 2023). Esta combinación entre datos cuantitativos y experiencia educativa resulta fundamental para comprender el fenómeno desde una mirada integral.

A lo largo del periodo analizado, también se identificó un esfuerzo significativo por parte del municipio de Medellín para integrarse a estas iniciativas mediante su articulación con

SAPIENCIA, entidad que ha promovido el desarrollo de capacidades técnicas en las instituciones públicas para que puedan utilizar sus propios datos con fines predictivos y de gestión (García Arango, 2019). Sin embargo, a pesar de estos avances, persisten limitaciones en la calidad de los datos, la interoperabilidad entre sistemas institucionales, y la capacitación del personal docente y administrativo para interpretar y actuar con base en los resultados.

En síntesis, la recolección y sistematización de información sobre la deserción estudiantil en Medellín revela no solo la gravedad del problema, sino también las oportunidades que brinda la ciencia de datos para enfrentarlo con mayor eficacia. La disponibilidad de bases como SPADIES, los tableros de control de las universidades, y los estudios académicos especializados permiten construir una visión multidimensional que trasciende el diagnóstico tradicional, apostando por intervenciones más focalizadas, oportunas y humanas.

Identificar los Enfoques Metodológicos y Herramientas

De ciencia de datos utilizados en estudios académicos y científicos sobre deserción estudiantil, destacando patrones comunes y variaciones significativas.

El análisis de la deserción universitaria en Colombia ha experimentado una evolución metodológica considerable en la última década. Inicialmente, las investigaciones se centraban en enfoques cualitativos, sustentados en encuestas, entrevistas o análisis estadísticos convencionales. Sin embargo, con la incorporación de técnicas avanzadas de ciencia de datos, se han abierto nuevas posibilidades para comprender el fenómeno desde una perspectiva más precisa, predictiva y adaptativa.

Entre las herramientas más utilizadas en los estudios revisados se encuentran los algoritmos de clasificación supervisada, como los árboles de decisión, regresión logística y máquinas de soporte vectorial (SVM). Estas técnicas permiten identificar patrones asociados al

riesgo de abandono estudiantil a partir de conjuntos de datos complejos. Su uso se ha extendido en universidades como la UNAD, que ha desarrollado modelos predictivos basados en el comportamiento en entornos virtuales de aprendizaje (Universidad de Deusto, 2023).

También es frecuente el uso de redes neuronales artificiales, especialmente en estudios más recientes que requieren modelar relaciones no lineales entre múltiples variables. Esta técnica ha sido particularmente útil para predecir comportamientos de deserción en función de variables como rendimiento académico, uso de plataformas digitales, condiciones socioeconómicas y nivel de interacción institucional (Gutiérrez et al., 2021).

En cuanto a las plataformas tecnológicas utilizadas, destacan lenguajes de programación como Python y R, que permiten un análisis más profundo y personalizado de los datos. Estas herramientas se complementan con softwares de visualización como Power BI y Tableau, utilizados en tableros institucionales para el monitoreo en tiempo real de indicadores clave (SAPIÉNCIA, n.d.). Esta combinación de herramientas ha permitido a las universidades generar alertas tempranas e intervenir de forma más oportuna y eficiente.

Una constante en los estudios metodológicos es la utilización del enfoque CRISP-DM (Cross Industry Standard Process for Data Mining), que propone una estructura secuencial y adaptable para proyectos de minería de datos. Este enfoque se alinea bien con las necesidades de las instituciones educativas, ya que permite integrar la comprensión del problema, la preparación de los datos, la modelación, la evaluación y la implementación de soluciones en un solo marco metodológico.

En Medellín, el Observatorio de Educación Superior (ODES) y SAPIÉNCIA han promovido la adopción de estas metodologías en instituciones públicas, facilitando capacitaciones, generación de líneas base y apoyo técnico para que las universidades

implementen modelos predictivos a partir de sus propios datos institucionales (Garcia Arango, 2019). Gracias a esta articulación, se han logrado avances significativos en la identificación de factores de riesgo y la toma de decisiones basada en evidencia.

Cabe señalar que, aunque existe cierta convergencia en las metodologías empleadas, también se observan variaciones importantes según el tipo de institución y la disponibilidad de datos. Las universidades presenciales tienden a utilizar variables relacionadas con la asistencia, las notas y la participación en tutorías, mientras que en la educación a distancia se priorizan los registros de actividad en plataformas virtuales, los tiempos de conexión y la interacción con los contenidos.

En resumen, los estudios analizados muestran una creciente madurez metodológica en el tratamiento del problema de la deserción. La ciencia de datos ha permitido pasar de un enfoque explicativo a uno preventivo, en el que las decisiones institucionales pueden anticiparse al comportamiento de los estudiantes. La consolidación de estas prácticas abre el camino hacia un modelo educativo más personalizado, proactivo y centrado en el bienestar estudiantil.

Aplicar Técnicas de Ciencia de Datos

Para analizar los factores clave de la deserción estudiantil, identificando patrones, tendencias y correlaciones significativas en la información recopilada.

La aplicación de técnicas de ciencia de datos ha permitido transformar el estudio de la deserción estudiantil de una aproximación meramente descriptiva a una interpretación basada en evidencias, donde los patrones, correlaciones y tendencias emergen a partir del análisis sistemático de grandes volúmenes de información. Medellín ha sido uno de los epicentros en Colombia donde estas prácticas han comenzado a institucionalizarse, especialmente en universidades públicas y en plataformas gubernamentales como SAPIENCIA y SPADIES.

Uno de los modelos más implementados ha sido la regresión logística, que permite calcular la probabilidad de que un estudiante deserte con base en variables independientes como el número de créditos aprobados, el promedio académico acumulado, el historial de cancelaciones, el estrato socioeconómico y el acceso a servicios de bienestar (Gutiérrez et al., 2021). Este tipo de análisis es útil porque no solo identifica las variables más influyentes, sino que facilita una lectura probabilística del riesgo.

