



# **Actividad 2:**

## **Comparación de métodos**

CE Inteligencia Artificial y Big Data

Sistemas de Big Data

2024/2025

Daniel Marín López

Carlos Chaves Hernández

# Índice

1. Enunciado.....	3
2. Aplicación de métodos.....	3
3. Ventajas y desventajas.....	6

## 1. Enunciado

Una empresa de servicios de salud quiere reducir las ausencias injustificadas de pacientes a sus citas médicas. Sospechan que hay patrones en los historiales de pacientes y el tipo de cita que podrían predecir si un paciente no asistirá. Te piden desarrollar un modelo para anticipar estas ausencias y tomar medidas preventivas (como enviar recordatorios o cambiar el tipo de cita).

## 2. Aplicación de métodos

Para ver cómo se ejecutarán los distintos métodos usaremos una tabla como se hizo anteriormente:

Método	Fase	Actividad
CRISP-DM	Comprensión del negocio	Reducir las ausencias injustificadas a las citas médicas mediante predicción proactiva, lo que tendría un impacto en la reducción de espera para una cita, mejora en la atención al paciente, Incremento de ingresos al reducir huecos no aprovechados.
	Comprensión de los datos	Demográficos: edad, sexo, ubicación, nivel socioeconómico.  Historial: asistencia previa, número de citas canceladas/fallidas.  Tipo de cita: especialidad, urgencia, presencial/virtual.  Información contextual: día y hora, distancia a la clínica, recordatorios enviados, clima.  Variables de comportamiento: llamadas previas, cambios de cita, comentarios administrativos.
	Preparación de los datos	Limpieza: imputar valores nulos, tratar outliers (e.g., edad > 120).  Ingeniería de variables: <ul style="list-style-type: none"> <li>tasa_asistencia = <math>\text{n}^\circ \text{asistencias} / \text{n}^\circ \text{total\_citas}</math></li> <li>día_semana, hora_pico (categorías de horario)</li> <li>recibió_recordatorio (sí/no)</li> <li>distancia_a_clínica (geocodificada si hay dirección)</li> </ul> Codificación de variables categóricas y normalización de variables numéricas.
	Modelado	Clasificación binaria (asiste vs no asiste).  Algoritmos candidatos pueden ser Árboles de decisión, Random Forest, XGBoost, Regresiones logísticas y Redes neuronales simples.

	Evaluación	Si se predicen con una precisión aceptable, por ejemplo $\geq 75\%$ los pacientes que no asistirán, se puede usar para enviar recordatorios u otras medidas.
	Despliegue	<p>Implementación:</p> <ul style="list-style-type: none"> <li>• API que reciba datos del paciente y devuelva probabilidad de ausencia.</li> <li>• Integración en el sistema de gestión de citas.</li> </ul> <p>Usuarios:</p> <ul style="list-style-type: none"> <li>• Personal administrativo para la planificación de llamadas.</li> <li>• Envío automático de recordatorios personalizados.</li> </ul>
TDSP	Comprensión del problema empresarial	<p><b>Objetivo de negocio:</b> Reducir las ausencias injustificadas de pacientes a sus citas médicas.</p> <p><b>Pregunta de ciencia de datos:</b> ¿Podemos predecir con antelación si un paciente no asistirá a su cita?</p> <p><b>Tipo de problema:</b> Clasificación binaria (asiste = 0, no asiste = 1).</p>
	Adquisición y comprensión de los datos	Podemos obtener los datos mediante sistemas de agendamiento, historiales médicos, CRM (para datos del paciente). Estos luego podemos procesarlos con programas como Spark y almacenarlos en HDFS o en Data Lakes. La calidad se puede sustentar mediante validaciones automáticas (formatos, valores nulos, duplicados), trazabilidad (mantener logs del pipeline y versionado de datasets) y verificación semántica (coherencia en fechas).
	Modelado	Se pueden probar varios modelos de clasificación como Regresión Logística, Random Forest o incluso Redes Neuronales. Podríamos revisar métricas como Recall y Precision para ver cómo de bueno es nuestro modelo. Usar validación cruzada para encontrar los mejores hiperparámetros, y que su interpretabilidad sea fácil para los médicos o administradores del hospital.
	Despliegue	Usando librerías como Flask se podría integrar una API Rest. El modelo se empaqueta y expone como endpoint para predicción en tiempo real o batch. Este se integraría con el sistema de agendamiento (EMR o CRM interno). Puede disparar alertas, recordatorios automáticos o reprogramaciones sugeridas.
	Aceptación del cliente	Los resultados pueden estar recogidos en un dashboard con las métricas de nuestro modelo. También se recogerán los casos reales acertados por el modelo y se puede mostrar qué variables son influyentes en el modelo. Se puede probar durante un tiempo corto y ver el feedback que el personal del hospital presenta ante las predicciones para ajustar el modelo.

