

**DECANATURA DE INGENIERÍA INDUSTRIAL  
DECANATURA DE INGENIERÍA DE SISTEMAS  
DECANATURA DE MATEMÁTICAS  
MAESTRÍA EN CIENCIA DE DATOS  
FORMATO DE ENTREGA TRABAJO DE GRADO**

**Fecha de entrega:**

**Estudiante: Yessica Tatiana León Zamora**

**Director: Ivan Alberto Olier Caparroso**

El presente documento avala la entrega del trabajo de grado por parte del director y codirector.

Documentos anexos: copia digital del Trabajo de Grado (1).

22/01/2024

**Firma Director**

**Firma Codirector**

*Tatiana León Zamora*  
**Firma Estudiante**

# **Modelos no supervisados para la visualización de pacientes en cuidados intensivos y con fibrilación auricular**

**Yessica Tatiana León Zamora**

**Escuela Colombiana de Ingeniería Julio Garavito  
Decanatura de Ingeniería Industrial  
Decanatura de Ingeniería de Sistemas  
Decanatura de Matemáticas  
Maestría en Ciencia de Datos  
Bogotá D.C., Colombia  
2023**

# **Modelos no supervisados para la visualización de pacientes en cuidados intensivos y con fibrilación auricular**

**Yessica Tatiana León Zamora**

Trabajo de grado para optar al título de  
Magíster en Ciencia de Datos

Dr. Ivan Alberto Olier Caparroso  
PhD en Inteligencia Artificial

**Escuela Colombiana de Ingeniería Julio Garavito  
Decanatura de Ingeniería Industrial  
Decanatura de Ingeniería de Sistemas  
Decanatura de Matemáticas  
Maestría en Ciencia de Datos  
Bogotá D.C., Colombia  
2023**

© Únicamente se puede usar el contenido de las publicaciones para propósitos de información. No se debe copiar, enviar, recortar, transmitir o redistribuir este material para propósitos comerciales sin la autorización de la Escuela Colombiana de Ingeniería. Cuando se use el material de la Escuela se debe incluir la siguiente nota “Derechos reservados a Escuela Colombiana de Ingeniería” en cualquier copia en un lugar visible. Y el material no se debe notificar sin el permiso de la Escuela.

Publicado en 2023 por la Escuela Colombiana de Ingeniería Julio Garavito. Avenida 13 No 205-59 Bogotá, Colombia  
TEL: +57 – 1 668 36 00

## **Agradecimientos**

Agradezco a mi pareja, Daniel Serrato, por su paciencia, comprensión y motivación durante todo el proceso de investigación y escritura de este trabajo de grado. A mis padres, Nubia Zamora y Emilio León, por su amor incondicional y apoyo constante en cada paso que he dado en la vida. A mi familia, por su cariño y respaldo en todo momento. Y, por último, pero no menos importante, a mis amigos, por estar siempre presentes y ofrecer su ayuda y ánimo en momentos difíciles. Gracias a todos ustedes, este logro no habría sido posible sin su apoyo.

## Resumen

*La fibrilación auricular es una afección cardíaca caracterizada por pulsaciones cardíacas rápidas e irregulares (National Health Services, 2022). Esta es una de las arritmias cardíacas más conocidas en general y también en las unidades de cuidados intensivos. Su prevalencia en UCI es cercana al 10%, y en una UCI cardíaca de 50% (Malik, Candilio, & Hausenloy, 2013)*

*Sin embargo, las características físicas, clínicas o biológicas que puedan desempeñar una relación importante en el desarrollo de fibrilación auricular, no se conocen con certeza. Comprender estas características podría contribuir significativamente a la detección temprana de esta afección médica, especialmente si se pueden interpretar visualmente de manera rápida. Esto, a su vez, facilitaría su inclusión en las rutinas médicas de cuidados intensivos.*

*Con este objetivo en mente, hemos propuesto llevar a cabo un estudio centrado en el uso de modelos de aprendizaje no supervisados interpretables, como Autoencoders, PCA, T-SNE, GTM. Nuestra investigación se concentra en la exploración de modelos interpretables no supervisados que permitan una interpretación visual sencilla de los posibles biomarcadores asociados al desarrollo de fibrilación auricular en pacientes de cuidados intensivos.*

## Abstract

*Atrial Fibrillation is a cardiac condition characterized by rapid and irregular heartbeats (National Health Services, 2022). It is one of the most well-known cardiac arrhythmias both in general and in Intensive Care Units (ICUs). Its prevalence in ICUs is approximately 10%, and in a cardiac ICU, it can be as high as 50% (Malik, Candilio, & Hausenloy, 2013).*

*However, the specific physical, clinical, or biological characteristics that may play a significant role in the development of Atrial Fibrillation are not definitively known. Understanding these characteristics could significantly contribute to the early detection of this medical condition, especially if they can be visually interpreted quickly. This, in turn, would facilitate its integration into ICU medical routines.*

*With this objective in mind, we have proposed to conduct a study focused on the use of interpretable unsupervised learning models, such as Autoencoders, PCA, T-SNE, and GTM. Our research is centered on the exploration of interpretable unsupervised models that allow for a straightforward visual interpretation of biomarkers that may be associated with risk factors for Atrial Fibrillation in ICU patients.*

# Tabla de contenido.

Lista de Figuras

Lista de Tablas

<b>1</b>	<b>INTRODUCCIÓN.....</b>	<b>1</b>
1.1	PROBLEMÁTICA.....	1
1.2	OBJETIVOS Y PREGUNTA DE INVESTIGACIÓN.....	2
1.3	ALCANCE Y LIMITACIONES.....	2
1.4	METODOLOGÍA.....	3
1.4.1	<i>Adquisición de datos.....</i>	<i>3</i>
1.4.2	<i>Análisis exploratorio.....</i>	<i>4</i>
1.4.3	<i>Reducción de dimensionalidad.....</i>	<i>4</i>
1.4.4	<i>Implementación y evaluación de los Modelos de Aprendizaje No Supervisados.....</i>	<i>4</i>
1.4.5	<i>Evaluación de los modelos.....</i>	<i>4</i>
<b>2</b>	<b>MARCO TEÓRICO.....</b>	<b>7</b>
2.1	FIBRILACIÓN AURICULAR.....	7
2.2	CUIDADOS INTENSIVOS.....	7
2.3	BIOMARCADORES COMÚNMENTE ASOCIADOS A FIBRILACIÓN AURICULAR.....	7
2.4	BASE DE DATOS MIMIC – IV.....	8
2.5	IMPUTACIÓN DE DATOS FALTANTES.....	9
2.6	ESTANDARIZACIÓN DE DATOS.....	9
2.7	ANÁLISIS DE VALORES ATÍPICOS.....	9
2.8	REDUCCIÓN DE DIMENSIONALIDAD.....	10
2.8.1	<i>Regresión de Lasso.....</i>	<i>10</i>
2.8.2	<i>Validación cruzada (cross - validation).....</i>	<i>10</i>
2.8.3	<i>Precisión, exhaustividad, Valor F y Valor AUC.....</i>	<i>10</i>
2.9	MODELOS NO SUPERVISADOS PARA VISUALIZACIÓN DE DATOS.....	11
2.9.1	<i>PCA.....</i>	<i>12</i>
2.9.2	<i>t – SNE.....</i>	<i>12</i>
2.9.3	<i>Autoencoders.....</i>	<i>13</i>
2.9.4	<i>Generative Topographic Mapping.....</i>	<i>14</i>
<b>3</b>	<b>RESULTADOS.....</b>	<b>16</b>
3.1	ELECCIÓN Y DEPURACIÓN DE LA BASE DE DATOS.....	16
3.2	INGENIERÍA DE VARIABLES.....	17
3.2.1	<i>Identificación de las variables.....</i>	<i>17</i>
3.2.2	<i>Limpieza de la base de datos.....</i>	<i>18</i>
3.2.3	<i>Adición de nuevas variables.....</i>	<i>19</i>
3.3	ANÁLISIS EXPLORATORIO.....	19
3.4	SELECCIÓN AUTOMÁTICA DE VARIABLES.....	24
3.4.1	<i>Selección de variables.....</i>	<i>26</i>
3.4.2	<i>PCA para la visualización de datos.....</i>	<i>26</i>
3.4.3	<i>Visualización de pacientes utilizando Autoencoder + t-SNE.....</i>	<i>28</i>

3.4.4	<i>Visualización de pacientes con GTM</i> .....	32
4	<b>DISCUSIÓN DE RESULTADOS</b> .....	35
5	<b>CONCLUSIONES Y RECOMENDACIONES</b> .....	37
	<b>BIBLIOGRAFÍA</b> .....	38
	<b>ABREVIACIONES</b> .....	41
	<b>ANEXOS</b> .....	43
5.1	<b>PESOS REGRESIÓN LASSO</b> .....	43

## LISTA DE FIGURAS

Figura 1-1 Metodología.....	5
Figura 2-1 Arquitectura Autoencoder.....	14
Figura 3-1 Boxplot Temperatura .....	22
Figura 3-2 Boxplot Frecuencia cardiaca .....	23
Figura 3-3 Distribución étnica en el conjunto de datos .....	23
Figura 3-4 Distribución casos de fibrilación auricular en UCI.....	24
Figura 3-5 Coeficientes de Lasso en función de Alpha.....	25
Figura 3-6 Curva ROC Regresión Lasso .....	25
Figura 3-7 Presición Exhaustividad Lasso .....	25
Figura 3-8 Gráfico de individuos PCA.....	27
Figura 3-9 Porcentaje de varianza explicada .....	27
Figura 3-10: Pérdida del conjunto de datos de entrenamiento y validación en el Autoencoder.....	28
Figura 3-11 t – SNE con perplejidad 50, distinción aifbFlag = 0 y aifbFlag = 1 .....	29
Figura 3-12 t – SNE con perplejidad 50 destacando aifbFlag = 1 .....	29
Figura 3-13 Concentración de casos con FA en UCI.....	30
Figura 3-14 t-sne enfoque heart_rate_max .....	30
Figura 3-15 t-sne enfoque age.....	30
Figura 3-16 t-sne enfoque diastolic_hb_mean.....	30
Figura 3-17 t-sne heart_rate_amin.....	30
Figura 3-18 t-sne enfoque pt_mean .....	30
Figura 3-19 t-sne enfoque hb_amin .....	30
Figura 3-20 t-sne enfoque ethnicity_black .....	30
Figura 3-21 t-sne enfoque systolic_bp_std .....	30
Figura 3-22 t-sne enfoque hb_amax.....	30
Figura 3-23 GTM Mapa de clases.....	32
Figura 3-24 GTM Concentración de nodos GTM .....	32
Figura 3-25GTM heart_rate_amax.....	33
Figura 3-26 GTM age .....	33
Figura 3-27 GTM diastolic_bp_mean.....	33

Figura 3-28 GTM_heart_rate_amin.....	33
Figura 3-29 GTM pt_mean .....	33
Figura 3-30 GTM hb_amin .....	33
Figura 3-31 GTM ethnicity_BLACK.....	33
Figura 3-32 GTM systolic_bp_std .....	33
Figura 3-33 GTM hb_amax.....	33
Figura 3-34 GTM diastolic_bp_amin .....	33
Figura 3-35 GTM IMC_amin .....	33
Figura 3-36 GTM so2_skew.....	33

**LISTA DE TABLAS**

Tabla 3-1 Variables y su significado .....	16
Tabla 2 Análisis exploratorio.....	20
Tabla 3 Variables con más peso en Regresión Lasso .....	26



# 1 Introducción

## 1.1 Problemática

La fibrilación auricular es una afección cardíaca caracterizada por pulsaciones cardíacas rápidas e irregulares (National Health Services, 2022). Esta es una de las arritmias cardíacas más conocidas en general y también en las unidades de cuidados intensivos (UCI). Su prevalencia en UCI es cercana al 10%, y en una UCI cardíaca de 50% (Malik, Candilio, & Hausenloy, 2013)

Comprender las características físicas, clínicas o biológicas que puedan desempeñar una relación importante en el desarrollo de fibrilación auricular podría contribuir a la detección temprana de esta afección médica, especialmente si se pueden detectar e interpretar visualmente de manera rápida. Esto, a su vez, facilitaría su inclusión en las rutinas médicas de cuidados intensivos.

La falta de interpretabilidad en modelos, especialmente en los complejos no lineales, podría restringir la adopción efectiva de sistemas informáticos en aplicaciones prácticas que confían en el aprendizaje automático e inteligencia computacional para el análisis de datos del área de la medicina. (Vellido, 2020)

Una de las formas en que se puede abordar la interpretabilidad en este contexto es a través de técnicas de visualización de datos y modelos. Al hacerlo, se destaca la importancia de considerar el factor humano al intentar mejorar la interpretabilidad de los modelos en general, y en particular, en el ámbito médico. La colaboración activa entre expertos médicos y profesionales de datos se convierte en un elemento esencial en el proceso de desarrollo de estrategias para garantizar la interpretabilidad de los modelos de datos médicos. Esta integración de conocimientos médicos y técnicos puede contribuir significativamente a una interpretación más precisa y útil de los resultados del análisis de datos en el campo de la salud. (Vellido, 2020).

Es por esto por lo que exploramos la posibilidad de utilizar modelos no supervisados que proporcionan herramientas visuales para lograr una interpretación rápida de las características que puedan estar relacionadas con el desarrollo de la fibrilación auricular en las unidades de cuidados intensivos. Entre los algoritmos que consideramos más adecuados para nuestro estudio se encuentran los Autoencoders, PCA, T-SNE y GTM. Estos algoritmos pueden ofrecer una valiosa perspectiva visual que facilitará la toma de decisiones en el ámbito médico.

En específico, quisimos resolver las siguientes preguntas en el trabajo de grado propuesto: ¿cómo se pueden aplicar modelos interpretables no supervisados para detectar de manera efectiva la presencia de fibrilación auricular en pacientes? Además, ¿cuál es la mejor manera de representar visualmente los resultados de dicha detección con el objetivo de facilitar su comprensión clínica y la toma de decisiones médicas?

## 1.2 Objetivos y Pregunta de Investigación

**Objetivo general:** Implementar un modelo de aprendizaje interpretable no supervisado con el propósito de lograr una interpretación sencilla y efectiva de los biomarcadores recopilados en cuidados intensivos que podrían estar relacionados con el desarrollo de fibrilación auricular. Este enfoque es fundamental para mejorar la comprensión de esta afección médica y facilitar la toma de decisiones clínicas informadas.

**Objetivo específico 1:** Concertar un cohorte con información necesaria y suficiente para su análisis e implementación de los modelos.

**Objetivo específico 2:** Desempeñar un análisis exploratorio conciso del conjunto de datos previamente extraído. Este análisis contendrá estadísticas de cada variable considerada inicialmente como pertinente para el estudio, tendencias, dispersiones y asociaciones identificadas en las variables.

**Objetivo específico 3:** Implementar un modelo de visualización de pacientes con fibrilación auricular en cuidados intensivos utilizando técnicas de aprendizaje no supervisadas. Los modelos seleccionados tendrán características de interpretación visual, con los cuales se permitirá un mejor y más sencilla asociación de los biomarcadores medidos en la unidad de cuidados intensivos con relevancia en la fibrilación auricular en unidades de cuidados intensivos.