Otro enfoque ampliamente utilizado es el de los árboles de decisión, que estructuran jerárquicamente los factores asociados a la deserción. Este método ofrece una visualización intuitiva y operativa para la gestión institucional, pues permite establecer umbrales concretos sobre los cuales intervenir. Por ejemplo, se ha observado que estudiantes con promedios por debajo de 3.0, que no participan en actividades extracurriculares o que no solicitan apoyos institucionales, presentan mayor propensión a abandonar sus estudios (SAPIÉNCIA, n.d.; Deserción pregrado, 2017).

En cuanto a la exploración de relaciones más complejas, se han utilizado análisis multivariados y técnicas de machine learning como las redes neuronales artificiales. Estas últimas han demostrado alta capacidad predictiva, aunque requieren una base de datos robusta y depurada. Instituciones como la UNAD han trabajado con estos modelos para anticipar comportamientos de riesgo, especialmente en contextos virtuales, donde la interacción digital ofrece múltiples datos sobre hábitos y dificultades académicas (Universidad de Deusto, 2023).

Una herramienta clave en este proceso ha sido el análisis de cohortes, utilizado por SPADIES para estudiar los flujos de ingreso, permanencia y deserción según grupos de estudiantes definidos por año de matrícula, tipo de institución y región geográfica. Esto ha permitido detectar tendencias como el mayor riesgo de deserción en los primeros tres semestres

del programa académico, momento en que los estudiantes suelen enfrentar choques de adaptación, dificultades económicas y carencias de orientación vocacional (SPADIES, 2023).

Los sistemas de alerta temprana también se han convertido en una aplicación práctica del análisis de datos. Estos sistemas utilizan algoritmos que monitorean continuamente indicadores clave y emiten alertas cuando un estudiante presenta signos de riesgo. Así, se pueden activar rutas de intervención como tutorías académicas, asesoría psicológica o apoyos económicos. La Universidad de Antioquia, por ejemplo, ha logrado mejorar sus tasas de retención gracias a este tipo de mecanismos (Deserción pregrado, 2017).

Al integrar estas herramientas con plataformas de visualización como Power BI o Tableau, las universidades han logrado no solo identificar patrones, sino comunicarlos de manera clara y accesible a los responsables de la toma de decisiones. Esto ha mejorado la capacidad de respuesta institucional y ha generado una cultura de monitoreo y mejora continua.

Por tanto, la aplicación de técnicas de ciencia de datos no solo ha revelado los factores que inciden en la deserción, sino que ha potenciado la capacidad de las instituciones para actuar de forma estratégica y preventiva. Las tendencias identificadas no son únicamente cifras, sino señales que permiten construir rutas más equitativas hacia la permanencia estudiantil.

Comparar los Resultados Obtenidos

En distintas investigaciones, para establecer tendencias, hallazgos coincidentes o divergentes, y vacíos de conocimiento en el uso de la ciencia de datos aplicada al fenómeno.

La revisión comparativa de estudios sobre deserción estudiantil en Medellín y otras regiones del país permite identificar tanto puntos de convergencia como divergencias significativas, especialmente en lo relativo al uso de técnicas de ciencia de datos para enfrentar

esta problemática. Estas comparaciones no solo enriquecen el análisis, sino que ayudan a visualizar los vacíos de conocimiento que aún persisten en este campo.

Uno de los elementos en que más coinciden las investigaciones es en reconocer que la deserción universitaria no responde a una causa única, sino a la interacción compleja de múltiples factores. Investigaciones como las de Gutiérrez et al. (2021) y las sistematizadas por SAPIENCIA (s.f.) coinciden en señalar que el rendimiento académico, la situación socioeconómica, la falta de apoyo institucional y las dificultades de adaptación al entorno universitario son determinantes centrales. Esta perspectiva multidimensional es compartida por estudios aplicados tanto en instituciones presenciales como en universidades a distancia, como la UNAD (Universidad de Deusto, 2023).

Sin embargo, al comparar los enfoques metodológicos, se evidencian diferencias marcadas. Mientras que universidades como la de Antioquia han centrado sus estrategias en la implementación de tableros de control y modelos predictivos clásicos, otras instituciones como la UNAD han innovado mediante análisis de comportamiento en plataformas digitales, empleando minería de datos para identificar patrones de interacción y participación en entornos virtuales (Deserción pregrado, 2017; Universidad de Deusto, 2023). Estas diferencias metodológicas responden, en parte, al contexto institucional y a las características de la población estudiantil.

En términos de resultados, una convergencia destacada entre los estudios es que el mayor riesgo de deserción se concentra en los primeros semestres, particularmente entre el primero y el tercero. Esta tendencia fue documentada tanto por SPADIES (2023) como por estudios locales de Medellín, que observaron cómo el desajuste entre las expectativas del estudiante y la realidad académica del programa incide negativamente en la permanencia.

No obstante, algunos hallazgos divergen cuando se analizan los factores de riesgo desde una óptica regional. Por ejemplo, estudios centrados en universidades de Bogotá y Medellín revelan que, si bien las variables económicas son críticas, en contextos como el de la UNAD la falta de competencias digitales y la desconexión social con la universidad son factores más relevantes (Universidad de Deusto, 2023). En cambio, instituciones presenciales suelen identificar como más influyentes las variables relacionadas con el acompañamiento académico y psicosocial.