ASUM-DM	Análisis	<p>Necesidades del usuario:</p> <ul style="list-style-type: none"> <li>● Reducir citas vacías</li> <li>● Mejorar la planificación y comunicación con pacientes</li> </ul> <p>Entregables esperados:</p> <ul style="list-style-type: none"> <li>● Modelo predictivo listo para producción</li> <li>● Documentación y panel de control con alertas</li> <li>● Plan de mejora continua</li> </ul>
	Diseño	<p>Recursos técnicos:</p> <ul style="list-style-type: none"> <li>● Científicos de datos, ingenieros de datos, equipo TI</li> <li>● Infraestructura cloud (Azure, AWS o GCP)</li> </ul> <p>Entorno:</p> <ul style="list-style-type: none"> <li>● Desarrollo y pruebas en entorno cloud (Databricks o JupyterHub)</li> <li>● Producción en contenedor (Docker + Kubernetes si es escalable)</li> </ul>
	Configuración y construcción	<p>Ciclo 1: EDA + primer modelo básico</p> <p>Ciclo 2: optimización del modelo y métricas</p> <p>Ciclo 3: integración en sistemas reales</p>
	Despliegue	<p>Despliegue progresivo por centros.</p> <p>Supervisión paralela con sistema actual para evitar impactos.</p> <p>Posibilidad de revertir si hay fallos.</p>
	Operación y optimización	<p>Monitorización del rendimiento del modelo.</p> <p>Reentrenamiento mensual o si cambia el patrón de comportamiento.</p> <p>Logging de predicciones erróneas para revisión.</p>

	Gestión del proyecto	<p>Roles definidos:</p> <ul style="list-style-type: none"> <li>• Project Manager: coordinación y planificación.</li> <li>• Data Scientist: desarrollo del modelo.</li> <li>• Data Engineer: ETL y arquitectura.</li> <li>• DevOps: despliegue y monitorización.</li> <li>• Stakeholders clínicos: validación y feedback.</li> </ul> <p>Tareas:</p> <ul style="list-style-type: none"> <li>• Planificación de entregas iterativas (Scrum o Kanban).</li> <li>• Reuniones quincenales con stakeholders.</li> <li>• Documentación y checklist de producción.</li> </ul>
--	----------------------	--

### 3. Ventajas y desventajas

Una vez vistas todas las metodologías, podemos ver unas claras ventajas y desventajas de cada una de ellas:

Metodología	Ventajas	Desventajas
<b>CRISP-DM</b>	<ul style="list-style-type: none"> <li>• Flexible y ampliamente adoptado: Ideal para proyectos exploratorios como el tuyo, sin requerimientos técnicos específicos.</li> <li>• Enfocado en el negocio: Comienza por comprender los objetivos del negocio, algo esencial cuando se quiere reducir ausencias médicas.</li> <li>• Iterativo: Permite refinar el modelo de forma continua conforme se aprende más del problema.</li> </ul>	<ul style="list-style-type: none"> <li>• No está orientado a la producción: No cubre detalles específicos sobre cómo desplegar modelos ni automatizarlos.</li> <li>• Poca formalización técnica: Requiere que el equipo complemente con prácticas de ingeniería para implementación moderna (APIs, DevOps, etc.).</li> </ul>
<b>TDSP</b>	<ul style="list-style-type: none"> <li>• Enfocado en implementación: Diseñado para proyectos que requieren desplegar modelos predictivos en producción.</li> <li>• Basado en herramientas modernas: Integra</li> </ul>	<ul style="list-style-type: none"> <li>• Más técnico: Puede ser excesivo si el equipo no tiene experiencia en MLOps o si el proyecto es exploratorio.</li> <li>• Ligero en exploración de negocio: Aunque contempla los objetivos, su foco está en</li> </ul>

	<p>fácilmente entornos colaborativos, control de versiones, pipelines y CI/CD.</p> <ul style="list-style-type: none"> <li>• Ideal para equipos técnicos: Estructura bien definida para equipos multidisciplinarios (científicos de datos, ingenieros, PMs).</li> </ul>	<p>el desarrollo técnico del modelo.</p>
ASUM-DM	<ul style="list-style-type: none"> <li>• Combina negocio y tecnología: Integra bien las fases de entendimiento del negocio, análisis de datos y planificación de despliegue.</li> <li>• Buen equilibrio entre agilidad y estructura: Es más moderno que CRISP-DM, con énfasis en gestión de proyectos y solución empresarial.</li> <li>• Soporta cambios: Diseñado para ambientes cambiantes y evolución de modelos.</li> </ul>	<ul style="list-style-type: none"> <li>• Menos conocido: Requiere más capacitación del equipo si no está familiarizado.</li> <li>• Más pesado que CRISP-DM: Puede implicar burocracia innecesaria si el equipo es pequeño o el caso es de baja complejidad.</li> </ul>

Lo que se recomienda para cada caso es lo siguiente:

- Si es la primera vez que abordan un problema predictivo con fines de mejora operativa, se puede usar CRISP-DM que es un buen punto de partida.
- Si el objetivo es crear un sistema productivo con alarmas, recordatorios automáticos o integración con sistemas clínicos, TDSP será más adecuado.
- Si se desea una visión de negocio bien alineada con procesos analíticos y potencial de escalabilidad futura, usar ASUM-DM puede ser el mejor balance.