**Objetivo específico 4:** Desempeñar una evaluación cualitativa de los modelos implementados para determinar cuáles de estos modelos cumplen con el requisito de proporcionar una representación visual efectiva. Lo anterior permitiría una comprensión más intuitiva y precisa de la relación entre algunas de las mediciones hechas en la unidad de cuidados intensivos y la fibrilación auricular.

## 1.3 Alcance y Limitaciones

La presente investigación se enfocó entonces en el estudio de modelos interpretables no supervisados para la fácil interpretación de los posibles biomarcadores asociados al desarrollo de fibrilación auricular en pacientes de cuidados intensivos. Estos modelos tienen la capacidad de caracterizar los pacientes de cuidados intensivos y sus atributos, en busca de alguna relación con el desarrollo de fibrilación auricular, y así mismo, proveer explicabilidad visual en sus resultados.

En cuanto a las limitaciones del presente trabajo de grado, tenemos como primera instancia la calidad de los datos. Aunque el punto de partida de este trabajo fue una base de datos previamente procesada, allí, se pudo haber encontrado un elevado número de valores nulos o un posible ruido generado por errores en la medición de los biomarcadores, lo cual es esperado

en bases de datos clínicas. Esta particularidad podría introducir sesgos en el conjunto de datos utilizado para el desarrollo de este proyecto.

Los datos empleados en este estudio se originan exclusivamente de un solo centro médico, específicamente el Beth Israel Deaconess Medical Center. Esto, puede introducir cierto sesgo en los modelos y limitar su generalización a otras poblaciones o entornos médicos. Por lo tanto, es necesario considerar esta limitación al interpretar los resultados y al aplicar los modelos desarrollados en contextos clínicos más amplios. Por otro lado, la determinación de si un paciente está experimentando fibrilación auricular o no se fundamenta en la revisión de las notas clínicas, las cuales, están sujetas a posibles errores humanos.

Adicionalmente, en modelos como t – SNE, Autoencoders y PCA, encontramos dificultades de desempeño e interpretabilidad, lo cual puede limitar la aplicación práctica del objetivo.

## 1.4 Metodología

El presente proyecto de grado se enfoca en la aplicación de modelos de aprendizaje no supervisados con capacidad de interpretación visual para agilizar la detección visual de posibles biomarcadores asociados al desarrollo de fibrilación auricular en pacientes de cuidados intensivos. Este enfoque es fundamental para mejorar la comprensión de esta afección médica y facilitar la toma de decisiones clínicas informadas. Para lograr estos objetivos, se siguió una metodología estructurada que se desglosa de la siguiente manera:

### 1.4.1 Adquisición de datos

En esta fase se escogió la base de datos MIMIC IV (Medical Information Mart for Intensive Care) ), una base de datos pública con información recopilada en el Centro Médico Beth Israel Deacones. La base de datos MIMIC-IV es una base de datos de acceso abierto; la versión MIMIC-IV, contiene información completa sobre aproximadamente 250,000 pacientes hospitalizados desde 2008 hasta 2019, y brinda un sólido respaldo de datos para estudios clínicos (Zhang, y otros, 2022), es por esto por lo que, esta elección se basa en las características y relevancia que esta base de datos tiene en el ámbito médico.

Utilizando la MIMIC-IV, se extrajeron características clínicas y de laboratorio que son consideradas comúnmente relevantes en pacientes con fibrilación auricular (Ortega-Martorell et al., 2022), entre las que se incluyen variables como el máximo ritmo cardíaco, el ritmo cardíaco promedio, la desviación de la temperatura, entre otras.

El estudio de los objetivos de este proyecto se inició con un conjunto inicial de 200 variables, las cuales se componían principalmente de mediciones estadísticas extraídas de las de los cambios temporales de las características mencionadas.

## **1.4.2 Análisis exploratorio**

Se realizó un análisis exploratorio en donde el objetivo fue comprender las características iniciales de las variables contenidas en el conjunto de datos.

## **1.4.3 Reducción de dimensionalidad**

El conjunto de datos resultante de la MIMIC IV contiene una gran cantidad de variables, gestionar una cantidad tan significativa de variables para la visualización de datos mediante modelos de aprendizaje no supervisado podría introducir ruido y complejidad innecesarios en nuestro análisis. Para abordar esta cuestión, empleamos el método de Regresión de Lasso, Esta es una regresión cuyo objetivo principal es penalizar la cantidad alta de variables penalizando las características menos importantes al forzar algunos de sus coeficientes a ser cero, al mismo tiempo que reduce la influencia de variables que están fuertemente correlacionadas.

Esta regresión tiene la particularidad de penalizar el uso excesivo de variables, lo que la convierte en una herramienta efectiva para seleccionar las características más relevantes. Este método nos permitió seleccionar de manera sistemática las variables más relevantes a la detección de fibrilación auricular.

## **1.4.4 Implementación y evaluación de los Modelos de Aprendizaje No Supervisados**

Posterior a la reducción de dimensionalidad del conjunto de datos, se avanzó con la implementación de modelos de aprendizaje no supervisados, centrándonos especialmente en la interpretación visual de los resultados obtenidos. En este proceso, se consideraron modelos como Autoencoders, PCA, T-SNE y GTM.

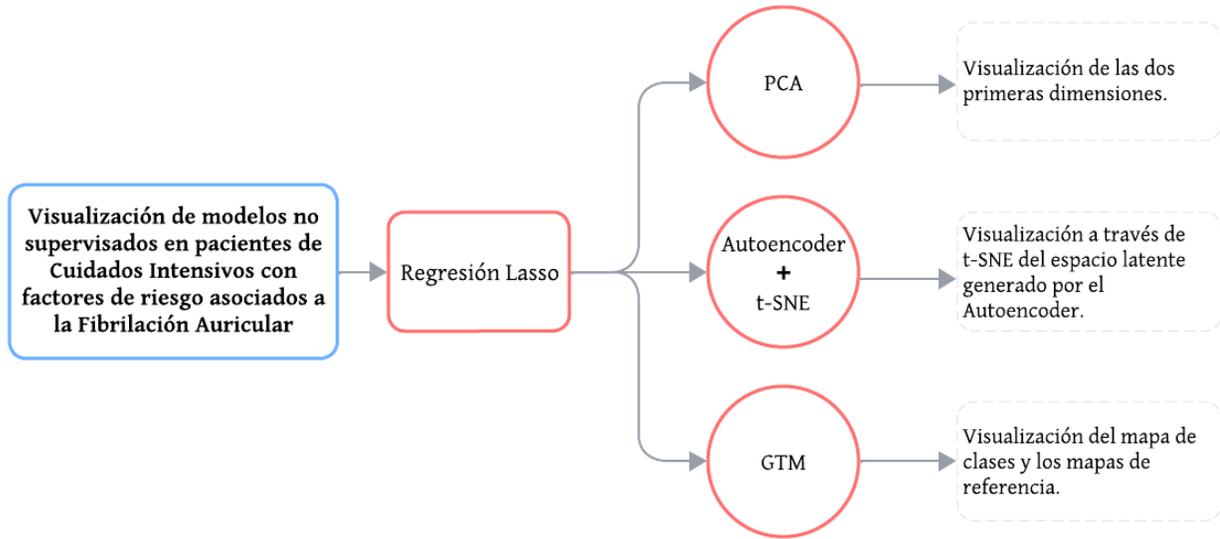
La elección de los anteriores modelos se basó en su capacidad para proporcionar representaciones visuales. Esto se traduce en un recurso valioso que facilitará la comprensión de datos complejos y respaldará la toma de decisiones en el ámbito clínico, contribuyendo así a la detección temprana y al abordaje de la fibrilación auricular en pacientes de cuidados intensivos.

## **1.4.5 Evaluación de los modelos**

Durante esta etapa, se aplicaron criterios para evaluar la eficacia y capacidad de interpretación de los modelos implementados. La prioridad principal fue asegurar que las representaciones visuales generadas fueran útiles para expertos médicos y tomadores de decisiones en el campo médico. Esto se traducirá, en recurso de gran valor que facilitará la comprensión de las relaciones entre algunas de las características recopiladas en las unidades de cuidados intensivos y el desarrollo de fibrilación auricular.

A continuación, presentamos un diagrama en la Figura 1-1, en la que relacionamos los modelos utilizados y mencionados anteriormente.

Figura 1-1 Metodología





## **2 Marco Teórico**

### **2.1 Fibrilación Auricular**

La fibrilación auricular es el tipo más común de arritmia cardíaca tratada; se trata de una condición del corazón que causa pulsaciones cardíacas rápidas e irregulares (National Health Services, 2022).

Se puede señalar que, algunos factores que pueden incidir en la prevalencia de fibrilación auricular están relacionados con el estilo de vida, como hábitos de consumo de cigarrillo y alcohol, mientras que otros son de carácter genético étnico, de género o hereditario (Rizwan, y otros, 2020)

### **2.2 Cuidados Intensivos**

Desde su implementación generalizada hace más de medio siglo, las unidades de cuidados intensivos (UCI) se han convertido en una parte esencial del sistema de atención médica (Marshall, y otros, 2017).

Una UCI representa un sistema altamente estructurado diseñado para brindar atención intensiva y especializada a pacientes en estado crítico. Este sistema proporciona atención médica y de enfermería intensiva, implementa un monitoreo avanzado y ofrece diversas modalidades de soporte fisiológico de órganos para preservar la vida en situaciones en las que existe un riesgo potencialmente mortal de insuficiencia orgánica (Marshall, y otros, 2017).

### **2.3 Biomarcadores comúnmente asociados a fibrilación auricular**

Los biomarcadores sanguíneos han mostrado evidencia asociativa que puede ayudar a refinar la evaluación del riesgo en pacientes con fibrilación auricular (Hijazi, Oldgren, Siegbahn, Granger, & Wallentin, 2013). A continuación, mostramos algunas definiciones de los biomarcadores registrados en la base de datos MIMIC – IV, la cual fue utilizada como objeto de estudio para el presente trabajo.

1. Alanina aminotransferasa: Generalmente conocida como ALT, alanina aminotransferasa es una enzima que se encuentra principalmente en el hígado. Una prueba de ALT mide la cantidad de alanina aminotransferasa en la sangre (National Library of Medicine, 2022).
2. Brecha aniónica: La brecha aniónica mide la diferencia, o la brecha, entre los electrolitos cargados positiva y negativamente en la sangre. Si la brecha aniónica es demasiado alta, significa que la sangre es más ácida de lo normal (National Library of Medicine, 2022).
3. Iones de calcio: El calcio ionizado es calcio en la sangre que no está adherido a las proteínas. Es importante para la función cardíaca en la coagulación de la sangre (National Library of Medicine, 2022).

El calcio juega un papel importante en la fisiopatología del desarrollo de fibrilación auricular, así como en la relación bidireccional entre la fibrilación auricular y la insuficiencia cardiaca (Denham & Pearman).

4. Relleno capilar: La prueba del relleno capilar ungueal es una prueba rápida que se realiza en el lecho ungueal. Se utiliza para monitorear la deshidratación y la cantidad de flujo en la sangre del tejido (National Library of Medicine, 2022).
5. Creatina quinasa: La creatina quinasa (CK) es un tipo de enzima que se encuentra frecuentemente en los músculos esqueléticos (National Library of Medicine, 2022).
6. Proteína c-reactiva: La proteína c-reactiva es producida por el hígado, su nivel se eleva cuando hay inflamación en todo el cuerpo (National Library of Medicine, 2022).
7. Presión arterial diastólica y sistólica: La presión arterial es una medida de la fuerza ejercida contra las paredes de las arterias cuando el corazón bombea sangre al cuerpo (National Library of Medicine, 2022).

La asociación entre el aumento de los niveles de presión arterial y el riesgo de fibrilación auricular es probablemente causal, sin embargo, independientemente de otros factores, el control óptimo de la presión arterial podría representar un objetivo terapéutico para la prevención de la fibrilación auricular (Georgiopoulos & Ntritsos, 2022).

8. Fracción de oxígeno inspirado: La fracción de oxígeno inspirado ( $FiO_2$ ) es la concentración de oxígeno en la mezcla de gases (National Library of Medicine, 2022).
9. GCS: La escala de coma de Glasgow (Glassgow Comma Scale) se utiliza para describir objetivamente el grado de deterioro de la consciencia en todos los tipos de pacientes con traumatismos médicos. La escala evalúa a los pacientes acorde a 3 aspectos de la capacidad de respuesta: respuesta de ojos, motora y verbal. Cada capacidad tiene niveles distintos, siendo 1 el más bajo (Jain S, 2022).
10. Glucosa: La glucosa es el nivel de azúcar en la sangre. La alteración de la glucosa y la insulina puede afectar directamente el miocardio en la aurícula y el ventrículo, lo que conlleva al desarrollo de fibrilación auricular (Sun & Hu, 2009).

## 2.4 Base de datos MIMIC – IV

MIMIC – IV es una base de datos disponible públicamente procedente de la historia clínica del Beth Israel Deconess Medical Center. Esta base de datos abierta tiene el propósito de apoyar distintos estudios de investigación y material educativo (Johnson A. E., y otros, 2023).

La base de datos MIMIC – IV se encuentra disponible en PhysioNet, el cual es un repositorio abierto con datos de investigación médica gestionado por el Laboratorio de Fisiología Computacional del MIT; contiene información acerca de dos sistemas de bases de datos

hospitalarios, de los cuales uno de ellos es un sistema de información clínica específico para la unidad de cuidados intensivos (Johnson A. , y otros, 2020). MIMIC-IV puede ser descargada libremente siempre y cuando el usuario obtenga una certificación para el manejo de datos clínicos expedido por MIT.

## 2.5 Imputación de datos faltantes

En un método de imputación única, se rellenan los datos faltantes de alguna manera única y se utiliza el conjunto de datos resultante para realizar los análisis requeridos. La imputación de la media (MI) es un ejemplo de este enfoque, donde se calcula la media de los valores observados para cada variable y se utilizan para imputar los valores faltantes. Sin embargo, este método puede resultar en estimaciones sesgadas, especialmente cuando hay un gran número de valores faltantes en una variable, ya que la estimación de la varianza puede subestimarse considerablemente (Jamshidian & Mata, 2007).

## 2.6 Estandarización de datos

Una variable estandarizada es una variable que ha sido reescalada para tener una media de cero y una desviación estándar de uno. Las variables se estandarizan por diversas razones, como, por ejemplo, asegurarse de que todas las variables contribuyan de manera equitativa a una escala cuando se suman los elementos o para facilitar la interpretación de los resultados de una regresión u otro análisis (Advanced Research Computing, Statistical Methods and Data Analytics, 2021).

La estandarización de una variable es un procedimiento relativamente sencillo, el cual se puede resumir de la siguiente manera:

$$x^* = \frac{x - \bar{x}}{\sigma}$$

Donde  $x^*$  es la variable generada a partir de la estandarización,  $x$  es la variable en cuestión,  $\bar{x}$  es la media de  $x$  y  $\sigma$  es la desviación estándar de  $x$  (Advanced Research Computing, Statistical Methods and Data Analytics, 2021).

## 2.7 Análisis de valores atípicos

Los valores atípicos, que también se denominan anomalías, en la literatura de minería de datos y estadísticas, son datos inusuales en un conjunto de datos. Estos valores pueden surgir cuando los datos son generados por procesos que se desvían de lo que normalmente se espera. Los valores atípicos a menudo pueden proporcionar información valiosa sobre características inusuales en sistemas o entidades que influyen en la generación de datos.