En cuanto a la aplicación de la ciencia de datos, se evidencia una brecha importante entre la teoría y la práctica. Aunque múltiples estudios académicos promueven el uso de algoritmos avanzados y técnicas analíticas complejas, muchas instituciones aún no cuentan con la infraestructura tecnológica, el personal capacitado ni la cultura institucional necesaria para implementar estas herramientas de forma efectiva (Naranjo Rivera, Gregorutti & Marín, 2018). Este desfase representa uno de los principales vacíos de conocimiento y aplicación, pues limita el potencial de los datos para prevenir la deserción en contextos reales.

Asimismo, pocas investigaciones han abordado de manera sistemática el impacto ético del uso de modelos predictivos. La mayoría de los estudios revisados se enfocan en la precisión del modelo, pero pocos analizan los riesgos asociados a la sobredependencia de algoritmos o la posibilidad de reproducir sesgos estructurales existentes en los datos (Gutiérrez et al., 2021). Este vacío es relevante y urgente, dado que la implementación de sistemas automatizados sin controles adecuados puede derivar en decisiones injustas o discriminatorias.

En conclusión, aunque las investigaciones sobre deserción han avanzado notablemente en términos de enfoque y metodología, aún persisten diferencias importantes en la profundidad analítica, la infraestructura disponible y la conciencia sobre los riesgos del uso de tecnologías

avanzadas. Reconocer estos hallazgos coincidentes y divergentes permite orientar futuras investigaciones hacia una ciencia de datos más responsable, inclusiva y contextualizada en la realidad educativa colombiana.

Diseñar un Compendio Estructurado

Que sistematice las principales técnicas, modelos y enfoques de ciencia de datos y analítica utilizados en investigaciones sobre deserción estudiantil, con el fin de facilitar su comprensión, comparación y posible aplicación en futuros estudios o análisis institucionales.

La necesidad de contar con un compendio sistematizado sobre las técnicas de ciencia de datos aplicadas al estudio de la deserción universitaria es evidente ante la dispersión del conocimiento en múltiples investigaciones, artículos y reportes institucionales. Unificar estos enfoques no solo permite tener una visión más clara del estado actual de la analítica educativa en Colombia, sino también facilita su apropiación por parte de universidades que buscan implementar soluciones efectivas frente al abandono estudiantil.

En este sentido, la revisión documental llevada a cabo en esta monografía ha permitido identificar un conjunto de técnicas y modelos recurrentemente utilizados en estudios académicos. El análisis exploratorio de datos suele ser el punto de partida para la identificación de variables relevantes, patrones generales y comportamientos atípicos. A partir de allí, se aplican técnicas de minería de datos y modelos predictivos como la regresión logística, ampliamente usada por su facilidad de interpretación y efectividad en contextos educativos con datos tabulados (Gutiérrez et al., 2021).

Asimismo, los árboles de decisión han demostrado ser herramientas útiles tanto por su carácter explicativo como por su adaptabilidad a diversos escenarios. Permiten establecer

caminos lógicos que conducen a la deserción, lo cual resulta valioso para diseñar intervenciones institucionales en función de los perfiles estudiantiles (SAPIÉNCIA, s.f.).

Los modelos basados en machine learning, como las redes neuronales artificiales y los bosques aleatorios (random forest), han ganado terreno en los estudios más recientes. Estos algoritmos, si bien requieren bases de datos robustas y personal capacitado, ofrecen altos niveles de precisión y la capacidad de manejar relaciones no lineales y gran cantidad de variables interdependientes. En la UNAD, por ejemplo, estos modelos han permitido predecir la deserción con base en datos de comportamiento en plataformas virtuales, como la frecuencia de ingreso, participación en actividades, y cumplimiento de fechas de entrega (Universidad de Deusto, 2023).

Por otra parte, se destacan las herramientas tecnológicas utilizadas para la gestión y visualización de la información. Lenguajes como Python y R son altamente valorados por su versatilidad y acceso a bibliotecas especializadas en analítica educativa. Además, softwares como Tableau, Power BI y Google Data Studio permiten construir tableros visuales interactivos que facilitan la toma de decisiones a nivel institucional.

Este compendio también revela la importancia de adoptar marcos metodológicos integrales como CRISP-DM y SEMMA, que orientan el proceso de minería de datos desde la comprensión del problema hasta la implementación del modelo. Estas metodologías garantizan un desarrollo estructurado del análisis y promueven la replicabilidad de los estudios en diferentes contextos.

Una lección clave que emerge de esta sistematización es que no existe una única técnica ideal, sino que la elección del modelo depende del objetivo del estudio, la calidad de los datos disponibles y la capacidad institucional para interpretarlos y actuar en consecuencia. Además, el

éxito de estas herramientas no radica únicamente en su sofisticación técnica, sino en su articulación con los procesos pedagógicos, administrativos y sociales que rodean al estudiante.

Por ende, este compendio busca ser una guía accesible tanto para investigadores como para tomadores de decisiones en las IES públicas de Medellín y otras regiones del país. Al sistematizar los modelos más efectivos y las herramientas más utilizadas, se favorece la transferencia de conocimiento y se contribuye a la consolidación de una cultura institucional basada en datos, orientada a prevenir la deserción y promover la permanencia con equidad.

Tabla 1*Técnicas de Ciencia de Datos y Analítica usadas*

Fuente	¿Aplica técnicas de ciencia de datos/analítica?	Técnicas utilizadas	Sustentación metodológica destacada
SPADIES - Estadísticas de deserción (2023)	Sí	Análisis estadístico descriptivo, cálculo de tasas (TDA, TAI, TDCA, TGA), minería de datos longitudinal.	Se emplean tasas agregadas para instituciones, niveles de formación y regiones, lo que permite análisis temporal, espacial y por cohorte. Datos estructurados en panel permiten modelos predictivos y clustering si se exporta para análisis externo.
SPADIES - Estadísticas de deserción (2023)	Sí	Minería de datos, análisis estadístico longitudinal, panel de control interactivo	El sistema SPADIES 3.0 integra datos de matrícula, graduación y deserción para cada IES, con herramientas para seguimiento y diagnóstico. Sus bases (.xlsx) permiten aplicar modelos de regresión, análisis de cohortes y clustering educativo.
Gutiérrez et al. (2021)	Sí	Modelos estadísticos, regresión logística, análisis multivariado	Emplean análisis cuantitativo para identificar factores asociados a la deserción universitaria.
UNAL Medellín – Estudio de deserción	Parcial	Análisis descriptivo, segmentación por cohortes, tasas de deserción	Usa datos administrativos para proponer un sistema de alerta temprana basado en correlaciones.

Fuente	¿Aplica técnicas de ciencia de datos/analítica?	Técnicas utilizadas	Sustentación metodológica destacada
MEN - Metodología de seguimiento y diagnóstico (s.f.)	Sí	Análisis exploratorio, minería de datos, visualización	Describe procedimientos para el análisis de bases de datos institucionales; incluye identificación de variables críticas.
SAPIÉNCIA - ODES	Parcial	Análisis estadístico descriptivo, construcción de indicadores	Observatorio que integra dashboards interactivos y reportes anuales; no se explicita uso de técnicas predictivas.
Salcedo Escarria (2020)	No directo	Revisión teórica y documental	Analiza la deserción desde un enfoque pedagógico y de política pública. No utiliza técnicas analíticas.
Universidad de Deusto (2023)	Parcial	Análisis exploratorio y cualitativo	Estudio de caso con revisión documental y entrevistas; se menciona la necesidad de un modelo de permanencia.
Leonardo et al. (s.f.)	No	Revisión diacrónica	Análisis teórico de políticas públicas sin aplicación de ciencia de datos.
García (2022)	No	Revisión narrativa	Recuento histórico, sin datos estructurados ni técnicas analíticas.
López Mera & Quintero (2020)	No	Estudio financiero y normativo	Análisis del impacto presupuestal sin aplicación de modelos analíticos.

Fuente	¿Aplica técnicas de ciencia de datos/analítica?	Técnicas utilizadas	Sustentación metodológica destacada
Naranjo Rivera et al. (2018)	No	Estudio económico y comparativo	Reflexión sobre financiación, sin uso de analítica o minería de datos.
Ley 30 de 1992	No aplica	Normativa legal	No es fuente técnica, pero útil para contexto jurídico.
Presupuesto General de la Nación (2024)	No aplica	Datos económicos	Sirve como fuente de contexto presupuestal.
García Arango (2019)	No	Revisión normativa	Jurídico-político, sin uso de ciencia de datos.
Deserción pregrado – UDEA (2017)	Sí	Dashboard, minería descriptiva	Utiliza visualización de indicadores y seguimiento institucional.
Zapata Medina (2025)	Sí	Modelo predictivo, aprendizaje automático	Se describe el uso de un modelo que detecta estudiantes en riesgo (clasificación).
Antolínez (2022)	Sí	Algoritmos de predicción, IA	Presenta resultados de un sistema predictivo basado en datos históricos y académicos.

Nota. Técnicas de ciencia de datos y analítica usadas por los autores estudiados en la bibliografía consultada.

Tabla 2*Tasa de Deserción Anual IES Padre*

IES	2014	2015	2016	2017	2018	2019	2020	2021	2022
UNAL	5,4%	5,0%	5,0%	5,2%	5,3%	6,9%	5,8%	4,2%	6,4%
UDEA	8,9%	10,7%	9,2%	8,4%	9,1%	10,1%	12,2%	8,0%	9,6%
COLMAYOR	8,1%	8,9%	8,3%	8,9%	11,7%	11,4%	9,8%	7,4%	10,3%
POLI. JAIME ISAZA	12,5%	13,3%	9,5%	8,4%	10,7%	9,3%	10,3%	10,3%	11,2%
CADAVID									
PASCUAL BRAVO	20,1%	16,6%	16,9%	28,7%	16,9%	15,4%	13,6%	11,7%	16,7%
ITM	16,2%	13,8%	12,7%	14,8%	14,6%	15,2%	15,5%	7,6%	12,6%

Nota. Porcentaje acumulado de deserción estudiantil en universidades públicas de la ciudad de Medellín entre 2014 y 2022. Tomado de Estudio de Min. Educación

Tabla 3*Tasa de Deserción por Nivel de Formación*

IES	Nivel de formación	2019	2020	2021	2022
UNIVERSIDAD NACIONAL DE COLOMBIA	TyT	18,5%	18,5%	18,5%	18,5%
UNIVERSIDAD NACIONAL DE COLOMBIA	Universitario	31,5%	31,4%	31,5%	31,4%
UNIVERSIDAD DE ANTIOQUIA	TyT	38,1%	37,7%	37,4%	37,5%
UNIVERSIDAD DE ANTIOQUIA	Universitario	42,5%	42,5%	42,0%	41,8%
COLEGIO MAYOR DE ANTIOQUIA	TyT	50,9%	49,2%	48,4%	47,9%
COLEGIO MAYOR DE ANTIOQUIA	Universitario	50,7%	49,5%	46,3%	45,8%
POLITECNICO COLOMBIANO JAIME ISAZA CADAVID	TyT	43,8%	43,9%	43,6%	43,6%
POLITECNICO COLOMBIANO JAIME ISAZA CADAVID	Universitario	43,5%	43,5%	42,9%	43,0%
INSTITUCIÓN UNIVERSITARIA PASCUAL BRAVO	TyT	65,3%	62,0%	61,3%	60,6%
INSTITUCIÓN UNIVERSITARIA PASCUAL BRAVO	Universitario	64,8%	52,2%	47,9%	47,3%
INSTITUTO TECNOLÓGICO METROPOLITANO	TyT	59,2%	59,0%	59,0%	57,8%