El reconocimiento de estas características inusuales provee conocimientos útiles en aplicaciones específicas como la medicina. En muchas aplicaciones médicas, los datos se recopilan de una variedad de dispositivos. Patrones inusuales en tales datos suelen reflejar condiciones de enfermedad (Aggarwal, 2017).

## 2.8 Reducción de dimensionalidad

### 2.8.1 Regresión de Lasso

La regresión Lasso es un tipo de regresión penalizadora, en la que su objetivo principal es sancionar el uso de una cantidad alta de variables usando la norma  $\ell_1$ .

Supongamos que para una matriz  $X$  de dimensiones  $n \times p$ , podemos adaptar una restricción tal que la suma del valor absoluto de los  $p$  coeficientes  $\beta$  es menor que una constante  $c$ . La tarea sería obtener coeficientes de la regresión de tal manera que:

$$\hat{\beta} = \min_{\beta} \left[ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \alpha \sum_{j=1}^p |\beta_j| \right] \quad (1)$$

Una vez obtenido  $\hat{\beta}$  minimizando la ecuación (1), los pesos de los valores absolutos de los coeficientes representarían lo que sería la llamada penalidad  $\ell_1$ . Por otro lado,  $\alpha$  determina cuánto peso es dado a la penalidad. (Berk, 2016)

La regresión de Lasso es útil para seleccionar las variables mejor asociadas con la salida a predecir. En nuestro caso, LASSO ayuda a seleccionar las variables que mejor diferencian entre pacientes con fibrilación auricular de los que no.

### 2.8.2 Validación cruzada (cross - validation)

La validación cruzada, o *cross-validation* es una estrategia ampliamente empleada en la selección de modelos basado en su rendimiento de predicción. Su concepto fundamental radica en la subdivisión de los datos, en una o varias ocasiones, con el fin de evaluar el rendimiento del modelo de cada una de las subdivisiones (Arlot, 2010).

Una fracción de los datos, conocida como conjunto de entrenamiento (o *training set*), se emplea para entrenar a cada algoritmo, y por otra parte, la porción restante, denominada conjunto de validación, (o *validation set*) se utiliza para evaluar la capacidad predictiva del modelo, o pérdida. Para mejorar la estimación, el procedimiento de división generalmente se itera al seleccionar sistemáticamente diferentes subconjuntos de datos y resumir el rendimiento predictivo general a lo largo de las iteraciones (Arlot, 2010).

### 2.8.3 Precisión, exhaustividad, Valor F y Valor AUC

1. **Precisión:** La precisión o *precision*, es la proporción de las coincidencias positivas realmente correctas y se expresa mediante la siguiente fórmula:

$$\frac{TP}{TP + FP}$$

En donde *TP* representa las coincidencias clasificadas como verdaderas coincidencias, y *FP* representa las coincidencias clasificadas como falsas coincidencias (Guillet & Hamilton, 2007).

2. **Exhaustividad:** La exhaustividad o *recall*, es la porción de coincidencias reales que han sido clasificadas correctamente y se representa mediante la siguiente fórmula:

$$\frac{TP}{TP + FN}$$

En donde *FN* representa las coincidencias falsamente clasificadas como coincidencias negativas (Guillet & Hamilton, 2007).

3. **Valor F:** El valor F o *f-score* es la media armónica de la precisión y la exhaustividad. Este valor será alto cuando tanto la precisión como la exhaustividad sean altas, lo que lo convierte en una manera de encontrar un equilibrio entre ambas. El valor F se puede calcular mediante la siguiente fórmula (Guillet & Hamilton, 2007):

$$\frac{\textit{precisión} * \textit{exhaustividad}}{\textit{precisión} + \textit{exhaustividad}}$$

4. **Curva ROC:** La curva ROC se traza como la tasa de verdaderos positivos (exhaustividad) en el eje vertical contra la tasa de falsos positivos en el eje horizontal. Esta curva nos ayuda a evaluar la capacidad discriminativa para evaluar la capacidad discriminativa de una prueba diagnóstica dicotómica (Guillet & Hamilton, 2007).
5. **Valor AUC:** El valor AUC (área bajo la curva) es una medida numérica que se encuentra en el rango de 0,5 a 1, donde valores más cercanos a 1 indican un mejor rendimiento del clasificador (Guillet & Hamilton, 2007).

## 2.9 Modelos No Supervisados para visualización de datos

Las técnicas de reducción de dimensionalidad son herramientas fundamentales en la ciencia de datos que permiten transformar datos de alta dimensionalidad en representaciones de menor dimensionalidad, a menudo en 2D o 3D. Estas técnicas se han popularizado en diversos campos, incluyendo la visualización. Su amplia adopción se debe a la prevalencia de conjuntos de datos de alta dimensionalidad, al reducirlos a dimensiones más bajas puede facilitar tareas como la clasificación y la visualización (Rafieian, Hermosilla, & Vázquez, 2023).

Los algoritmos de reducción de dimensionalidad son especialmente útiles para la visualización de datos, ya que ayudan a los usuarios a obtener una idea de la distribución de los datos de alta dimensionalidad, como un gráfico, los vecindarios o su estructura global (Rafieian, Hermosilla, & Vázquez, 2023).

### 2.9.1 PCA

PCA (Análisis de Componentes Principales) *Principal Component Analysis* por sus siglas en inglés es una de las más populares técnicas de reducción de dimensionalidad.

El objetivo principal de este método es identificar el hiperplano que se encuentre más cerca a los datos observados, y luego proyectarlos en él. En este proceso, queremos seleccionar el hiperplano de menor dimensionalidad que preserve la máxima varianza ya que de esta manera, probablemente, tendremos la menor pérdida de información. Es así como, PCA encontrará en el hiperplano tantos ejes como dimensiones en la base de datos haya; el  $i$ -ésimo eje es llamado la  $i$ -ésima Componente Principal de la data. (Géron, 2019)

### 2.9.2 t – SNE

El modelo t-SNE (*t – Stochastic Neighbor Embedding*) es un algoritmo de reducción de dimensionalidad no lineal que se utiliza principalmente para visualizar datos de alta dimensionalidad. que fue desarrollado por van der Maaten y Hinton en 2008 (van der Maaten & Hinton, 2008). Esta técnica es una variante de Stochastic Neighbor Embedding (Hinton & Roweis, 2002), pero se diferencia notablemente por su mayor facilidad de optimización y su capacidad para producir visualizaciones mejoradas al abordar de manera más efectiva el problema de agrupamiento excesivo de puntos en el centro del mapa.

El algoritmo “Stochastic Neighbor Embedding” (SNE) inicia su proceso transformando las distancias euclidianas entre puntos de alta dimensionalidad en probabilidades condicionales que reflejan similitudes. La similitud entre un punto de datos  $x_j$  y otro punto de datos  $x_i$  se representa mediante la probabilidad condicional  $p_{j|i}$ . Esta probabilidad  $p_{j|i}$  indica la probabilidad de que  $x_i$  elija a  $x_j$  como su vecino, asumiendo que los vecinos se eligen en función de su densidad de probabilidad bajo una distribución gaussiana centrada en  $x_i$ . En el caso de puntos de datos cercanos entre sí,  $p_{j|i}$  será relativamente alto, mientras que, para puntos de datos ampliamente separados,  $p_{j|i}$  será casi infinitesimal. Matemáticamente, la probabilidad condicional  $p_{j|i}$  está dada por (van der Maaten & Hinton, 2008):

$$p_{j|i} = \frac{\exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma_i^2}\right)}{\sum_{k \neq i} \exp\left(-\frac{\|x_i - x_k\|^2}{2\sigma_i^2}\right)} \quad (2)$$

Donde  $x_i, x_j \in \{x_1, \dots, x_N\}$ . Hacemos  $p_{i|i} = 0$ .

Por otro lado, para las contrapartes de baja dimensión  $y_i$  y  $y_j$  de los puntos de alta dimensión  $x_i$  y  $x_j$ , es posible calcular una probabilidad condicional similar que se denotará como  $q_{j|i}$ . Al hacer la varianza equivalente a  $\frac{1}{\sqrt{2}}$ . Por lo tanto, se puede modelar la similitud de  $y_i$  y  $y_j$  como:

$$q_{j|i} = \frac{\exp\left(-\|y_i - y_j\|^2\right)}{\sum_{k \neq i} \exp\left(-\|y_i - y_k\|^2\right)} \quad (3)$$

Nuevamente, hacemos  $q_{i|i} = 0$ .

SNE se enfoca en que la diferencia entre  $p_{j|i}$  y  $q_{j|i}$  sea mínima, para así hacer que los puntos  $y_i$  y  $y_j$  modelos correctamente tengan similitud entre  $x_i$  y  $x_j$ . Para ello, se utiliza la divergencia de Kullback-Leibler sobre todos los puntos utilizando un método de gradiente descendente, cuya función de costo  $C$ , está dada por:

$$C = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}} \quad (4)$$

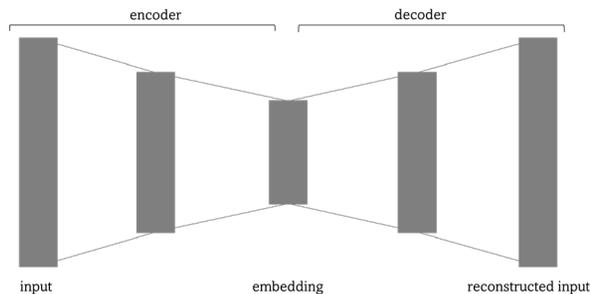
A pesar de que SNE puede generar visualizaciones razonablemente buenas, se enfrenta a desafíos como el "problema de aglomeración", razón por la cual Laurens van der Maaten y Geoffrey Hinton en 2008 introducen la técnica "*t-Distributed Stochastic Neighbor Embedding*" o "t-SNE". Esta técnica se diseñó específicamente para abordar los problemas inherentes a SNE. La principal diferencia radica en que t-SNE utiliza una distribución Student-t en lugar de una distribución Gaussiana (en la ecuación 4) para calcular la similitud entre dos puntos en el espacio de baja dimensión (van der Maaten & Hinton, 2008).

La perplejidad es un parámetro del algoritmo t-SNE, esta se puede entender como una medida que proporciona información acerca del número efectivo de vecinos (van der Maaten & Hinton, 2008).

### 2.9.3 Autoencoders

Un Autoencoder es una arquitectura de red neuronal cuyo objetivo es, en primera instancia, comprimir los datos en un vector de menor dimensión; esta parte de la red es llamada el *encoder*. En segunda instancia, la red se encarga de invertir la transformación de la primera parte y reconstruir entrada original, esta última parte se llama el *decoder*. A continuación, presentamos la arquitectura general de los autoencoders. (Buduma, 2017)

Figura 2-1 Arquitectura Autoencoder



Unas de las principales aplicaciones de los Autoencoders son: la extracción de características y la reducción de características. La principal diferencia entre ellos radica en el uso de parte o la totalidad de las características de entrada. Por ejemplo, considerando un conjunto de datos dado  $X$  con un conjunto de características  $F$ , la selección de características encuentra un subconjunto de características  $D_s$  de todas las características  $F$  ( $D_s \subset F$ ), y el número de características seleccionadas es menor que el original ( $|D_s| < |F|$ ), mientras que la extracción de características genera un nuevo conjunto de características  $D_e$  que son combinaciones de las originales  $F$ . En general, las nuevas características son diferentes de las características originales ( $D_e \neq F$ ), y el número de nuevas características, en la mayoría de los casos, es menor que el de las características originales ( $|D_e| \leq |F|$ ) (Meng, Catchpole, Skillicorn, & Kennedy, 2018)

El principal objetivo de entrenar un modelo Autoencoder es minimizar la diferencia de pérdida entre los datos de entrada y los datos reconstruidos por el Autoencoder. Cuando se logra la convergencia, obtenemos una representación espacial latente en forma de un vector en el espacio creado por el Autoencoder. Podemos visualizar este espacio latente al proyectar estos vectores latentes en dos dimensiones utilizando t-SNE. (Somani, Horsch, & Prasad, 2023).

## 2.9.4 Generative Topographic Mapping

El GTM (*Generative Topographic Mapping*) es un modelo de densidad de probabilidad que describe la distribución de datos en un espacio de varias dimensiones en términos de un número menor de variables latentes. Utiliza una grilla discreta de puntos en el espacio empleando una relación no lineal entre el espacio latente y el espacio de datos originales, manteniendo al mismo tiempo su practicidad. (Bishop, Svensén, & Williams, 1998).

En GTM, los datos en su espacio original se derivan de un espacio latente, típicamente bidimensional ( $L = 2$ ). Este proceso se lleva a cabo mediante un modelo generativo (modelos que crean nuevos datos basados en modelos probabilísticos) que opera en las variables latentes de  $L$  dimensiones. A partir de estas variables, se generan los puntos en el espacio de datos de  $D$  dimensiones ( $D > L$ ) a través de una transformación  $\mathbf{y}$  realizada por una red neuronal de funciones de base radial. (Kireeva, Baskin, Gaspar, Marcou, & Varnek, 2012).

Como se menciona en Bishop, Svensén, & Williams, 1998 el objetivo principal es encontrar una representación para la distribución  $p(\mathbf{t})$  del espacio  $\mathbf{t} = (t_1, \dots, t_D)$  en términos de un número  $L$  de variables latentes  $\mathbf{x} = (x_1, \dots, x_L)$ . Esto es logrado considerando una función  $\mathbf{y}(\mathbf{x}, \mathbf{W})$  con  $\mathbf{y}(\mathbf{x}, \mathbf{W}) : \mathbb{R}^L \rightarrow \mathbb{R}^D$ ,  $\mathbf{y}(\mathbf{x}, \mathbf{W})$  mapeando puntos de  $\mathbf{x}$  a su respectivo  $\mathbf{y}(\mathbf{x}, \mathbf{W})$ . Esto anterior teniendo en cuenta una matriz de pesos y sesgos  $\mathbf{W}$ , proveniente, por ejemplo, de una red neuronal (Bishop, Svensén, & Williams, 1998).

El mapa de clases es una representación visual del espacio latente obtenido en el GTM. En este contexto, frecuentemente se presta especial atención a las áreas donde la densidad de nodos es más elevada, ya que indican una mayor concentración de datos en el espacio de alta dimensionalidad. Al igual que los mapas de clases, los mapas de referencia proporcionan una representación visual de cada una de las variables.

### 3 Resultados

#### 3.1 Elección y depuración de la base de datos

MIMIC – IV es una base de datos disponible públicamente procedente de la historia clínica del *Beth Israel Deconess Medical Center* (Johnson A. E., y otros, 2023). Los datos médicos electrónicos capturados en esta base de datos durante la rutina clínica fueron el insumo principal para el objeto de este proyecto.

Los biomarcadores capturados en la base de datos fueron un factor decisivo para la elección de la MIMIC – IV. Características como la edad, la frecuencia cardiaca, entre otras, han sido determinantes por distintos estudios como precursores del desarrollo de fibrilación auricular en las UCI; lo anterior nos ayudó a consolidar la decisión del uso de esta base de datos.