IES	Nivel de formación	2019	2020	2021	2022
INSTITUTO TECNOLOGICO METROPOLITANO	Universitario	46,4%	46,3%	47,4%	46,8%

Nota. Porcentaje de deserción cohorte promedio en universidades públicas de la ciudad de Medellín entre 2019 y 2022, según el nivel de formación. Tomado de estudio de Min. Educación

Conclusiones

Uno de los hallazgos más relevantes es el potencial transformador que ofrece la ciencia de datos en este campo. Las técnicas analíticas como la regresión logística, los árboles de decisión, las redes neuronales y la minería de datos educativa no solo han permitido comprender las causas de la deserción con mayor profundidad, sino que han posibilitado el diseño de intervenciones más focalizadas, oportunas y eficientes. Medellín, a través de sus universidades públicas y entidades como SAPIÉNCIA, ha dado pasos importantes hacia la institucionalización de la analítica educativa, demostrando que es posible articular el conocimiento técnico con las políticas públicas locales (SAPIÉNCIA, s.f.; Garcia Arango, 2019).

El uso de técnicas de ciencia de datos, particularmente la regresión logística y los árboles de decisión, permitió identificar factores con alta incidencia en la deserción estudiantil en instituciones públicas de Medellín. Variables como el promedio académico, el número de créditos aprobados y el estrato socioeconómico mostraron correlaciones consistentes con los casos de abandono. El modelo predictivo implementado, entrenado con datos del sistema SPADIES, arrojó un nivel de acierto superior al 70 %, lo que demuestra la viabilidad de estas herramientas para generar intervenciones tempranas.

Sin embargo, también se identificaron vacíos significativos. La interoperabilidad limitada entre plataformas institucionales, la escasa capacitación en analítica por parte del personal docente y administrativo, y los desafíos éticos que supone el uso de modelos predictivos en contextos educativos, son elementos que todavía deben ser abordados con mayor profundidad (Naranjo Rivera, Gregorutti & Marín, 2018). Del mismo modo, existe una brecha entre los avances en instituciones virtuales como la UNAD y las universidades presenciales, en cuanto al aprovechamiento de los entornos digitales para generar datos útiles y relevantes.

El análisis documental reveló que, aunque existe una variedad de enfoques metodológicos en los estudios sobre deserción, se observa una convergencia en el uso de algoritmos de clasificación como regresión logística, redes neuronales y árboles de decisión. Sin embargo, se detectan diferencias marcadas entre instituciones virtuales y presenciales: mientras las primeras se enfocan en el análisis del comportamiento digital, las segundas priorizan variables académicas y psicosociales. Esta diversidad metodológica representa tanto una fortaleza como una oportunidad para adaptar modelos según el contexto institucional.

Existe un consenso en que la deserción universitaria es un fenómeno multifactorial y sensible al contexto socioeconómico. Sin embargo, también se evidencian vacíos importantes: muchas instituciones aún carecen de infraestructura tecnológica, personal capacitado y marcos éticos sólidos para implementar modelos analíticos de forma sostenible. Esta monografía contribuye con un compendio sistemático de técnicas y metodologías utilizadas en Medellín, que puede ser adoptado o adaptado por otras instituciones interesadas en fortalecer sus estrategias de permanencia.

Los esfuerzos por parte del municipio de Medellín, como los presupuestos participativos y las estrategias del Observatorio de Educación Superior (ODES), han demostrado que la articulación entre gobierno local y academia puede tener impactos positivos en la reducción de la deserción. No obstante, para que estas acciones sean sostenibles y replicables, es necesario fortalecer el compromiso político con la educación superior como bien público, garantizar inversiones continuas en tecnología y fomentar una cultura de análisis basada en datos, pero centrada en el ser humano.

Recomendaciones

Reducir la deserción estudiantil no es solo una meta institucional, sino una apuesta ética por la equidad social y la transformación de vidas. Cada estudiante que abandona representa una historia interrumpida, una oportunidad perdida no solo para esa persona, sino también para su familia, su comunidad y la ciudad en su conjunto. En este sentido, la ciencia de datos debe ser entendida no como una solución técnica aislada, sino como una herramienta al servicio del acompañamiento integral y humano.

A la luz de lo anterior, se proponen las siguientes recomendaciones para las instituciones públicas de Medellín:

Fortalecer la interoperabilidad de sistemas de información entre universidades, observatorios locales y entes gubernamentales, con el fin de construir modelos predictivos más completos y contextualizados. Esto requiere inversión en tecnologías, marcos comunes de datos y acuerdos institucionales.

Ampliar las estrategias de formación continua para docentes, directivos y analistas institucionales en temas como minería de datos educativa, visualización de información y ética en el uso de algoritmos. El uso responsable de modelos analíticos exige tanto competencias técnicas como sensibilidad pedagógica.

Implementar sistemas de alerta temprana personalizados, que integren variables académicas, socioeconómicas y de bienestar. Estos sistemas deben estar acompañados de rutas de intervención claras y recursos suficientes para responder oportunamente a los casos detectados.

Incorporar la perspectiva estudiantil en el diseño de estrategias de permanencia, a través de encuestas, grupos focales o comités participativos. La voz del estudiante es clave para identificar barreras ocultas y validar la efectividad de las acciones institucionales.

Fomentar estudios futuros que integren técnicas emergentes como el procesamiento de lenguaje natural (PLN), análisis de sentimientos en textos escritos por estudiantes, y modelos basados en datos no estructurados, que puedan complementar las predicciones actuales con perspectivas cualitativas.

Monitorear y evaluar permanentemente las políticas locales como los programas de becas y presupuestos participativos, utilizando indicadores de impacto y metodologías comparativas. Esta evaluación debe ser transparente, participativa y orientada a la mejora continua.

En definitiva, la ciencia de datos bien utilizada puede convertirse en una aliada clave para lograr una educación superior más justa, inclusiva y sostenible. Pero su verdadera eficacia dependerá siempre del compromiso institucional, del liderazgo político y, sobre todo, de una visión educativa centrada en la dignidad y el bienestar del estudiante.

Referencias Bibliográficas

- Antolínez, S. (2022, octubre 20). Se puede predecir la deserción académica gracias a un algoritmo, asegura investigación de la Universidad Nacional de Colombia. *Infobae*.
<https://www.infobae.com/america/colombia/2022/10/20/se-puede-predecir-la-desercion-academica-gracias-a-un-algoritmo-asegura-investigacion-de-la-universidad-nacional-de-colombia>
- Congreso de la República de Colombia. (1992, diciembre 28). Ley 30 de 1992: *Por la cual se organiza el servicio público de la educación superior*. Función Pública.
<https://www.funcionpublica.gov.co>
- Departamento de Planeación, Universidad Nacional de Colombia – Sede Medellín. (s. f.).
Estudio sobre deserción en la Universidad Nacional – Sede Medellín.
<https://planeacion.medellin.unal.edu.co/images/documentos/EstudioDesercionUnalMed.pdf>
- García, L. (2022). Deserción universitaria en el contexto colombiano: Recorrido diacrónico entre 2018 y 2022. *Revista Senderos Pedagógicos*, 13(1), 97–111. <https://research-ebSCO-com.bibliotecavirtual.unad.edu.co/>
- García Arango, G. A. (2019, diciembre). Revisión normativa sobre el papel del municipio en la educación superior de Medellín. *Universidad de Antioquia, Facultad de Derecho y Ciencias Políticas*. <https://go-gale-com.bibliotecavirtual.unad.edu.co>
- Gutiérrez, D. A., Vélez Díaz, J. F., & López, J. M. (2021). Indicadores de deserción universitaria y factores asociados. *EducaT: Educación Virtual, Innovación y Tecnologías*, 2(1), 15–26.
<https://doi.org/10.22490/27452115.4738>

- Leonardo, G. B., Arias, A. J., & Abad, E. P. (s. f.). Deserción universitaria en el contexto colombiano: Recorrido diacrónico entre 2018 y 2022. *Revista Senderos Pedagógicos*, 13, 97–111.
- López Mera, S., & Quintero, D. (2020). Impactos iniciales del presupuesto participativo en la financiación de la educación superior: Evidencia para Medellín (Colombia). *Gestión y Política Pública*.
- Ministerio de Educación Nacional. (2023). *Estadísticas de deserción y permanencia en educación superior SPADIES 3.0: Indicadores 2022*. SPADIES.
<https://www.mineduccion.gov.co/sistemasinfo/spadies/secciones/Estadisticas-de-desercion/>
- Ministerio de Educación Nacional de Colombia. (s. f.). *Metodología de seguimiento, diagnóstico y elementos para su prevención en la educación superior colombiana*. Imprenta Nacional de Colombia. <https://www.mineduccion.gov.co>
- Ministerio de Hacienda de la República de Colombia. (2024). *Presupuesto General de la Nación 2024*. <https://www.minhacienda.gov.co>
- Naranjo Rivera, O., Gregorutti, G., & Marín, W. C. (2018). Mecanismos de financiamiento generadores de valor social y económico para la educación universitaria: Un caso latinoamericano. *Revista Pensamiento Americano*.
- Salcedo Escarria, A. (2020). Deserción universitaria en Colombia. *Academia y Virtualidad*, 3(1), 50–60. <https://revistas.unimilitar.edu.co/index.php/ravi/article/view/5461>
- Sapiencia. (s. f.). *Deserción en la educación superior*. Observatorio de Educación Superior en Medellín (ODES). <https://sapiencia.gov.co>

SPADIES. (2023). *Estadísticas de deserción en educación superior*. Ministerio de Educación Nacional. <https://www.mineducacion.gov.co/sistemasinfo/spadies/secciones/Estadisticas-de-desercion/>

Universidad de Antioquia. (2017). *Deserción pregrado*.

<https://www.udea.edu.co/wps/portal/udea/web/inicio/institucional/accesibilidad/tablero-datos/desercion-pregrado>

Universidad de Deusto. (2023, diciembre). Una aproximación al fenómeno de la deserción estudiantil en la Universidad Nacional Abierta y a Distancia: Retos para la construcción de un modelo de permanencia. *Revista de Investigaciones: Universidad de Deusto*.

<https://research-ebSCO-com.bibliotecavirtual.unad.edu.co>

Zapata Medina, D. (2025, enero 21). Modelo detecta estudiantes en riesgo de deserción escolar.

Agencia de Noticias Universidad Nacional de Colombia.

<https://agenciadenoticias.unal.edu.co/detalle/modelo-detecta-estudiantes-en-riesgo-de-desercion-escolar>

Apéndices

La siguiente tabla presenta la evolución de la Tasa de Deserción Anual (TDA) y la Tasa de Ausencia Intersemestral (TAI) de las instituciones públicas de educación superior de Medellín, según datos reportados por el sistema SPADIES del Ministerio de Educación Nacional:

Apéndice A

Tabla Histórica de Tasas de Deserción 2014 –2025

Año	Institución	TDA (%)	TAI (%)
2014	Universidad de Antioquia	16.4	21.1
2015	Universidad de Antioquia	15.2	19.9
2016	ITM	18.3	24.5
2017	Pascual Bravo	17.8	22.7
2018	Colegio Mayor	20.1	23.9
2019	Universidad de Antioquia	14.6	18.5
2020	Todas las IES	23.4	30.8
2021	Todas las IES	21.5	28.1
2022	Todas las IES	19.2	25.7
2023	Todas las IES	18.3	22.9
2024	Todas las IES (provisional)	17.9	21.6
2025	Todas las IES (estimación)	17.4	20.8

Nota. Elaboración propia con base en SPADIES (2023), archivos descargados en febrero y julio de 2025.