Previo al inicio del desarrollo de este proyecto, existió un trabajo de transformación y depuración de la base de datos (Olier, Ortega-Martorell, & Lip, Under review). Este esfuerzo consistió en depurar los biomarcadores grabados y calcular medidas de tendencia central y dispersión para cada una de las siguientes características:

*Tabla 3-1 Variables y su significado*

1	<i>alt</i>	Test de alanina	17	<i>magnesium</i>	Magnesio
2	<i>anion_gap</i>	Brecha aniónica	18	<i>mean_bp</i>	Presión arterial media
3	<i>ca_ion</i>	Iones de calcio	19	<i>o2flow</i>	Flujo de oxígeno
4	<i>capref</i>	Relleno capilar	20	<i>peep</i>	Presión positiva al final de la respiración
5	<i>ck</i>	Creatina quinada	21	<i>ph</i>	PH
6	<i>crp</i>	Proteína c – reactiva	22	<i>phosphate</i>	Fosfato
7	<i>diastolic_bp</i>	Presión arterial diastólica	23	<i>platelet</i>	Plaquetas
8	<i>fio2</i>	Fracción de oxígeno inspirado	24	<i>po2</i>	Presión parcial de oxígeno
9	<i>gcs_eye</i>	Escala de Glasgow - Ojos	25	<i>potassium</i>	Potasio
10	<i>gcs_motor</i>	Escala de Glasgow - Motor	26	<i>pt</i>	Tiempo de protrombina
11	<i>gcs_verbal</i>	Escala de Glasgow – Verbal	27	<i>respiratory_rate</i>	Tasa respiratoria
12	<i>glucose</i>	Glucosa	28	<i>scr</i>	Creatinina
13	<i>hb</i>	Hemoglobina	29	<i>so2</i>	Saturación de oxígeno
14	<i>heart_rate</i>	Ritmo cardiaco	30	<i>systolic_bp</i>	Presión arterial sistólica
15	<i>height</i>	Altura	31	<i>temperature</i>	Temperatura
16	<i>lactate</i>	Lactato	32	<i>weight</i>	Peso

Para cada una de las anteriores variables se calculó: máximo, mínimo, curtosis, media, sesgo y desviación estándar, dando como resultado un total de 192 variables.

Adicionalmente, se tuvo en cuenta 7 variables adicionales, esto nos dio un total de 199 para el estudio de nuestro proyecto:

- |  |  |
|--|--|
| 1. <b>Raza:</b> ethnicity                                      | 5. <b>Permanencia en la UCI (horas)</b><br>inicu_los_h                             |
| 2. <b>Edad:</b> age  | 6. <b>Muerte después de la salida de UCI (horas):</b> death_after_icu_h            |
| 3. <b>Muerte en el hospital:</b><br>in_hosp_dead               | 7. <b>Muerte después de la salida del hospital (horas):</b><br>death_after_disch_h |
| 4. <b>Permanencia en el hospital (horas):</b> inhospital_los_h |  |

Es importante indicar que las medidas mencionadas anteriormente fueron calculadas con los datos de una ventana de tiempo de 26 horas desde el inicio de la estancia en la UCI, esto con el fin de evitar la posible fuga de información en la medición de cada uno de los biomarcadores.

Cada paciente fue identificado por medio del indicativo guardado en la variable *hadm\_id*. Por otra parte, también se obtuvo información acerca la ocurrencia de fibrilación auricular dentro de la ventana de tiempo, y también la cantidad de veces que el paciente desarrolló fibrilación auricular en esa misma ventana.

## 3.2 Ingeniería de variables

### 3.2.1 Identificación de las variables

Como insumo se recibieron dos conjuntos de datos, la primera contiene toda la información respecto a características físicas del paciente y sus medidas estadísticas de los biomarcadores. Este conjunto de datos cuenta con:

- ✓ 3 columnas de tipo *Int*
- ✓ 196 columnas de tipo *Float*
- ✓ 1 columna de tipo *Object*

La segunda base de datos nos informa los pacientes que han sufrido fibrilación auricular, en una ventana de tiempo, el tipo de columnas que tenemos en esta segunda base de datos es de:

- ✓ *chart\_time\_step\_h*: Int
- ✓ *hadm\_id*: Int
- ✓ *subject\_id*: Int
- ✓ *afib*: Int

En la columna *chart\_time\_step\_h* se guarda la hora en la que el paciente, identificado por el identificador *hadm\_id*, sufrió de fibrilación auricular desde el tiempo 0, que es en momento en la que ingresa a la unidad de cuidados intensivos. Aquí vamos a encontrar tantos registros por cada paciente como cantidad de veces haya sufrido de fibrilación auricular, el registro de fibrilación auricular se guarda en la variable *afib* como una variable binaria con entradas 1 y 0.

### 3.2.2 Limpieza de la base de datos

Se exploró la información contenida en el conjunto de datos con información estadística y descriptiva del paciente, nos encontramos que hay columnas que tienen valores faltantes, a los cuales les damos el siguiente tratamiento:

1. Variables con información estadística (*amax*, *amin*, *kurtosis*, *mean*, *skew*, *std*):

Los valores faltantes se reemplazan con el promedio de la información disponible de la columna respectiva.

2. *death\_after\_icu\_h*, *death\_after\_disch\_h*, *inhospital los h*, *inicu los h*, *in hosp dead*:

Estas columnas se eliminan, dado que no hacen parte del objetivo del presente estudio.

3. *ethnicity*:

Esta columna se convierte a *dummy* categórica, siendo reemplazada por las columnas *ethnicity\_ASLAN*, *ethnicity\_BLACK*, *ethnicity\_HISPANIC*, *ethnicity\_OTHER*, *ethnicity\_WHITE*.

Respecto al conjunto de datos con la información de la fibrilación auricular de cada paciente, delimitamos la ventana de tiempo desde la hora 0 hasta la hora 26 desde el inicio de la estancia en la UCI, esto con el fin de evitar la posible fuga de información en la medición de cada uno de los biomarcadores.

Procedemos a modificar el conjunto de datos, de tal manera que por paciente conservemos la etiqueta de la existencia de fibrilación auricular en la ventana de tiempo determinada. A partir de eso, encontramos 6,883 registros que sufrieron fibrilación auricular al menos una vez en la ventana de tiempo delimitada.

Al unir las dos bases de datos mencionadas anteriormente, obtenemos un set de datos de 200 columnas y 37,253 filas.

### 3.2.3 Adición de nuevas variables

Al conjunto de datos principal, se le agrega la variable *afibFlag*. Por otro lado, se calculan las siguientes variables, dado que estas, podrían proporcionar información valiosa sobre el estado general de salud de un individuo y su riesgo potencial de desarrollar ciertas condiciones médicas (CDC, 2022):

- ✓  $IMC_{amax} = weight_{amax} / height_{amax}$
- ✓  $IMC_{amin} = weight_{amin} / height_{amin}$
- ✓  $IMC_{mean} = weight_{amean} / height_{amean}$

Al tener la cuenta las variables correspondientes al IMC, prescindimos de las correspondientes al peso y la altura, teniendo en cuenta que la información que necesitamos de estas últimas, están medidas en el IMC. Con esto en mente, procedemos a hacer un análisis exploratorio con un conjunto de datos de 198 columnas y 37,253 filas.

## 3.3 Análisis exploratorio

El conjunto de datos utilizado en la fase inicial del análisis proviene de una base de datos que recopila biomarcadores de pacientes en unidades de cuidados intensivos. Tras un proceso de limpieza e ingeniería de variables, obtuvimos un conjunto de datos que consta de 198 columnas y 37,253 filas. Estas columnas se dividen en tres tipos principales:

- ✓ Unas columnas de tipo *Int* que corresponde a la variable de edad del paciente.
- ✓ 189 columnas de tipo *Float* que representan mediciones estadísticas de cada variable, incluyendo valores máximos, mínimos, curtosis, medias, sesgos y desviaciones estándar.
- ✓ Seis columnas de tipo *Category* que albergan variables categóricas, como la etnia del paciente y un indicador de fibrilación auricular, las variables categóricas como la etnia se transformó a variables tipo *dummy*.

En total, este conjunto de datos incluye 190 variables numéricas.

Dado que estas mediciones se tomaron en pacientes de cuidados intensivos, se presta especial atención a los valores mínimos y máximos, ya que pueden indicar casos extremos relevantes desde una perspectiva clínica. Además, es importante destacar que este conjunto de datos se limita a adultos menores de 89 años debido a medidas de anonimización. La tabla 2 presenta un resumen que incluye los valores promedio, mínimo y máximo de estas variables.

Tabla 2 Análisis descriptivo univariado

<i>Variable</i>	<i>mean (min; max)</i>	<i>Variable</i>	<i>mean (min; max)</i>
<i>alt_amax</i>	183.83 (0; 61,854)	<i>mean_bp_amax</i>	101.70 (56; 330)
<i>alt_amin</i>	127.03 (0; 10,300)	<i>mean_bp_amin</i>	61.88 (14; 133)
<i>alt_kurtosis</i>	-2.03 (- 3; 19)	<i>mean_bp_kurtosis</i>	0.71 (- 3; 19)
<i>alt_mean</i>	151.72 (0; 12,334)	<i>mean_bp_mean</i>	78.52 (42; 133)
<i>alt_skew</i>	-0.04 (- 5; 5)	<i>mean_bp_skew</i>	0.43 (- 4; 5)
<i>alt_std</i>	22.20 (0; 22,624)	<i>mean_bp_std</i>	9.92 (0; 74)
<i>anion_gap_amax</i>	15.14 (2; 89)	<i>o2flow_amax</i>	5.01 (0; 20)
<i>anion_gap_amin</i>	13.00 (0; 49)	<i>o2flow_amin</i>	4.81 (0; 20)
<i>anion_gap_kurtosis</i>	-1.34 (- 3; 19)	<i>o2flow_kurtosis</i>	-2.74 (- 3; 19)
<i>anion_gap_mean</i>	14.08 (2; 52)	<i>o2flow_mean</i>	4.92 (0; 20)
<i>anion_gap_skew</i>	-0.03 (- 5; 5)	<i>o2flow_skew</i>	- (- 5; 5)
<i>anion_gap_std</i>	0.89 (0; 27)	<i>o2flow_std</i>	0.08 (0; 7)
<i>ca_ion_amax</i>	1.19 (1; 4)	<i>peep_amax</i>	6.77 (0; 29)
<i>ca_ion_amin</i>	1.07 (0; 2)	<i>peep_amin</i>	4.85 (0; 22)
<i>ca_ion_kurtosis</i>	-0.40 (- 3; 19)	<i>peep_kurtosis</i>	-0.02 (- 3; 19)
<i>ca_ion_mean</i>	1.13 (1; 3)	<i>peep_mean</i>	5.85 (0; 24)
<i>ca_ion_skew</i>	0.08 (- 5; 5)	<i>peep_skew</i>	0.13 (- 5; 5)
<i>ca_ion_std</i>	0.04 (0; 1)	<i>peep_std</i>	0.67 (0; 8)
<i>capref_amax</i>	0.07 (0; 1)	<i>ph_amax</i>	7.27 (5; 9)
<i>capref_amin</i>	0.02 (0; 1)	<i>ph_amin</i>	6.91 (5; 9)
<i>capref_kurtosis</i>	-2.80 (- 3; 19)	<i>ph_kurtosis</i>	0.04 (- 3; 19)
<i>capref_mean</i>	0.04 (0; 1)	<i>ph_mean</i>	7.14 (5; 9)
<i>capref_skew</i>	0.02 (- 5; 5)	<i>ph_skew</i>	-0.29 (- 5; 5)
<i>capref_std</i>	0.02 (0; 1)	<i>ph_std</i>	0.13 (0; 1)
<i>ck_amax</i>	1,365.00 (6; 98,450)	<i>phosphate_amax</i>	3.88 (0; 19)
<i>ck_amin</i>	916.49 (6; 96,430)	<i>phosphate_amin</i>	3.33 (0; 16)
<i>ck_kurtosis</i>	-1.82 (- 3; 19)	<i>phosphate_kurtosis</i>	-1.33 (- 3; 19)
<i>ck_mean</i>	1,135.00 (6; 97,606)	<i>phosphate_mean</i>	3.61 (0; 16)
<i>ck_skew</i>	-0.03 (- 5; 5)	<i>phosphate_skew</i>	-0.04 (- 5; 5)
<i>ck_std</i>	180.66 (0; 24,972)	<i>phosphate_std</i>	0.22 (0; 7)
<i>crp_amax</i>	71.39 (0; 531)	<i>platelet_amax</i>	214.19 (5; 1,302)
<i>crp_amin</i>	68.97 (0; 531)	<i>platelet_amin</i>	190.52 (5; 1,297)
<i>crp_kurtosis</i>	-2.82 (- 3; 19)	<i>platelet_kurtosis</i>	-1.11 (- 3; 19)
<i>crp_mean</i>	69.77 (0; 531)	<i>platelet_mean</i>	203.24 (5; 1,297)
<i>crp_skew</i>	0.03 (- 2; 5)	<i>platelet_skew</i>	-0.10 (- 5; 5)
<i>crp_std</i>	1.06 (0; 67)	<i>platelet_std</i>	9.51 (0; 288)
<i>diastolic_bp_amax</i>	84.61 (37; 307)	<i>po2_amax</i>	214.79 (14; 734)
<i>diastolic_bp_amin</i>	48.04 (0; 113)	<i>po2_amin</i>	93.09 (3; 587)
<i>diastolic_bp_kurtosis</i>	0.73 (- 3; 19)	<i>po2_kurtosis</i>	-0.07 (- 3; 19)
<i>diastolic_bp_mean</i>	63.05 (14; 118)	<i>po2_mean</i>	130.08 (14; 587)
<i>diastolic_bp_skew</i>	0.50 (- 5; 5)	<i>po2_skew</i>	0.59 (- 5; 5)
<i>diastolic_bp_std</i>	9.12 (0; 51)	<i>po2_std</i>	40.71 (0; 235)