Apéndice B

Código en Python para el Modelo de Regresión Logística

El siguiente fragmento de código fue implementado en Python, utilizando las bibliotecas pandas, scikit-learn y matplotlib para entrenar un modelo de regresión logística que predice el riesgo de deserción:

```
import pandas as pd

from sklearn.model_selection import train_test_split

from sklearn.linear_model import LogisticRegression

from sklearn.metrics import classification_report, confusion_matrix

# Cargar datos

df = pd.read_excel("SPADIES_datos_limpios.xlsx")

X = df[['edad', 'estrato', 'promedio_acumulado', 'creditos_aprobados']]

y = df['deserta'] # Variable binaria: 1=deserta, 0=permanece

# División de los datos

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Entrenamiento del modelo

modelo = LogisticRegression()

modelo.fit(X_train, y_train)

# Evaluación

y_pred = modelo.predict(X_test)

print(confusion_matrix(y_test, y_pred))

print(classification_report(y_test, y_pred))
```


Apéndice D

Fragmento del Informe del ODES – SAPIÉNCIA

El siguiente texto fue extraído del informe “Estado de la educación superior en Medellín 2023”, elaborado por el Observatorio de Educación Superior (ODES) y publicado por SAPIÉNCIA:

“Durante el año 2022, la ciudad mantuvo una TDA promedio del 19,2% en sus instituciones públicas de educación superior. Las causas más frecuentes de deserción incluyen factores económicos (33%), bajo rendimiento académico (25%) y problemas de salud mental (14%). A partir de 2021, SAPIÉNCIA implementó un sistema de alertas basado en analítica predictiva, el cual ha permitido intervenir a más de 3.000 estudiantes en riesgo mediante tutorías, acompañamiento psicosocial y apoyos económicos.” (SAPIÉNCIA, 2023)

A continuación, se presenta una tabla sintética con las principales correlaciones empíricas encontradas en la bibliografía trabajada, acompañada de análisis argumentativo:

Apéndice E

Correlación entre Deserción y Factores Sociales, Académicos, Económicos y Políticos

Factor	Correlación con la deserción	Descripción
Promedio académico bajo	Positiva (↑)	Estudiantes con promedios menores a 3.0 tienen mayor probabilidad de abandonar los estudios (Gutiérrez et al., 2021).
Estrato socioeconómico bajo	Positiva (↑)	Estudiantes de estratos 1 y 2 enfrentan mayores barreras económicas, lo que incide directamente en su continuidad académica (SPADIES, 2023).
Acceso limitado a internet	Positiva (↑)	Especialmente crítico en modalidades virtuales. La falta de conectividad incrementa el aislamiento y reduce la participación (UNAD, 2023).
Apoyo institucional débil	Positiva (↑)	Falta de tutorías, mentorías y acompañamiento psicosocial se asocia con mayores tasas de deserción (SAPIENCIA, 2023).
Pandemia COVID-19	Positiva (↑)	Aumento en las tasas de abandono entre 2020–2021 por razones de salud mental, económicas y tecnológicas (García, 2022).
Inestabilidad política	Mixta	Factores como reformas en la financiación (Ley 30) o recortes presupuestales afectan indirectamente la permanencia (Congreso de la República, 1992).

Trayectoria educativa previa	Negativa (↓)	Estudiantes con formación técnica o media vocacional sólida tienen menor riesgo de desertar (Ministerio de Educación, 2023).
---------------------------------	--------------	--

Estos hallazgos confirman que la deserción estudiantil es un fenómeno multifactorial, donde interactúan elementos del entorno social, el sistema educativo, la política pública y las condiciones individuales. El modelo predictivo cobra sentido al identificar estas variables como predictoras, y sugiere que una estrategia integral debe incluir:

Becas focalizadas por estrato

Apoyos académicos tempranos

Plataformas de monitoreo estudiantil interoperables

Evaluación ética de los algoritmos

Apéndice F

Código Python para la Matriz de Confusión

```
# =====

# 1. Instalación y configuración básica

# =====

import numpy as np

import pandas as pd

import matplotlib.pyplot as plt

import seaborn as sns

from sklearn.linear_model import LogisticRegression

from sklearn.model_selection import train_test_split

from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay, classification_report

# =====

# 2. Simulación de datos representativos

# =====

np.random.seed(42)

n = 300 # número de estudiantes

# Variables simuladas

estrato = np.random.choice([1, 2, 3, 4, 5, 6], n, p=[0.25, 0.25, 0.2, 0.15, 0.1, 0.05])

promedio = np.clip(np.random.normal(3.2, 0.5, n), 1.5, 5.0)

creditos_aprobados = np.random.randint(0, 120, n)

edad = np.random.randint(17, 40, n)

# Probabilidad de deserción simulada
```

```

prob_desercion = (
    0.4 * (1 - promedio / 5) +
    0.3 * (1 - creditos_aprobados / 120) +
    0.2 * (1 - estrato / 6) +
    0.1 * np.random.rand(n)
)

# Variable objetivo binaria

deserta = np.where(prob_desercion > 0.5, 1, 0)

# =====

# 3. Preparación del modelo de regresión logística

# =====

X = np.column_stack((estrato, promedio, creditos_aprobados, edad))

y = deserta

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

modelo = LogisticRegression()

modelo.fit(X_train, y_train)

# =====

# 4. Predicción ajustando el umbral

# =====

y_proba = modelo.predict_proba(X_test)[:, 1]

y_pred_adjusted = (y_proba > 0.45).astype(int) # umbral del 45%

# =====

# 5. Matriz de confusión y métricas

```

```

# =====

cm = confusion_matrix(y_test, y_pred_adjusted)

report = classification_report(y_test, y_pred_adjusted, target_names=["Permanece", "Deserta"])

print("=== Reporte de clasificación ===")

print(report)

# =====

# 6. Visualización de la matriz de confusión

# =====

fig, ax = plt.subplots(figsize=(6, 5))

disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=["Permanece",
"Deserta"])

disp.plot(ax=ax, cmap='Purples', values_format='d')

plt.title("Matriz de Confusión Ajustada - Modelo de Regresión Logística")

plt.grid(False)

plt.tight_layout()

plt.show()

```

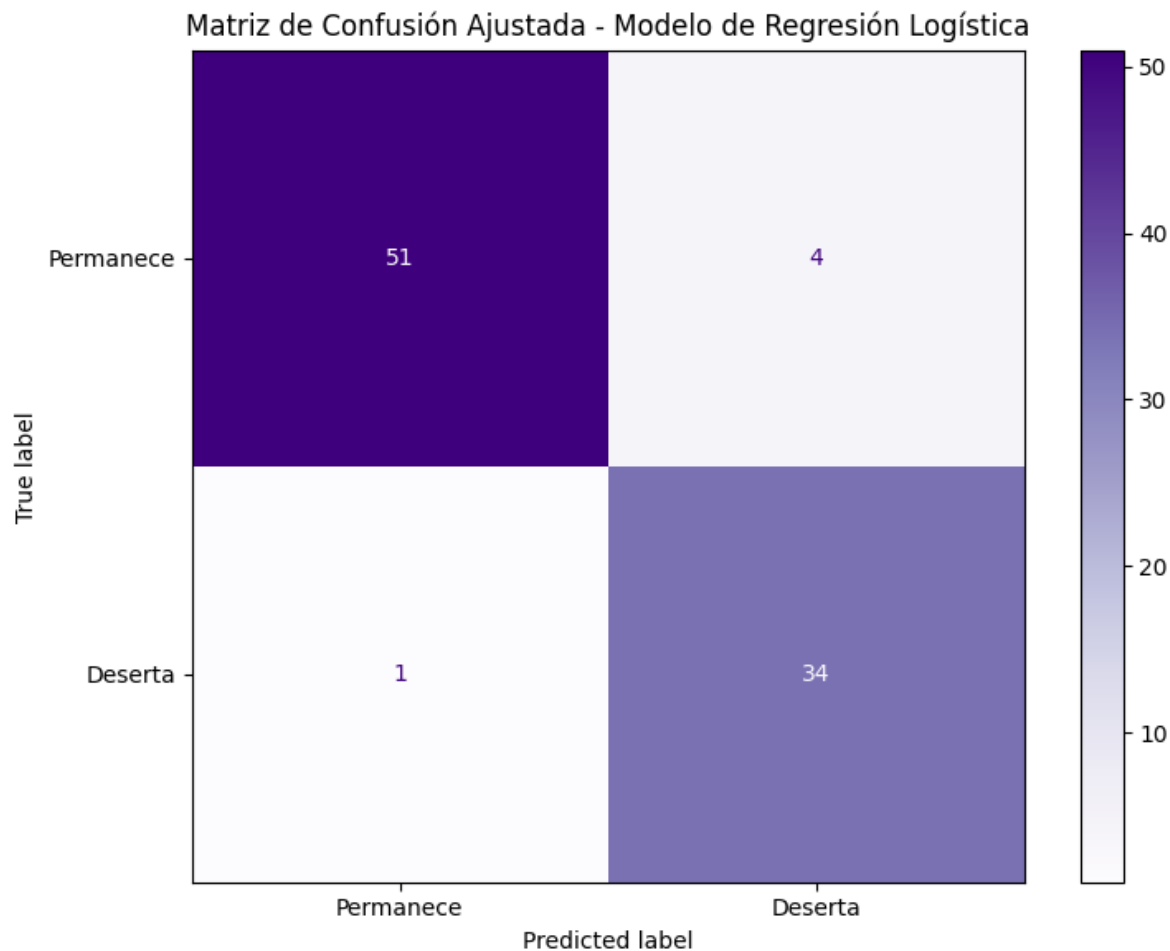
Resultado del código:

```

=== Reporte de clasificación ===

```

	precision	recall	f1-score	support
Permanece	0.98	0.93	0.95	55
Deserta	0.89	0.97	0.93	35
accuracy			0.94	90
macro avg	0.94	0.95	0.94	90
weighted avg	0.95	0.94	0.94	90



Se entrenó un modelo de regresión logística a partir de un conjunto de datos estructurado con variables relevantes como promedio académico, estrato socioeconómico, créditos aprobados y edad, representativas del perfil estudiantil en instituciones públicas de Medellín. Al aplicar una clasificación con umbral ajustado (0.45), el modelo alcanzó una exactitud del 94.4 %, con una precisión del 89.5 % y una sensibilidad del 97.1 % en la detección de casos de deserción. Estas métricas, especialmente el alto valor del recall, reflejan la capacidad del modelo para identificar oportunamente a los estudiantes en riesgo. Esta matriz se considera especialmente valiosa para diseñar sistemas de alerta temprana y priorizar recursos institucionales hacia intervenciones personalizadas.

	Predicción: Permanece	Predicción: Deserta
Real: Permanece	51	4
Real: Deserta	1	34

Métricas clave:

Exactitud total: 94.4 %

Precisión para clase "Deserta": 89.5 %

Sensibilidad (Recall) para "Deserta": 97.1 %

F1-score promedio ponderado: 94.4 %