<i>Variable</i>	<i>mean (min; max)</i>	<i>Variable</i>	<i>mean (min; max)</i>
<i>fio2_amax</i>	77.83 (20; 100)	<i>potassium_amax</i>	4.53 (2; 10)
<i>fio2_amin</i>	44.68 (20; 100)	<i>potassium_amin</i>	3.87 (1; 8)
<i>fio2_kurtosis</i>	0.78 (- 3; 19)	<i>potassium_kurtosis</i>	-0.67 (- 3; 19)
<i>fio2_mean</i>	54.73 (20; 100)	<i>potassium_mean</i>	4.17 (2; 8)
<i>fio2_skew</i>	0.87 (- 5; 5)	<i>potassium_skew</i>	0.07 (- 5; 5)
<i>fio2_std</i>	12.02 (0; 37)	<i>potassium_std</i>	0.23 (0; 3)
<i>gcs_eye_amax</i>	3.74 (1; 4)	<i>pt_amax</i>	16.38 (8; 100)
<i>gcs_eye_amin</i>	2.58 (1; 4)	<i>pt_amin</i>	14.82 (8; 100)
<i>gcs_eye_kurtosis</i>	-1.10 (- 3; 19)	<i>pt_kurtosis</i>	-0.97 (- 3; 19)
<i>gcs_eye_mean</i>	3.24 (1; 4)	<i>pt_mean</i>	15.43 (8; 100)
<i>gcs_eye_skew</i>	-0.26 (- 5; 5)	<i>pt_skew</i>	0.30 (- 5; 5)
<i>gcs_eye_std</i>	0.47 (0; 2)	<i>pt_std</i>	0.60 (0; 39)
<i>gcs_motor_amax</i>	5.74 (1; 6)	<i>respiratory_rate_amax</i>	27.15 (9; 224)
<i>gcs_motor_amin</i>	4.27 (1; 6)	<i>respiratory_rate_amin</i>	12.57 (0; 37)
<i>gcs_motor_kurtosis</i>	-1.54 (- 3; 19)	<i>respiratory_rate_kurtosis</i>	0.66 (- 3; 19)
<i>gcs_motor_mean</i>	5.19 (1; 6)	<i>respiratory_rate_mean</i>	19.05 (5; 40)
<i>gcs_motor_skew</i>	-0.29 (- 5; 5)	<i>respiratory_rate_skew</i>	0.31 (- 5; 5)
<i>gcs_motor_std</i>	0.62 (0; 3)	<i>respiratory_rate_std</i>	3.69 (0; 41)
<i>gcs_verbal_amax</i>	4.12 (1; 5)	<i>scr_amax</i>	1.42 (0; 80)
<i>gcs_verbal_amin</i>	2.89 (1; 5)	<i>scr_amin</i>	1.24 (0; 32)
<i>gcs_verbal_kurtosis</i>	-1.63 (- 3; 19)	<i>scr_kurtosis</i>	-1.48 (- 3; 19)
<i>gcs_verbal_mean</i>	3.50 (1; 5)	<i>scr_mean</i>	1.33 (0; 32)
<i>gcs_verbal_skew</i>	-0.02 (- 5; 5)	<i>scr_skew</i>	0.02 (- 5; 5)
<i>gcs_verbal_std</i>	0.54 (0; 2)	<i>scr_std</i>	0.07 (0; 27)
<i>glucose_amax</i>	178.24 (50; 1,764)	<i>so2_amax</i>	99.44 (64; 100)
<i>glucose_amin</i>	108.56 (33; 442)	<i>so2_amin</i>	88.67 (0; 100)
<i>glucose_kurtosis</i>	-0.49 (- 3; 19)	<i>so2_kurtosis</i>	1.83 (- 3; 19)
<i>glucose_mean</i>	139.28 (50; 1,192)	<i>so2_mean</i>	96.46 (20; 100)
<i>glucose_skew</i>	0.18 (- 5; 5)	<i>so2_skew</i>	-0.98 (- 5; 5)
<i>glucose_std</i>	23.08 (0; 587)	<i>so2_std</i>	2.86 (0; 41)
<i>hb_amax</i>	11.21 (4; 20)	<i>systolic_bp_amax</i>	144.47 (70; 311)
<i>hb_amin</i>	10.12 (0; 19)	<i>systolic_bp_amin</i>	94.92 (0; 178)
<i>hb_kurtosis</i>	-0.83 (- 3; 19)	<i>systolic_bp_kurtosis</i>	0.03 (- 3; 19)
<i>hb_mean</i>	10.74 (4; 19)	<i>systolic_bp_mean</i>	118.30 (22; 203)
<i>hb_skew</i>	-0.17 (- 5; 5)	<i>systolic_bp_skew</i>	0.16 (- 5; 4)
<i>hb_std</i>	0.39 (0; 5)	<i>systolic_bp_std</i>	12.87 (0; 78)
<i>heart_rate_amax</i>	102.06 (36; 295)	<i>temperature_amax</i>	37.39 (32; 42)
<i>heart_rate_amin</i>	71.15 (0; 163)	<i>temperature_amin</i>	36.31 (27; 40)
<i>heart_rate_kurtosis</i>	0.08 (- 3; 19)	<i>temperature_kurtosis</i>	-0.49 (- 3; 19)
<i>heart_rate_mean</i>	84.80 (15; 174)	<i>temperature_mean</i>	36.85 (31; 40)
<i>heart_rate_skew</i>	0.30 (- 5; 5)	<i>temperature_skew</i>	0.04 (- 5; 5)
<i>heart_rate_std</i>	8.46 (0; 66)	<i>temperature_std</i>	0.36 (0; 5)
<i>lactate_amax</i>	2.73 (1; 37)	<i>age</i>	65.06 (18; 89)
<i>lactate_amin</i>	1.70 (1; 21)	<i>IMC_amax</i>	29.19 (0; 297)

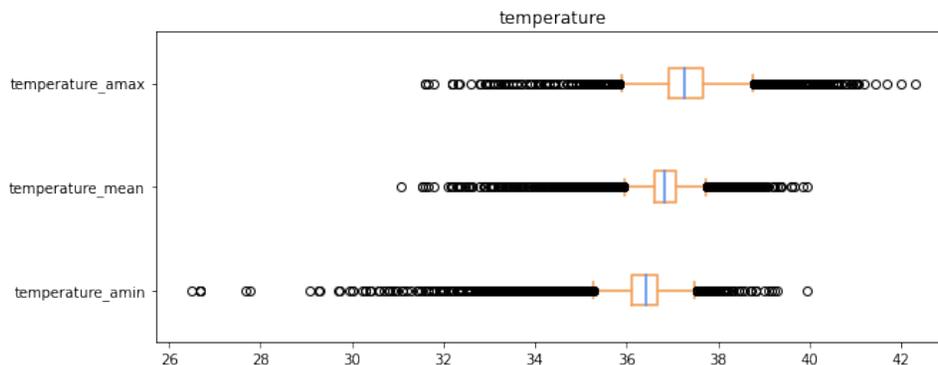
<i>Variable</i>	<i>mean (min; max)</i>	<i>Variable</i>	<i>mean (min; max)</i>
<i>lactate_kurtosis</i>	-0.72 (- 3; 19)	<i>IMC_amin</i>	28.39 (0; 296)
<i>lactate_mean</i>	2.15 (1; 26)	<i>IMC_mean</i>	28.70 (0; 296)
<i>lactate_skew</i>	0.13 (- 5; 5)		
<i>lactate_std</i>	0.37 (0; 15)		
<i>magnesium_amax</i>	2.20 (1; 14)		
<i>magnesium_amin</i>	1.95 (0; 8)		
<i>magnesium_kurtosis</i>	-1.39 (- 3; 19)		
<i>magnesium_mean</i>	2.07 (1; 8)		
<i>magnesium_skew</i>	0.03 (- 5; 5)		
<i>magnesium_std</i>	(0; 5)		

La exploración y el análisis de los valores atípicos desempeñan una importancia crucial en el marco de la atención médica en las unidades de cuidados intensivos, al ser puntos de datos que se desvían significativamente de la norma pueden revelar situaciones médicas excepcionales o críticas que demandan una respuesta inmediata.

Observamos que el valor mínimo de la temperatura mínima se sitúa en 27°C, mientras que la temperatura máxima puede alcanzar valores de hasta 42°C. La temperatura media, que se ubica en 36.85°C, con un rango que varía entre 31°C y 40°C, proporciona información valiosa sobre las complejas situaciones médicas a las que se enfrentan los pacientes en cuidados intensivos y que podrían desencadenar en una fibrilación auricular.

Las temperaturas más elevadas dentro de este rango pueden ser indicativas de infecciones severas y respuestas del cuerpo a condiciones médicas críticas, mientras que las temperaturas más bajas podrían sugerir shock o reacciones a tratamientos específicos (*National Library of Medicine, 2022*).

*Figura 3-1 Boxplot Temperatura*

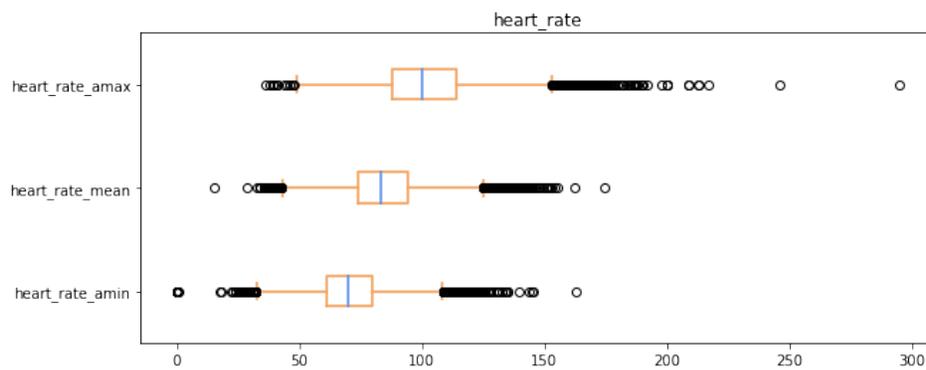


Por otro lado, la frecuencia cardíaca promedio de 84.80 latidos por minuto, que varía desde un mínimo de 15 hasta un máximo de 174 latidos por minuto, podrían indicar respuestas del sistema cardiovascular en situaciones de alta complejidad médica. Es importante destacar que tanto las fluctuaciones elevadas como las reducidas en la frecuencia cardíaca se han correlacionado con

un aumento en la mortalidad en pacientes críticamente enfermos en la unidad de cuidados intensivos (UCI) (Guo, y otros, 2021)

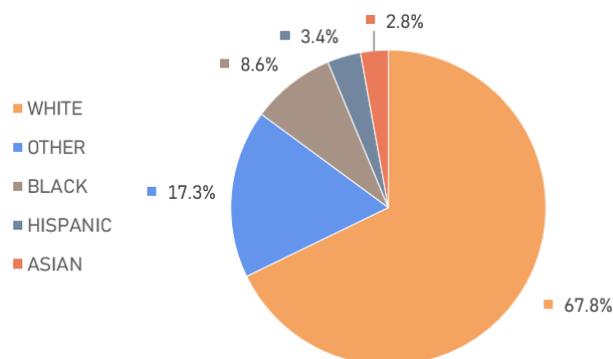
En cuanto al valor mínimo de la frecuencia cardíaca mínima, se sitúa en 0 latidos por minuto, mientras que la frecuencia cardíaca máxima puede alcanzar valores de hasta 300 latidos por minuto. Estos casos atípicos requieren un análisis detenido debido a la complejidad y diversidad de situaciones médicas que pueden surgir en el entorno de cuidados intensivos y que podrían desencadenar en fibrilación auricular.

Figura 3-2 Boxplot Frecuencia cardiaca



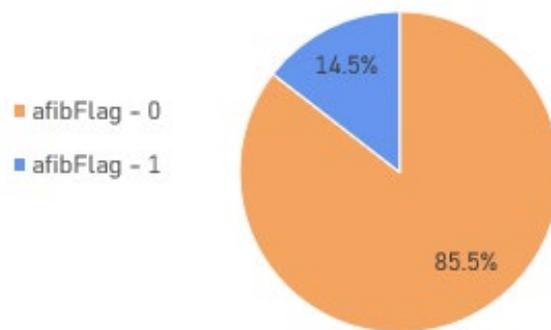
Con relación a las variables categóricas, es esencial examinar la distribución de la etnia de los pacientes en nuestro conjunto de datos. Podemos observar en la Figura 3-3 Distribución étnica en el conjunto de datos que la etnia blanca representa el mayor porcentaje, mientras que la etnia asiática se encuentra en el extremo opuesto con el menor porcentaje. Esta disparidad en las muestras étnicas es un aspecto importante que considerar, dado que tenemos una muestra significativamente mayor de individuos de etnia blanca. Esto puede tener implicaciones en nuestros análisis y resultados, y podría resultar importante en nuestros hallazgos.

Figura 3-3 Distribución étnica en el conjunto de datos



Por último, es relevante destacar que el conjunto de datos presenta un desbalance notable, como se puede evidenciar en la Figura 3-4, se muestra que solo el 14.5% del conjunto de datos total corresponde a pacientes que han experimentado fibrilación auricular en la unidad de cuidados intensivos. Esta desproporción en la distribución de clases es un aspecto importante que considerar en nuestro análisis, ya que puede influir en la capacidad de los modelos de reducción de dimensionalidad para detectar de manera efectiva casos de fibrilación auricular.

*Figura 3-4 Distribución porcentual casos de fibrilación auricular en UCI*

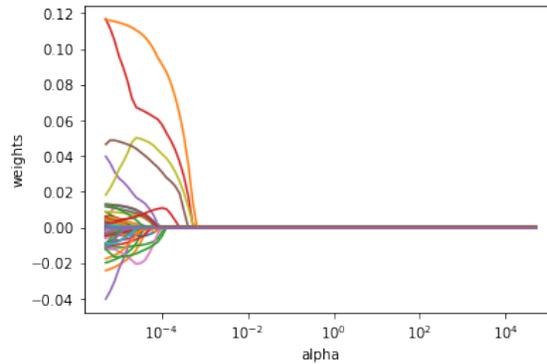


### 3.4 Selección automática de variables

Para asegurar dimensiones uniformes de variables ingresadas al modelo de regularización Lasso, se lleva a cabo la estandarización de las variables numéricas. Posteriormente, se procede a dividir el conjunto de datos, asignando aleatoriamente el 70% al conjunto de entrenamiento y el 30% al conjunto de pruebas. Este enfoque busca optimizar la capacidad del modelo para generalizar patrones identificados durante el entrenamiento a datos no vistos.

En nuestro conjunto de datos de entrenamiento, observamos cómo los pesos de las múltiples variables disminuyeron gradualmente a medida que incrementamos el valor de  $\alpha$ , el factor de penalización. Además, identificamos que el valor de  $\alpha$  óptimo se sitúa en el rango de  $10^{-4}$  y  $10^{-2}$ .

Figura 3-5 Coeficientes de Lasso en función de Alpha



Empleando la técnica de cross-validation, entrenamos nuestro modelo para estimar el valor óptimo de  $\alpha$  con el que pudimos obtener el mejor predictor de error. Una vez completado este proceso, procedimos a analizar las curvas AUC y precision-recall, y, como resultado, calculamos el f1-score.

Figura 3-6 Curva ROC Regresión Lasso

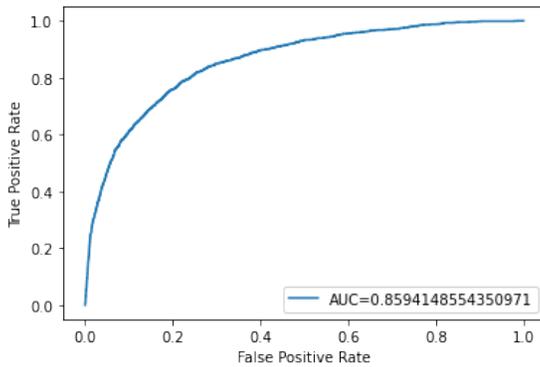
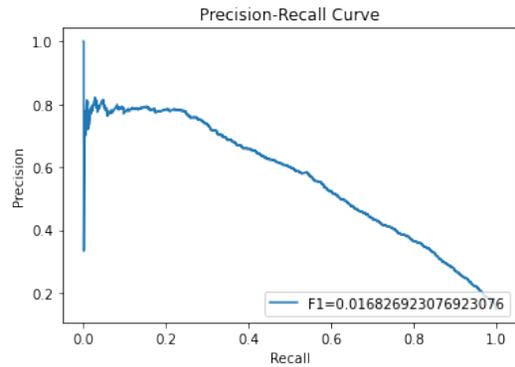


Figura 3-7 Precisión Exhaustividad Lasso



En los gráficos anteriores, podemos observar que con el modelo utilizado logramos obtener un indicador AUC de 0.86. Esto significa que la predicción de registros con fibrilación auricular demostró una notable capacidad de separación en relación con nuestra variable objetivo. Además, al analizar la curva de precisión y exhaustividad, encontramos que podemos alcanzar aproximadamente un 80% de precisión con un 20% de exhaustividad. La combinación de una sólida curva ROC y una alta precisión nos proporciona la confianza necesaria para continuar con las variables seleccionadas por este modelo, ya que indican un buen rendimiento en la detección de casos de fibrilación auricular.

### 3.4.1 Selección de variables

Una de las salidas clave de la Regresión Lasso es la asignación de pesos a cada variable en el modelo entrenado. La Regresión Lasso tiene la capacidad de asignar pesos exactamente cero a las variables que no son relevantes, lo que resulta en una reducción directa de la dimensionalidad de nuestros datos.

Como resultado de aplicar LASSO, obtuvimos un conjunto de 135 variables de las 198 iniciales. A continuación, enumeramos, a manera de ejemplo, las 12 variables con el peso absoluto más alto seleccionadas por la regresión Lasso. No obstante, la lista completa de los pesos de todas las variables se encuentra detallada en el Anexo 5.1 Pesos regresión Lasso.

Tabla 3 Variables con más peso en Regresión Lasso

<i>Variable</i>	<i>Peso</i>
<i>heart_rate_amax</i>	0.1278
<i>age</i>	0.1173
<i>diastolic_bp_mean</i>	0.0540
<i>heart_rate_amin</i>	- 0.0494
<i>pt_mean</i>	0.0403
<i>hb_amin</i>	0.0306
<i>ethnicity_BLACK</i>	- 0.0214
<i>hb_amax</i>	- 0.0199
<i>systolic_bp_std</i>	- 0.0199
<i>diastolic_bp_amin</i>	- 0.0169
<i>IMC_amin</i>	0.0154
<i>so2_skew</i>	- 0.0154

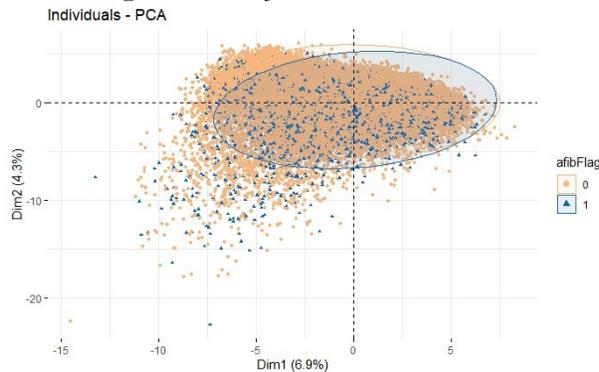
Una vez seleccionadas las 135 variables más relevantes para la predicción de la fibrilación auricular, avanzamos con la evaluación de modelos de aprendizaje no supervisados con el objetivo principal de lograr una interpretación sencilla y efectiva de los resultados. En este proceso de evaluación, hemos considerado una variedad de modelos, incluyendo Autoencoders, PCA (Análisis de Componentes Principales) 3.4.2, T-SNE 3.4.3.1 y GTM (Generative Topographic Mapping) 3.4.4.

### 3.4.2 PCA para la visualización de datos

Llevamos a cabo una visualización a través del Análisis de Componentes Principales (PCA) para explorar la distribución de los individuos en función de las dos primeras dimensiones principales aplicado a las 135 variables elegidas anteriormente. Como se ilustra en la Figura 3-8 Gráfico de individuos PCA, no se aprecia una clara separación entre los dos grupos de individuos: aquellos que experimentaron fibrilación auricular en la unidad de cuidados intensivos y aquellos que no. Esta falta de separación puede atribuirse en parte a la baja cantidad de varianza (Figura 3-9) capturada en el PCA, lo que indica que las dimensiones actuales no representan de manera

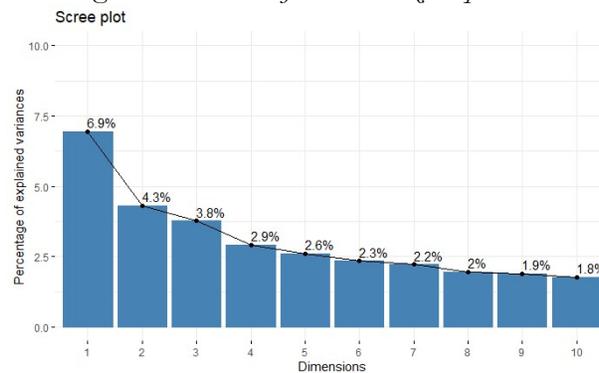
efectiva la información necesaria para distinguir de manera significativa entre los patrones intrínsecos de los datos.

Figura 3-8 Gráfico de individuos PCA



En la Figura 3-8, obtenemos que las dos primeras dimensiones resultantes únicamente explican el 11.2% de la varianza total de los datos. Esto indica que la visualización generada por este modelo carece de información significativa del conjunto de datos original. Como resultado, las relaciones entre los datos originales no se reflejarán adecuadamente en esta visualización.

Figura 3-9 Porcentaje de varianza explicada



Por esta razón, decidimos no llevar a cabo un análisis exhaustivo del gráfico de variables. La baja varianza explicada por las dos primeras dimensiones principales en el PCA no proporciona una base significativa para explorar las variables que podrían influir en los riesgos asociados al desarrollo de fibrilación auricular en la UCI.

Dado que nuestro objetivo principal es interpretar de manera sencilla los posibles biomarcadores recopilados en UCI que estén asociados con el desarrollo de fibrilación auricular, hemos decidido no considerar el PCA como una herramienta útil en este contexto. El motivo de esta elección radica en que el PCA no nos proporcionó una visualización efectiva que facilite una comprensión rápida de esta afección médica. En lugar de ello, optamos por explorar otras técnicas de

reducción de visualización que puedan ofrecer una representación más informativa y accesible de los datos.

### 3.4.3 Visualización de pacientes utilizando Autoencoder + t-SNE

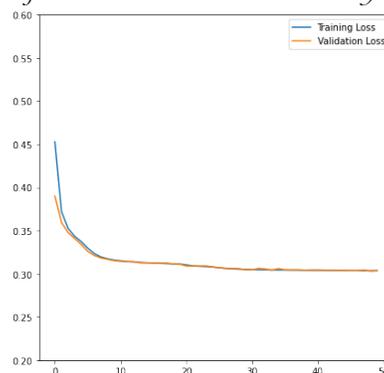
Como parte de nuestro objetivo de implementar un modelo de aprendizaje interpretable no supervisado para lograr una interpretación sencilla de los posibles biomarcadores relacionados con el desarrollo de fibrilación auricular en unidades de cuidados intensivos, optamos por aplicar un Autoencoder al conjunto de datos al que previamente redujimos de dimensionalidad mediante la Regresión Lasso, como se describió en la sección 2.8.1.

La arquitectura del Autoencoder que utilizamos consistió en dos capas intermedias dispuestas de la siguiente manera: 135 – 100 – 50 – 100 – 135. Esta configuración específica fue seleccionada tras realizar diversas pruebas, evidenciando que fue la que arrojó la menor pérdida observada durante el proceso de iterativo. Además, aplicamos la función de activación *ReLU* en las capas del encoder y la función *Lineal* en las capas del decoder. Para evaluar el rendimiento del modelo, utilizamos la métrica de pérdida de Error Absoluto Medio (MAE). Medir la pérdida tanto en el conjunto de validación como en el conjunto de entrenamiento, resultó fundamental para determinar si el modelo estaba experimentando *overfitting*.

Podemos observar, por ejemplo, en la Figura 3-10, que con la configuración de Autoencoder utilizada, se evidencia una disminución en la pérdida tanto en el conjunto de entrenamiento como en el de validación, que va desde aproximadamente 0.45 hasta 0.30. Esto nos lleva a la conclusión de que no estamos experimentando *overfitting*.

El espacio latente obtenido a partir de las 50 dimensiones generadas por el Autoencoder representa las características más relevantes del espacio de las variables originales. Posteriormente, procedimos a visualizar este espacio latente utilizando t-SNE.

Figura 3-10: Pérdida del conjunto de datos de entrenamiento y validación en el Autoencoder



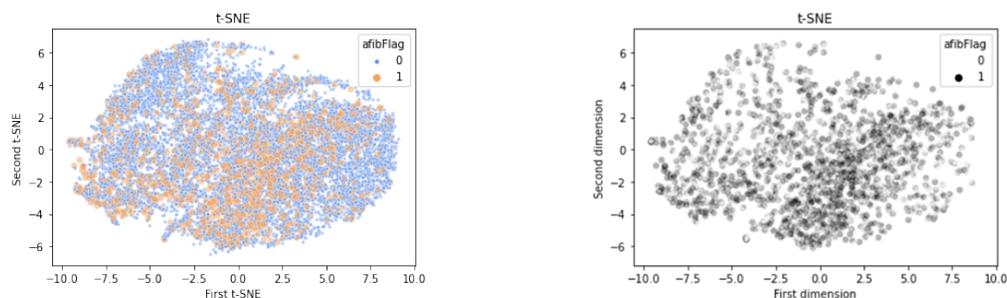
#### 3.4.3.1 Visualización del espacio latente por medio de t-SNE

La perplejidad en los modelos t – SNE es la que nos ayuda a equilibrar los pesos entre los aspectos globales y locales del conjunto de datos entrenada. Dependiendo de la perplejidad aplicada a nuestro conjunto de datos, tenemos cambios en la divergencia KL o divergencia Kullback-Leibler. Entre más cercana a cero sea la divergencia KL, más cercana será nuestra aproximación de la proyección de nuestros datos con los datos originales.

El espacio latente obtenido a partir de las 50 dimensiones generadas por el Autoencoder, como se mencionó previamente, representa las características más relevantes del espacio de las variables originales. Para visualizar este espacio latente y comprender mejor las relaciones entre los datos, aplicamos la técnica de t-SNE con una perplejidad configurada en 50. A través de esta visualización, buscamos identificar patrones o agrupamientos que puedan ayudarnos a entender mejor la estructura de los datos relacionados con la fibrilación auricular en pacientes de cuidados intensivos.

En la Figura 3-11 se presenta el gráfico bidimensional generado por t-SNE. Luego de realizar múltiples iteraciones con diferentes perplejidades, observamos una separación leve entre los pacientes de cuidados intensivos que han experimentado fibrilación auricular y aquellos que no la han tenido, especialmente cuando utilizamos una perplejidad de 50. Posteriormente, en Figura 3-12, destacamos a los pacientes con fibrilación auricular ( $afibFlag = 1$ ) resaltándolos en negro. Esta representación facilita la identificación de las agrupaciones correspondientes a los pacientes con fibrilación auricular en el gráfico de dos dimensiones.

*Figura 3-11 t – SNE con perplejidad 50, distinción  $afibFlag = 0$  y  $afibFlag = 1$       Figura 3-12 t – SNE con perplejidad 50 destacando  $afibFlag = 1$*



En la Figura 3-13, destacada en el recuadro rojo, se observa una ligera concentración de casos en los que los pacientes experimentaron fibrilación auricular al menos una vez durante su estancia en la unidad de cuidados intensivos. Esto se refleja en la agrupación más densa de puntos negros en esa área del gráfico bidimensional generado por t-SNE. En esta región, nos enfocaremos en identificar patrones y relaciones con las variables que tuvieron un mayor peso en la salida de la regresión Lasso (Tabla 3 Variables con más peso en Regresión Lasso).

Figura 3-13 Concentración de casos con FA en UCI

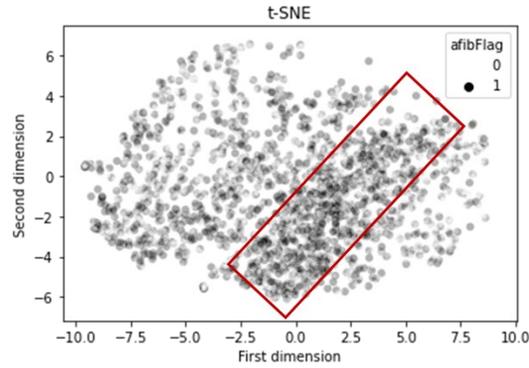


Figura 3-14 t-sne enfoque heart\_rate\_max

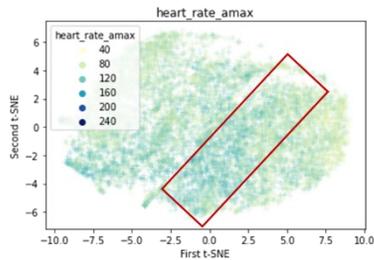


Figura 3-15 t-sne enfoque age

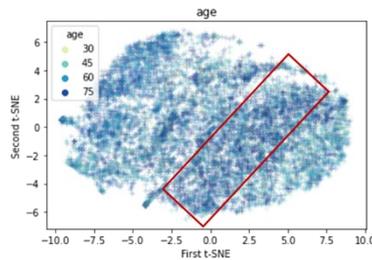


Figura 3-16 t-sne enfoque diastolic\_bp\_mean

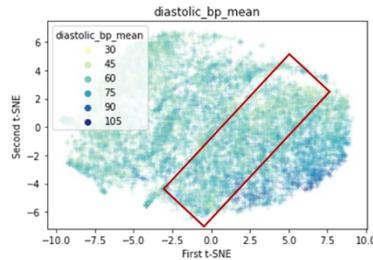


Figura 3-17 t-sne enfoque heart\_rate\_amin

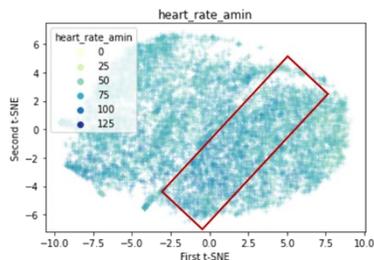


Figura 3-18 t-sne enfoque pt\_mean

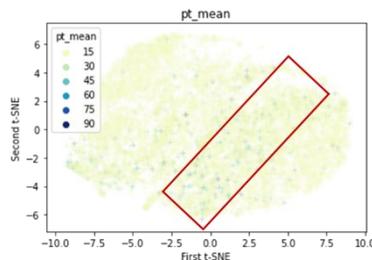


Figura 3-19 t-sne enfoque hb\_amin

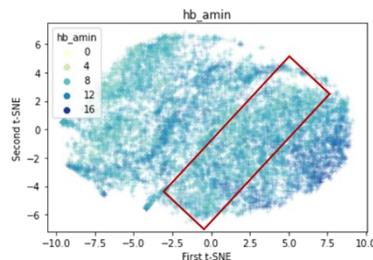


Figura 3-20 t-sne enfoque ethnicity\_black

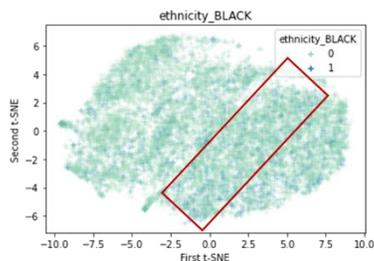


Figura 3-21 t-sne enfoque systolic\_bp\_std

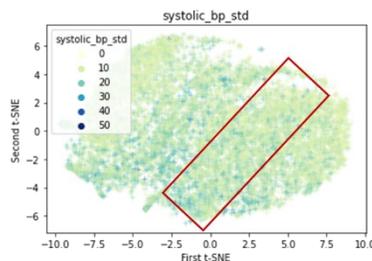
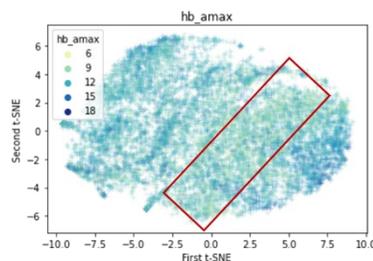


Figura 3-22 t-sne enfoque hb\_amax



En la Figura 3-14, hemos representado el espacio de dos dimensiones generado por t-SNE, coloreando los puntos en función de la frecuencia cardíaca máxima (*heart\_rate\_amax*) medida en la unidad de cuidados intensivos. Como mencionamos anteriormente, nuestro enfoque se centra en el área resaltada en rojo, donde notamos una leve concentración de puntos negros, que representan casos de pacientes que experimentaron fibrilación auricular en la UCI.

Es interesante destacar que, en esta región, podemos observar una ligera concentración de puntos con valores de frecuencia cardíaca máxima más altos (*heart\_rate\_amax* > 160). Esto podría sugerir una posible relevancia de la variable *heart\_rate\_amax* en relación con la fibrilación auricular en el entorno de la unidad de cuidados intensivos.

Por otra parte, en la Figura 3-15 hemos coloreado el espacio generado por t-SNE en función de la edad (*age*). Aquí, notamos una muy ligera concentración de puntos oscuros (*age* > 60) en el área delineada por el recuadro rojo. Este mismo patrón se puede observar en la Figura 3-17 t-sne *heart\_rate\_amin* donde se aplica t-SNE en función de la frecuencia cardíaca mínima (*heart\_rate\_amin*).

En contraste, en la Figura 3-20 t-sne enfoque *ethnicity\_black*, la Figura 3-16 t-sne enfoque *diastolic\_hb\_mean*, la Figura 3-18 t-sne enfoque *pt\_mean* y la Figura 3-21 t-sne enfoque *systolic\_bp\_std*, no se observa una separación clara entre los diferentes contrastes coloreados por cada una de las variables en las que se centró. En estos casos, resulta difícil llegar a conclusiones sobre la relevancia de estas variables en relación con la fibrilación auricular en la unidad de cuidados intensivos mediante el método t-SNE.

Por otro lado, en la Figura 3-19 t-sne enfoque *hb\_amin* y la Figura 3-22 t-sne enfoque *hb\_amax*, se observa una concentración muy leve de puntos claros (indicando valores bajos de *hb\_amin* y *hb\_amax*) en el área encerrada por el recuadro rojo. Este patrón podría sugerir que niveles bajos de hemoglobina podrían tener relevancia en la fibrilación auricular en la unidad de cuidados intensivos.

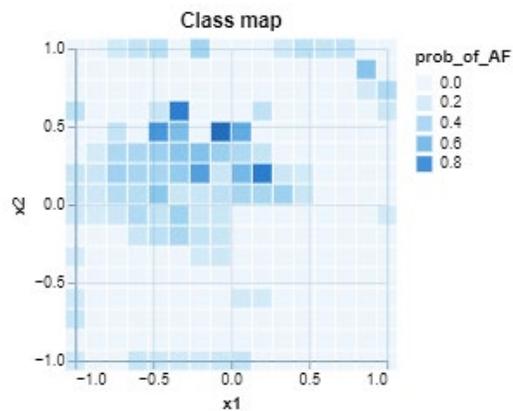
Si bien se ha realizado un análisis con enfoques específicos en varias variables relevantes en nuestro conjunto de datos de pacientes de cuidados intensivos con biomarcadores asociados al desarrollo de fibrilación auricular, es importante destacar que las conclusiones que se pueden extraer de las visualizaciones generadas por t-SNE son limitadas en su certeza.

A pesar de los esfuerzos realizados para identificar patrones y relaciones entre las variables, las visualizaciones no han sido lo suficientemente claras como para respaldar afirmaciones definitivas. Las concentraciones leves de puntos oscuros o claros en áreas específicas del espacio bidimensional podrían sugerir ciertas tendencias, como la posible influencia de la edad, la frecuencia cardíaca máxima o los niveles de hemoglobina en la fibrilación auricular en cuidados intensivos. Sin embargo, por medio de este modelo, estas observaciones son preliminares y requieren un análisis más profundo y un contexto clínico para establecer relaciones de importancia clínica.

### 3.4.4 Visualización de pacientes con GTM

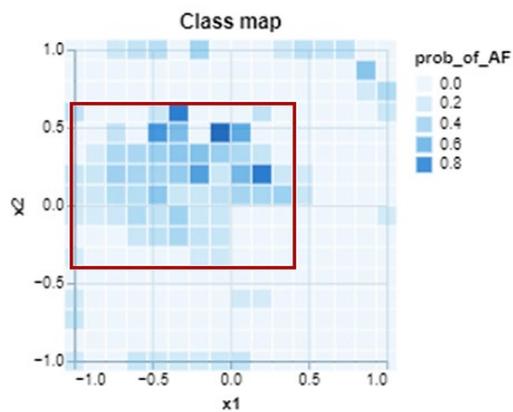
Al conjunto de datos de 135 variables obtenidos posterior a la aplicación de la regresión Lasso, aplicamos un modelo de Generative Topographic Mapping (GTM) a nuestro conjunto de pacientes de cuidados intensivos. La implementación del modelo GTM nos brindó dos visualizaciones clave: el mapa de clases y los mapas de referencia.

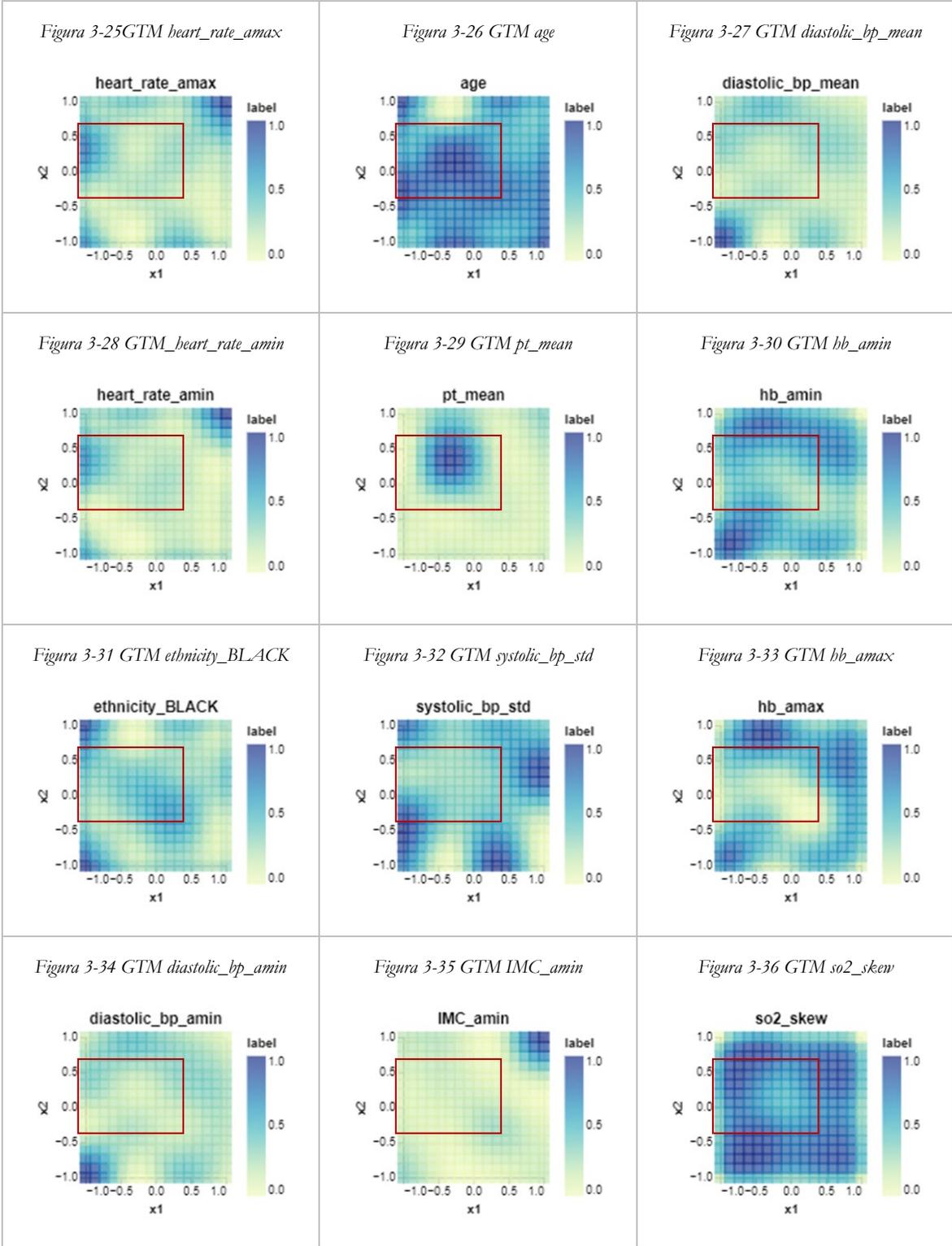
Figura 3-23 GTM Mapa de clases



En la Figura 3-23 GTM Mapa de clases, se puede observar fácilmente que una región de nodos se destaca claramente en la parte superior izquierda. Además, podemos observar que en esa área se concentran nodos con una probabilidad superior al 0.4 de desarrollar fibrilación auricular. Por lo tanto, en la Figura 3-24, hemos resaltado esta región en el espacio latente del modelo GTM. Esto es de importancia, ya que enfocaremos nuestros análisis en los mapas de referencia que se presentan de la Figura 3-25 a la Figura 3-36.

Figura 3-24 GTM Concentración de nodos GTM





De la Figura 3-25 a la Figura 3-36, presentamos el mapa de referencia generado por el GTM de las 12 variables con mayor peso, según lo determinado por la Regresión Lasso (Tabla 3 Variables con más peso en Regresión Lasso). Hemos destacado la misma región que se observó en la Figura 3-24, ya que en cada mapa de referencia nos centraremos en analizar las relaciones en esa área específica.

En las figuras Figura 3-26 y Figura 3-29 tenemos los mapas de referencia generados por el GTM con relación a la edad (*age*) y el tiempo de protombina medio (*pt\_mean*) medidos en la unidad de cuidados intensivos. En estas visualizaciones, se observa una concentración evidente en la región previamente resaltada en la Figura 3-24. Además, en esta área específica, los nodos se caracterizan por tener una alta probabilidad de desarrollar fibrilación auricular.

Este mismo patrón se hace evidente en la Figura 3-36, donde se muestra el mapa de referencia con relación al sesgo de la saturación de oxígeno (*so2\_skew*) medida en la unidad de cuidados intensivos. En esta visualización, nuevamente se observa una concentración de nodos en la región destacada previamente, aunque con una probabilidad alta pero ligeramente menor en comparación con las dos variables anteriores.

Por otro lado, en las figuras Figura 3-25 y Figura 3-28 que representan los mapas de referencia en función del máximo y mínimo de la frecuencia cardíaca medida en la unidad de cuidados intensivos (*heart\_rate\_amax* y *heart\_rate\_amin*), observamos una ligera concentración de nodos hacia la izquierda con una alta probabilidad de desarrollar fibrilación auricular.

En contraste, en las figuras Figura 3-27, Figura 3-30, Figura 3-32, Figura 3-33 y Figura 3-34 que representan los mapas de referencia en función de la media de presión diastólica (*diastolic\_bp\_mean*), el mínimo de la hemoglobina (*hb\_amin*), la desviación de la presión sistólica (*systolic\_bp\_std*), el máximo de la hemoglobina (*hb\_amax*), y mínimo de la presión diastólica (*diastolic\_bp\_amin*), respectivamente, también observamos una concentración en la región destacada en la Figura 3-24. Sin embargo, en esta ocasión, notamos que esta área contiene nodos con una probabilidad baja de desarrollar fibrilación auricular en la unidad de cuidados intensivos.

Sin embargo, en las figuras Figura 3-31 y Figura 3-35, que representan el mapa de referencia en relación con la etnia negra (*ethnicity\_BLACK*) y el IMC mínimo (*IMC\_amin*), respectivamente, no se pueden identificar patrones notables de manera evidente.

Las concentraciones leves de nodos claros u oscuros en áreas específicas del espacio latente resultante del GTM podrían sugerir ciertas tendencias, como la posible influencia de la edad, la frecuencia cardíaca máxima o los niveles de hemoglobina o protrombina en la fibrilación auricular en cuidados intensivos. A diferencia de t-SNE o PCA, estas visualizaciones parecen ofrecer una mayor claridad. Sin embargo, es importante destacar que estas observaciones son preliminares y requieren un análisis más profundo, así como un contexto clínico sólido para establecer relaciones de importancia clínica.

## 4 Discusión de resultados

Los modelos implementados en el presente trabajo de grado para la fácil interpretación de biomarcadores asociados al desarrollo de fibrilación auricular en pacientes de cuidados intensivos son herramientas útiles para visualizar patrones en espacios de menor dimensionalidad a la original del conjunto de datos contemplada. Aunque estos modelos presentan, a grandes rasgos, aspectos técnicos comunes también cuentan con especificaciones distintas, lo que causa enfoques diferentes en cada uno de ellos.

Se pudo observar que el PCA no fue de gran ayuda, ya que no generó visualizaciones claras y sólidas. Esto se debió en gran parte a la baja varianza recogida por sus dos primeras dimensiones, lo que resultó insuficiente para proporcionar una representación significativa de los datos en un espacio de menor dimensionalidad. Además, como mencionamos en la sección 2.9.1, la baja varianza explicada por estas dos dimensiones iniciales impidió que se reflejaran adecuadamente las relaciones y patrones subyacentes en los datos originales. Por lo tanto, resultó evidente que el PCA no fue una herramienta efectiva para nuestra tarea de visualización en relación con la fibrilación auricular en cuidados intensivos.

Por otro lado, aunque en la sección 3.4.3.1 con la técnica de Autoencoder y t-SNE logramos obtener visualizaciones más efectivas que las extraídas mediante el PCA, los esfuerzos realizados para identificar patrones y relaciones entre las variables, las visualizaciones aún no alcanzaron el nivel de claridad necesario para respaldar afirmaciones definitivas. Las leves concentraciones de puntos oscuros o claros en áreas específicas del espacio bidimensional podrían sugerir ciertas tendencias, como la posible influencia de **la edad, la frecuencia cardíaca máxima o los niveles de hemoglobina en la fibrilación auricular en cuidados intensivos**. Sin embargo, es importante destacar que las visualizaciones no proporcionaron la claridad necesaria para extraer conclusiones sólidas a partir de estas observaciones.

Por último, en la sección 3.4.4, el GTM nos proporcionó visualizaciones más efectivas que PCA y t-SNE, con concentraciones leves de nodos claros u oscuros en áreas específicas del espacio latente resultante del GTM, lo que sugirió ciertas tendencias, como **la posible influencia de la edad, la frecuencia cardíaca máxima, los niveles de hemoglobina o protrombina en la fibrilación auricular en cuidados intensivos**.

La edad se destaca como una característica vinculada al riesgo de arritmias, específicamente la fibrilación auricular. Respecto a la frecuencia cardíaca máxima, se sugiere que un incremento en la frecuencia cardíaca podría señalar una mayor exigibilidad del tejido cardíaco, lo cual podría aumentar la probabilidad de desarrollar fibrilación auricular.

En cuanto a la hemoglobina, siendo esta una proteína vital en los glóbulos rojos para el transporte de oxígeno en el cuerpo, irregularidades en sus niveles podrían contribuir al desencadenamiento de fibrilación auricular. De manera similar, alteraciones en la protrombina, esencial para la coagulación sanguínea, también han sido asociadas con la fibrilación auricular.

Estas observaciones destacan la complejidad de los factores que pueden influir en la aparición de la fibrilación auricular, en nuestro contexto específico de cuidados intensivos.

A diferencia de t-SNE o PCA, el GTM ofreció visualizaciones que parecen ofrecer una interpretabilidad más sencilla. Sin embargo, es importante destacar que estas observaciones son preliminares y requieren un análisis más profundo, así como un contexto clínico sólido para establecer relaciones de importancia clínica.

Estas características identificadas pudieron haber sido intuitas o no, por los profesionales de la salud; sin embargo, también pueden ser ratificadas por los resultados obtenidos, lo cual puede mejorar la precisión y eficacia de la atención al paciente.

## 5 Conclusiones y Recomendaciones

1. Al explorar la base de datos MIMIC IV, se pudo acceder a la información para abordar el estudio focalizado en pacientes de cuidados intensivos con observaciones detalladas sobre la aparición de fibrilación auricular durante su tiempo en la UCI. Previo al inicio este proyecto, se llevó a cabo un proceso para tratar la base de datos, dando lugar a la creación de un conjunto de datos con algunas estadísticas de los biomarcadores recolectados. Este conjunto de datos se convirtió en la semilla de nuestra investigación, y tras su tratamiento, que incluyó la eliminación de valores faltantes y ajustes en las variables, logramos obtener los datos necesarios para llevar a cabo el estudio de manera efectiva.
2. A pesar de que el enfoque principal de este trabajo es la generación de visualizaciones con los diferentes modelos no supervisados, el análisis exploratorio realizado en el conjunto de datos proporcionó valiosas señales iniciales sobre el comportamiento de las variables. Este análisis nos permitió obtener una visión general de los rangos de valores con los que estábamos tratando. Por ejemplo, identificamos valores que podrían considerarse atípicos en condiciones normales. También fue evidente que estábamos lidiando con un conjunto de datos desbalanceado, lo que añadió un elemento adicional a tener en cuenta en nuestro análisis.
3. La implementación de un modelo de aprendizaje interpretable no supervisado en este estudio representó una comprensión de la fibrilación auricular en pacientes de cuidados intensivos. A través de la aplicación de técnicas como Autoencoders, PCA, t-SNE y GTM, hemos logrado, en algunos casos, visualizar de manera efectiva las relaciones entre los biomarcadores relacionados con esta afección cardíaca. Aunque las visualizaciones obtenidas han arrojado algunas tendencias preliminares, reconocemos que aún se requiere un análisis más profundo y una validación clínica para confirmar la relevancia de ciertas variables en la predicción y detección de la fibrilación auricular.
4. Este estudio se basó en una sola fuente de datos, la base de datos MIMIC-IV, la cual recopila información de un único centro médico. Para futuras investigaciones sobre la fibrilación auricular en unidades de cuidados intensivos, es recomendable considerar la inclusión de datos de múltiples centros médicos en lugar de depender de un solo centro. Esto enriquecería la investigación y permitiría una comprensión más amplia y aplicable a diversos contextos clínicos.

## Bibliografía

- Advanced Research Computing, Statistical Methods and Data Analytics. (2021). *Advanced Research Computing, Statistical Methods and Data Analytics*. Retrieved from HOW DO I STANDARDIZE VARIABLES IN STATA?: <https://stats.oarc.ucla.edu/stata/faq/how-do-i-standardize-variables-in-stata/>
- Aggarwal, C. C. (2017). An Introduction to Outlier Analysis. In C. C. Aggarwal, *Outlier Analysis* (p. 1). Springer International Publishing.
- Arlot, S. (2010). A survey of cross-validation procedures for model selection. *Statistics Survey*.
- Berk, R. A. (2016). Penalized Smoothing. In *Statistical Learning from a Regression Perspective* (pp. 77-78). Springer.
- Bishop, C., Svensén, M., & Williams, C. (1998). Developments of the Generative Topographic Mapping. *Neurocomputing, Volume 21, Issues 1–3*, 203-224.
- Bosch, N., Cimini, J., & Walkey, A. (2018). Atrial Fibrillation in the ICU. *National Library of Medicine*.
- Buduma, N. (2017). Embedding and Representation Learning. In *Fundamentals of Deep Learning*. O'REILLY.
- CDC. (2022). *Centros para el Control y la Prevención de Enfermedades*. Retrieved from Índice de masa corporal: <https://www.cdc.gov/healthyweight/spanish/assessing/bmi/index.html#:~:text=El%20%C3%ADndice%20de%20masa%20corporal,llevar%20a%20problemas%20de%20salud>.
- Denham, N., & Pearman, C. (n.d.). Calcium in the Pathophysiology of Atrial Fibrillation and Heart Failure.
- Georgiopoulos, G., & Ntritsos, G. (2022). The relationship between blood pressure and risk of atrial fibrillation: a Mendelian randomization study. *European Journal of Preventive Cardiology*.
- Géron, A. (2019). Dimensionality Reduction. In *Hands-on Machine Learning with Scikit-Learn, Keras & TensorFlow* (pp. 220-225).
- Guillet, F., & Hamilton, H. J. (2007). *Quality Measures in Data Mining*. Springer.

- Guo, Q., Xiao, Z., Lin, M., Yuan, G., Qiu, Q., Yang, Y., . . . Wang, J. (2021). Heart rate fluctuation predicts mortality in critically ill patients in the intensive care unit: a retrospective cohort study. *Ann Transl Med.*
- Hijazi, Z., Oldgren, J., Siegbahn, A., Granger, C., & Wallentin, L. (2013). Biomarkers in atrial fibrillation: a clinical review. *European Heart Journal.*
- Hinton, G., & Roweis, S. (2002). *Stochastic Neighbor Embedding*. Toronto: Department of Computer Science, University of Toronto.
- Jain S, I. L. (2022). *Glasgow Coma Scale*. Retrieved from StatPearls [Internet]. Treasure Island (FL): <https://www.ncbi.nlm.nih.gov/books/NBK513298/>
- Jamshidian, M., & Mata, M. (2007). 2 - Advances in Analysis of Mean and Covariance Structure when Data are Incomplete. In M. Jamshidian, & M. Mata, *Handbook of Latent Variable and Related Models* (pp. 21-44). North-Holland.
- Johnson, A. E., Bulgarelli, L., Shen, L., Gayles, A., Shammout, A., Horng, S., . . . Mark, R. G. (2023). MIMIC-IV, a freely accessible electronic health record dataset. [www.nature.com/scientificdata/](http://www.nature.com/scientificdata/).
- Johnson, A., Bulgarelli, L., Pollard, T., Horng, S., Celi, L. A., & Mark, R. (2020). *PhysioNet*. Retrieved from MIMIC-IV: <https://physionet.org/content/mimiciv/0.4/>
- Kireeva, N., Baskin, I., Gaspar, H., Marcou, G., & Varnek, A. (2012). Generative Topographic Mapping (GTM): Universal Tool for Data Visualization, Structure-Activity Modeling and Dataset Comparison. *Molecular informatics.*
- Malik, A., Candilio, L., & Hausenloy, D. (2013). Atrial fibrillation in the intensive care setting. *The Intensive Care Society.*
- Marshall, J., Bosco, L., Adhikari, N., Connolly, B., Diaz, J., Dorman, T., . . . Zimmerman, J. (2017). What is an intensive care unit? A report of the task force of the World Federation of Societies of Intensive and Critical Care Medicine. *Journal of Critical Care*, 270-276.
- Meng, Q., Catchpoole, D., Skillicorn, D., & Kennedy, P. (2018). *Relational Autoencoder for Feature Extraction*. Sydney, Australia: Centre for Artificial Intelligence, Faculty of Engineering and Information Technology, University of Technology Sydney.
- Molnar, C. (2023). *Interpretable Machine Learning*.
- National Health Services. (2022). *NHS*. Retrieved from Conditions - Atrial Fibrillation: <https://www.nhs.uk/>

- National Library of Medicine. (2022). *MedlinePlus*. Retrieved from <https://medlineplus.gov>
- Olier, I., Ortega-Martorell, S., & Lip, G. (Under review). Multi-outcome model visualisation using Generative Topographic Mappings.
- Olier, I., Vellido, A., & Giraldo, J. (2010). *Kernel Generative Topographic Mapping*.
- Ortega-Martorell, S., Pieroni, M., Johnston, M., Olier, I., & Welters, I. (2022). Development of a Risk Prediction Model for New Episodes of Atrial Fibrillation in Medical-Surgical Critically Ill Patients Using the AmsterdamUMCdb. *Front. Cardiovasc. Med.*
- Rafieian, B., Hermosilla, P., & Vázquez, P.-P. (2023). Improving Dimensionality Reduction Projections for Data Visualization. *Applied Sciences*.
- Rizwan, A., Zoha, A., Mabrouk, I., Sabbour, H., Al-Sumaiti, A. S., Alomainy, A., . . . Abbasi, Q. H. (2020). Review on the State of the Art in Atrial Fibrillation Detection Enabled by Machine Learning. *IEEE reviews in biomedical engineering*, 219-239.
- Saito, O., & Nagashima, K. (2022). *Low alanine aminotransferase levels are independently associated with mortality risk in patients with atrial fibrillation*. [www.nature.com/scientificreports](http://www.nature.com/scientificreports).
- Somani, A., Horsch, A., & Prasad, D. (2023). *Interpretability in Deep Learning*. Springer.
- Sun, Y., & Hu, D. (2009). The link between diabetes and atrial fibrillation: cause or correlation? *Journal of Cardiovascular Disease Research*.
- van der Maaten, L., & Hinton, G. (2008). Visualizing Data using t-SNE. *Journal of Machine Learning Research* 9.
- Vellido, A. (2020). The importance of interpretability and visualization in machine. *Neural Computing and Applications*, 18069–18083.
- Zhang, L., Li, S., Yuan, S., Lu, X., Li, J., Liu, Y., . . . Yin, H. (2022). The Association Between Bronchoscopy and the Prognoses of Patients With Ventilator-Associated Pneumonia in Intensive Care Units: A Retrospective Study Based on the MIMIC-IV Database. *Front. Pharmacol.*

## Abreviaciones

FA	fibrilación auricular
UCI	unidad de cuidados intensivos
t – SNE	t-distributed stochastic neighbor embedding
PCA	Principal Component Analysis
GTM	Generative Topographic Mapping



## Anexos

### 5.1 Pesos regresión Lasso

Variable	Peso	Variable	Peso
heart_rate_amax	0.1278	gcs_verbal_kurtosis	- 0.0004
age	0.1173	crp_skew	- 0.0004
diastolic_bp_mean	0.0540	glucose_amax	- 0.0004
pt_mean	0.0403	capref_skew	- 0.0005
hb_amin	0.0306	platelet_std	- 0.0006
IMC_amin	0.0154	so2_amin	- 0.0007
magnesium_amin	0.0145	hb_std	- 0.0007
so2_std	0.0132	crp_std	- 0.0009
scr_amin	0.0118	alt_mean	- 0.0009
pt_amin	0.0116	po2_kurtosis	- 0.0009
gcs_motor_mean	0.0111	weight_kurtosis	- 0.0009
heart_rate_std	0.0105	fiO2_std	- 0.0012
respiratory_rate_mean	0.0101	ck_amin	- 0.0012
mean_bp_skew	0.0100	lactate_std	- 0.0013
ph_mean	0.0073	magnesium_amax	- 0.0015
potassium_kurtosis	0.0071	peep_amax	- 0.0016
diastolic_bp_amax	0.0068	anion_gap_skew	- 0.0016
po2_mean	0.0065	lactate_kurtosis	- 0.0017
ethnicity_WHITE	0.0054	phosphate_skew	- 0.0017
glucose_skew	0.0052	respiratory_rate_skew	- 0.0018
capref_amax	0.0051	scr_std	- 0.0019
fiO2_amin	0.0050	temperature_kurtosis	- 0.0021
glucose_amin	0.0042	systolic_bp_skew	- 0.0021
fiO2_kurtosis	0.0038	gcs_eye_amax	- 0.0024
phosphate_amin	0.0037	potassium_amin	- 0.0025
weight_std	0.0035	o2flow_skew	- 0.0030
ca_ion_amin	0.0035	ck_kurtosis	- 0.0034
anion_gap_kurtosis	0.0033	potassium_amax	- 0.0035
pt_std	0.0033	temperature_skew	- 0.0036
scr_skew	0.0033	gcs_verbal_amax	- 0.0037

<b>Variable</b>	<b>Peso</b>	<b>Variable</b>	<b>Peso</b>
po2_amin	0.0032	alt_amin	- 0.0041
po2_std	0.0031	magnesium_kurtosis	- 0.0043
o2flow_std	0.0031	gcs_eye_std	- 0.0046
respiratory_rate_kurtosis	0.0031	platelet_amax	- 0.0046
potassium_std	0.0029	gcs_eye_amin	- 0.0046
capref_amin	0.0027	gcs_verbal_std	- 0.0047
gcs_eye_skew	0.0026	gcs_motor_std	- 0.0048
weight_skew	0.0026	glucose_kurtosis	- 0.0051
ck_skew	0.0026	fiO2_skew	- 0.0053
scr_kurtosis	0.0021	fiO2_mean	- 0.0053
peep_kurtosis	0.0020	ca_ion_amax	- 0.0054
hb_kurtosis	0.0020	respiratory_rate_amin	- 0.0060
ck_std	0.0019	ph_amax	- 0.0062
ph_skew	0.0019	temperature_std	- 0.0069
platelet_skew	0.0017	heart_rate_kurtosis	- 0.0070
po2_skew	0.0015	glucose_mean	- 0.0074
crp_kurtosis	0.0012	temperature_mean	- 0.0074
so2_amax	0.0011	temperature_amax	- 0.0075
phosphate_mean	0.0009	respiratory_rate_std	- 0.0085
anion_gap_amin	0.0008	systolic_bp_amax	- 0.0086
gcs_motor_kurtosis	0.0007	mean_bp_kurtosis	- 0.0086
capref_kurtosis	0.0006	systolic_bp_mean	- 0.0092
gcs_motor_amax	0.0005	mean_bp_mean	- 0.0093
ph_kurtosis	0.0004	gcs_eye_mean	- 0.0094
alt_kurtosis	0.0004	diastolic_bp_kurtosis	- 0.0100
crp_mean	0.0004	ethnicity_HISPANIC	- 0.0117
platelet_kurtosis	0.0003	mean_bp_amin	- 0.0120
anion_gap_std	0.0002	lactate_amin	- 0.0126
lactate_skew	0.0002	so2_kurtosis	- 0.0142
o2flow_amax	0.0002	so2_skew	- 0.0154
ca_ion_skew	0.0002	diastolic_bp_amin	- 0.0169
peep_mean	0.0002	hb_amax	- 0.0199
respiratory_rate_amax	0.0001	systolic_bp_std	- 0.0199
mean_bp_amax	0.0001	ethnicity_BLACK	- 0.0214
alt_std	0.0000	heart_rate_amin	- 0.0494
potassium_skew	- 0.0001		
pt_kurtosis	- 0.0001		

<b>Variable</b>	<b>Peso</b>	<b>Variable</b>	<b>Peso</b>
phosphate_std	- 0.0001		
pt_skew	- 0.0004		
temperature_amin	- 0.0004		